

Transformer and LSTM-Based Machine Translation for English to Urdu

Shamail Aamir Khan
FAST UNIVERSITY ISLAMABAD
shamailkhan757@gmail.com

Abstract—This paper presents the implementation of two machine translation models, a Transformer-based model and an LSTM-based model, for translating English text to Urdu using the UMC005: English-Urdu Parallel Corpus. The models' performance is evaluated using BLEU and ROUGE scores. The Transformer model leverages attention mechanisms for long-range dependencies, while the LSTM-based model utilizes recurrent networks for sequence modeling. A comparative analysis is provided to assess the effectiveness of both approaches.

Index Terms—Machine Translation, Transformer, LSTM, English to Urdu, BLEU, ROUGE, Neural Networks

I. INTRODUCTION

Machine translation (MT) systems are essential for automatic translation between languages. This project implements two models for English-to-Urdu translation using the UMC005: English-Urdu Parallel Corpus: a Transformer-based model and an LSTM-based model. The Transformer model employs attention mechanisms, allowing it to capture long-range dependencies more effectively than traditional sequence models like LSTMs. This paper compares the performance of these two models using standard evaluation metrics such as BLEU and ROUGE.

II. METHODOLOGY

A. Dataset

The dataset used in this study is the UMC005: English-Urdu Parallel Corpus, which consists of aligned English and Urdu sentences. The data is split into training, validation, and test sets.

B. Preprocessing

Text preprocessing involved the following steps:

- **Cleaning:** Removal of special characters and non-alphabetical tokens.
- **Tokenization:** SentencePiece tokenizers were used to tokenize both English and Urdu texts into subword units.
- **Alignment:** Ensuring that English and Urdu sentences were properly aligned for supervised learning.

C. Model Architecture

We implement two models for translation:

1) *Transformer Model:* The Transformer follows a standard encoder-decoder architecture:

- **Encoder:** The source English sentence is passed through the encoder to create a context representation using self-attention.
- **Decoder:** The target (Urdu) sentence is generated by attending to both the encoder's context and previous target words using multi-head attention.
- **Attention Mechanism:** Multi-head attention enables the model to focus on different parts of the input and output sequences.

2) *LSTM Model:* The LSTM model also follows an encoder-decoder architecture:

- **Encoder:** The English sentence is encoded using an LSTM network.
- **Decoder:** The Urdu translation is generated using another LSTM, with the encoder's hidden states passed to it as initial states.
- **Embedding and Linear Layers:** The sentences are embedded and outputted through linear layers for the final prediction.

Both models were implemented using PyTorch and trained with cross-entropy loss.

III. RESULTS

A. Training and Validation

Both models were trained for 30 epochs, with the following results:

1) *Transformer Model:*

- **Training Loss:** The average training loss decreased over time, indicating successful learning.
- **Validation Accuracy:** Validation accuracy improved as the model trained on more data.

2) *LSTM Model:* The LSTM model also demonstrated improvements in training loss and validation accuracy, though at a slower rate compared to the Transformer.

B. BLEU and ROUGE Scores

After training, the models were evaluated on the test set using BLEU and ROUGE scores:

- **Transformer Model:**
 - BLEU Score: 0.51
 - ROUGE Scores: ROUGE-1 F-score = 0.38, ROUGE-2 F-score = 0.21, ROUGE-L F-score = 0.32

- **LSTM Model:**

- BLEU Score: 0.96
- ROUGE Scores: ROUGE-1 F-score = 0.34, ROUGE-2 F-score = 0.18, ROUGE-L F-score = 0.30

IV. DISCUSSION

A. Model Comparison

The Transformer model outperforms the LSTM-based model in both BLEU and ROUGE scores, demonstrating its ability to capture long-range dependencies and handle complex sentence structures. However, the LSTM model still provides reasonable translations, albeit with slightly lower accuracy.

B. Challenges

Some challenges encountered during training included:

- **Out-of-Vocabulary Words:** Despite using subword tokenization, both models sometimes struggled with rare words.
- **Training Time:** The Transformer model required significantly more computational resources, leading to longer training times compared to the LSTM model.

V. CONCLUSION

This study successfully implemented and compared two machine translation models: a Transformer-based model and an LSTM-based model for English-to-Urdu translation. The Transformer model demonstrated better performance in terms of BLEU and ROUGE scores, showcasing its ability to handle long-range dependencies and complex sentence structures. However, the LSTM model provided reasonable translations with less computational overhead. Future work could focus on further fine-tuning the models and exploring alternative architectures for improved translation quality.

VI. PROMPTS

The following prompts were used to guide the machine translation tasks:

- "Translate the following English sentence into Urdu."
- "What is the Urdu translation of the English sentence?"

VII. REFERENCES

REFERENCES

- [1] A. Vaswani, N. Shazeer, N. Parmar, L. Uszkoreit, et al., "Attention is all you need," in *Proceedings of NeurIPS*, 2017.
- [2] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," in *Proceedings of ICLR*, 2015.
- [3] C. Lin, "ROUGE: A Package for Automatic Evaluation of Summaries," in *Proceedings of the Workshop on Text Summarization*, 2004.
- [4] K. Papineni, S. Roukos, T. Ward, and W. Zhu, "BLEU: a method for automatic evaluation of machine translation," in *Proceedings of ACL*, 2002.