

NPC Chatbot for Skyrim using Bidirectional LSTMs

Mohan Chandra S S, Karthik S M, Mohith N Reddy

School of Computer Science and Engineering, RV University, Bangalore, India
USN: 1RVU23CSE282, 1RVU22CSE215, 1RVU22CSE283

Abstract—Natural language interaction in role-playing games like Skyrim significantly enhances immersion and user experience. This project presents an intelligent NPC chatbot trained on immersive dialogue data, aiming to replicate the in-game conversational style of Skyrim. Using a deep learning approach based on Bidirectional LSTM networks, the chatbot is trained to generate context-aware responses. The system utilizes a custom dialogue dataset and undergoes preprocessing, tokenization, and training using TensorFlow/Keras frameworks. Experimental results highlight the chatbot's ability to generate thematically consistent dialogue, laying the groundwork for future integration into gaming engines and extended conversational agents.

Index Terms—NPC Chatbot, Skyrim, Natural Language Processing, Deep Learning, Dialogue Systems, Interactive Storytelling

I. INTRODUCTION

The ability of non-player characters (NPCs) to interact using natural language significantly enhances the realism and immersion in modern video games. Skyrim, an open-world role-playing game, has inspired many AI projects due to its deep lore and expansive dialogue trees. Traditional methods for chatbot implementation are often rigid and lack contextual adaptability. With the advent of deep learning, especially Recurrent Neural Networks (RNNs) and their variants like LSTM (Long Short-Term Memory), there is now an opportunity to create dynamic, learning-based dialogue systems.

This project aims to develop an LSTM-based NPC chatbot that mimics the conversational style of Skyrim NPCs. Leveraging dialogue data and deep learning architectures, the chatbot is trained to respond in a contextually appropriate and lore-consistent manner. This work explores the design, training, and evaluation of such a model, offering insights into its feasibility for in-game deployment and potential extensions for multi-turn dialogue capabilities.

II. RELATED WORK

Dialogue systems and chatbots have been central to the field of Artificial Intelligence (AI), with advancements in machine learning significantly enhancing their effectiveness. Early approaches primarily relied on rule-based systems and hand-crafted scripts, which limited the flexibility and naturalness of the responses. Recent progress has shifted toward neural-based architectures, which can learn from data and produce more natural and dynamic conversations.

In the context of NPC chatbots, there has been considerable interest in games like Skyrim, which feature expansive dialogue trees and allow for a rich, open-world experience. Several projects have utilized deep learning techniques to model NPC behavior, using datasets containing dialogues from games. Studies have shown promising results in generating coherent, contextually relevant responses for in-game NPCs.

However, many existing implementations still face challenges in maintaining long-term coherence in multi-turn dialogues and ensuring that responses align with the game's lore. This work aims to address these issues by focusing on an LSTM-based architecture, which is well-suited to sequential data and capable of capturing contextual information over long dialogues.

III. METHODOLOGY

A. Dataset and preprocessing

The dataset used for training the NPC chatbot was generated using an AI model designed to replicate the conversational style and narrative tone of Skyrim NPCs. It originally consisted of 7,565 unique dialogue interactions, which were expanded to over 50,000 through data augmentation techniques that introduced conversational variations. This ensured broader coverage of topics such as quests, lore, and world-building. The dataset was guided by content from the imperial.json file, which contains Skyrim's core lore, ensuring that the chatbot's responses remain contextually and thematically accurate.

Prior to training, the data undergoes several preprocessing steps:

- **Tokenization:** The raw text data is tokenized using Keras' Tokenizer, which converts each word into a unique integer.
- **Padding:** Sequences are padded to a uniform length to ensure consistent input size to the model.
- **Cleaning:** The dataset is cleaned to remove unnecessary characters, special symbols, and noise that could disrupt the learning process.

The final processed dataset had a vocabulary size of 6,086 tokens and was split into training, validation, and test sets, with the majority allocated for training the LSTM-based models.

B. Model Architecture

To assess different modeling approaches for Skyrim-inspired chatbot dialogues, three distinct neural architectures were implemented and evaluated:

Model 1: Baseline BiLSTM (Learned Embeddings)

This architecture serves as the baseline for comparing the effectiveness of other models. The model uses learned embeddings (not pre-trained), which are initialized randomly and then optimized during training. The architecture is designed for processing sequential data, and it employs Bidirectional LSTM layers to capture contextual information from both past and future words in a sentence.

Architecture Breakdown:

- **Embedding Layer:** This layer is responsible for transforming each word in the sequence into a dense vector representation. The embeddings are learned from scratch during training.
- **Bidirectional LSTM (2 layers):** The two Bidirectional LSTM layers each contain 512 units, allowing the model to capture dependencies in the data from both directions (forward and backward). This helps the model better understand context by leveraging future and past words in the sentence.
- **Dropout Layers:** Dropout is applied between LSTM layers to help prevent overfitting by randomly setting a fraction of the input units to zero during each update cycle. The rate is set to 0.3 in this case.
- **Dense Layer:** The final dense layer produces a probability distribution over the vocabulary using the softmax activation function. This layer outputs a predicted word or token from the vocabulary for each timestep.
- **Optimizer:** The Adam optimizer is used with a learning rate of 0.001, enabling adaptive optimization of the model weights.
- **Loss Function:** Sparse Categorical Cross-Entropy is used as the loss function to compare the predicted and true word sequences.

Model 2: BiLSTM with GloVe Embeddings

In this architecture, we enhance the baseline BiLSTM model by using pre-trained GloVe embeddings (300-dimensional). The GloVe embeddings are precomputed on a large corpus of text and represent a more semantically rich representation of words. The embedding layer is frozen (not trainable) during the model's training phase to retain the semantic relationships present in the GloVe embeddings.

Architecture Breakdown:

- **GloVe Embedding Layer:** The GloVe embeddings are loaded into the model's embedding layer, but the layer is set to non-trainable. This ensures that the pretrained embeddings, which capture semantic word relationships, are preserved without updating during training.
- **Bidirectional LSTM (2 layers):** Like the first model, this one includes two Bidirectional LSTM layers with 512 units each. This setup ensures that the model captures

context in both directions (forward and backward).

- **Dropout Layers:** Dropout is applied with a rate of 0.3 to avoid overfitting, which is especially important when using pretrained embeddings to prevent the model from memorizing specific patterns.
- **Dense Layer:** A dense layer with a softmax activation is applied to predict the next word or token in the sequence, based on the output from the LSTM layers.
- **Optimizer:** The Adam optimizer is used with the same learning rate of 0.001 as the baseline model.
- **Loss Function:** The model uses Sparse Categorical Cross-Entropy to compute the error between the predicted word sequences and the ground truth.

Model 3: DistilGPT2 (Transformer-based)

The third model uses DistilGPT2, a distilled version of the GPT-2 architecture that is more lightweight while retaining much of the performance benefits of the full GPT-2 model. This transformer-based model uses self-attention mechanisms to process the entire sequence in parallel, allowing it to capture long-range dependencies better than LSTM-based models. Fine-tuning is performed on the Skyrim dialogue dataset to adapt the pre-trained model for generating coherent responses.

Architecture Breakdown:

- **Pre-trained DistilGPT2:** This model starts with a DistilGPT2 decoder, a transformer model trained on large-scale text corpora. The model is fine-tuned using the Skyrim-inspired dialogue data to adjust the weights according to the specific language and context of the game.
- **Fine-Tuning Process:** Using the HuggingFace Trainer API, the model is fine-tuned on the tokenized dataset, which consists of dialogues from Skyrim. This helps the model learn how to generate contextually relevant responses based on the training data.
- **Self-Attention Mechanism:** DistilGPT2 utilizes self-attention to process input sequences in parallel, which allows the model to consider the entire sequence at once and understand the relationships between words across long distances. This makes it effective at generating fluent and coherent text.
- **Output Layer:** The model generates token predictions (next words or responses) using the language modeling objective.
- **Training Settings:** The model is trained using the Adam optimizer, and various training settings (e.g. batch size, learning rate, number of epochs) are specified using the TrainingArguments class from the HuggingFace library.

IV. RESULTS

A. Evaluation

The models Baseline BiLSTM, BiLSTM with GloVe embeddings, and DistilGPT2 were evaluated based on a set of performance metrics to assess their ability to generate coherent and contextually appropriate dialogue for Skyrim NPCs.

1. Accuracy

Accuracy was used as a primary evaluation metric during both training and validation to understand how well each

model predicted the correct words and responses. It represents the percentage of correct predictions out of the total predictions made by the model.

$$\text{Accuracy} = \frac{\text{Number of Correct Predictions}}{\text{Total Number of Predictions}} \times 100$$

Training Accuracy: Measures how well the model performed on the training dataset.

Validation Accuracy: Measures how well the model performed on new, unseen data.

These values help to understand if the model is overfitting (high training accuracy but low validation accuracy) or underfitting (low accuracy on both training and validation).

2. Perplexity

Perplexity is a metric commonly used to evaluate language models. It measures how well the model predicts the next word in a sequence, giving an indication of how coherent the generated text is. Lower perplexity indicates better performance.

The formula for perplexity is:

$$\text{Perplexity} = e^{\frac{1}{N} \sum_{i=1}^N \log P(w_i)}$$

Where:

- N is the number of words in the sequence.
- $P(w_i)$ is the probability assigned by the model to the i-th word.

B. Result Analysis

The performance of the three models Baseline BiLSTM, BiLSTM with GloVe embeddings, and DistilGPT2 was evaluated based on training accuracy, validation accuracy, perplexity, training time, and human evaluation. These models were trained on a Skyrim-themed dialogue dataset generated using AI-driven methods, which simulate the conversational patterns and narrative style of Skyrim NPCs. The dataset underwent preprocessing steps, including tokenization, padding, and cleaning, and was split into training, validation, and test sets for model evaluation.

The models were trained on a Google Colab T4 GPU, with the following configurations:

- LSTM and BiLSTM with GloVe: Trained for 50 epochs with a batch size of 64.
- DistilGPT2: Trained for 1 epochs with a batch size of 2

Model	Training Accuracy (%)	Validation Accuracy (%)	Perplexity Training	Perplexity Validation
Baseline BiLSTM	92.52	91.17	1.38	1.50
BiLSTM with GloVe	91.42	90.26	1.50	1.65

The Baseline BiLSTM model achieved the highest performance among the tested architectures, with a training accuracy of 92.52%, validation accuracy of 91.17%, and the lowest perplexity scores 1.38 on the training set and 1.50 on the validation set. These metrics indicate that the model was able to learn the conversational structure of the dataset effectively and generalize well to unseen data.

The BiLSTM model enhanced with GloVe embeddings achieved slightly lower performance, with a training accuracy of 91.42% and validation accuracy of 90.26%. The corresponding perplexity values were 1.50 (training) and 1.65 (validation), suggesting a minor decline in the model's ability to predict the next word accurately. This may be due to the fixed nature of GloVe embeddings, which may not fully capture the nuanced and context-sensitive nature of NPC dialogues.

The DistilGPT2 model, although trained for only 1 epoch with a batch size of 2, showed promising results. The model's loss decreased steadily throughout the training, from 1.092600 at step 50 to 0.236900 at step 25000. While its training accuracy and validation accuracy were not explicitly provided, the model demonstrated competitive performance, especially considering the very limited training duration.

Overall, while both models performed well, the Baseline BiLSTM outperformed the GloVe-enhanced version in all metrics, making it more suitable for this specific dialogue generation task. However, the GloVe model remains a viable option for environments where pre-trained embeddings are preferred for semantic alignment or domain adaptation.

V. CONCLUSION

This project demonstrates the feasibility of developing an immersive NPC chatbot for Skyrim using deep learning techniques. By evaluating three models Baseline BiLSTM, BiLSTM with GloVe embeddings, and DistilGPT2 we observed that each offers unique advantages suited to different deployment contexts.

The BiLSTM with GloVe embeddings achieved stable performance with lower computational requirements, making it suitable for offline or resource-constrained environments. In contrast, DistilGPT2 generated more coherent and contextually rich dialogues, making it ideal for real-time applications in high-end gaming engines where computational resources are readily available.

Future work may include integrating the chatbot into an actual game engine, extending the dataset to support multi-turn conversations for more dynamic interactions, and experimenting with larger transformer models such as GPT-3 or LLaMA to further enhance contextual understanding and dialogue realism.

VI. REFERENCES

- Pennington, J., Socher, R., & Manning, C. D. (2014). **GloVe: Global Vectors for Word Representation**. *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. <https://nlp.stanford.edu/projects/glove>
- Wolf, T., Debut, L., Sanh, V., et al. (2020). **Transformers: State-of-the-Art Natural Language Processing**. *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, 38–45. <https://huggingface.co/transformers>
- 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>
- Li, J., Monroe, W., & Jurafsky, D. (2016). *A Persona-Based Neural Conversation Model*. Proceedings of the 34th International Conference on Machine Learning (ICML 2016). <https://arxiv.org/abs/1603.06155>
- Zhou, L., Gao, J., Li, D., & Shum, H. Y. (2018). *The Design and Implementation of XiaoIce, an Empathetic Social Chatbot*. *Computational Linguistics*, 46(1), 53–93. https://doi.org/10.1162/coli_a_00368