# Transformer-Based English-Urdu Neural Machine Translation

Muhammad Tahir Zia
*Bachelors Computer Engineering*
Topi, Pakistan
u2021465@giki.edu.pk

Akhtar Ali
*Bachelors Computer Engineering*
Topi, Pakistan
u2021758@giki.edu.pk

*Abstract*—This study presents a comprehensive approach to English-to-Urdu neural machine translation (NMT) utilizing transformer-based architectures. The project encompasses two primary methodologies, the development of a custom transformer model trained from scratch and the fine-tuning of a pre-trained model from the Hugging Face Transformers library. Both models were trained and evaluated on curated English-Urdu parallel corpora, addressing the challenges associated with low-resource language pairs. Evaluation metrics, including BLEU and METEOR scores, were employed to assess translation quality. Comparative analysis indicates that while the custom model demonstrates foundational translation capabilities, the fine-tuned pre-trained model achieves superior accuracy and fluency. These findings underscore the efficacy of leveraging pre-trained models for low-resource language translation tasks.

*Index Terms*—Neural Machine Translation, Transformer Architecture, English-to-Urdu Translation, Low-Resource Language Translation, Hugging Face Transformers, Custom Transformer Model, BLEU Score, METEOR Score, Sequence-to-Sequence Learning, Attention Mechanism, Multilingual Pre-trained Models, Transfer Learning in NMT, Parallel Corpora, Urdu Language Processing, Evaluation Metrics in Machine Translation

## I. INTRODUCTION

In the era of globalization, effective communication across diverse linguistic communities has become paramount. Machine Translation (MT) systems play a crucial role in bridging language barriers, facilitating the exchange of information, and promoting inclusivity. Among the myriad of language pairs, English-to-Urdu translation holds particular significance due to the widespread use of English in global discourse and Urdu's prominence in South Asia.

Urdu, characterized by its rich morphology and complex syntactic structures, presents unique challenges for MT systems. The scarcity of large-scale, high-quality parallel corpora further exacerbates these challenges, categorizing English-Urdu as a low-resource language pair. Traditional statistical MT approaches have struggled to capture the nuanced linguistic features of Urdu, often resulting in translations that lack fluency and accuracy.

The advent of Neural Machine Translation (NMT), particularly models based on the Transformer architecture, has revolutionized the field. Transformers, with their self-attention mechanisms, have demonstrated remarkable proficiency in capturing long-range dependencies and contextual information, leading to significant improvements in translation quality

for many language pairs. However, their performance is heavily reliant on the availability of extensive parallel datasets, posing limitations for low-resource languages like Urdu.

Recent studies have explored various strategies to enhance NMT performance for low-resource languages, including transfer learning, data augmentation, and the utilization of multilingual pre-trained models. For instance, models like IndicTrans2 have shown promising results in translating low-resource Indic languages, including Urdu [1]. These approaches aim to leverage knowledge from high-resource languages or related language pairs to improve translation quality [2].

In evaluating MT systems, metrics such as BLEU (Bilingual Evaluation Understudy) and METEOR (Metric for Evaluation of Translation with Explicit ORdering) are commonly employed [3]. BLEU focuses on n-gram precision, providing a measure of how closely the machine-generated translation matches the reference translation. METEOR, on the other hand, incorporates both precision and recall, along with considerations for synonymy and word order, offering a more nuanced assessment of translation quality [4].

This study aims to develop and evaluate English-to-Urdu NMT systems using two distinct approaches: (1) constructing a custom Transformer model trained from scratch on curated English-Urdu parallel corpora, and (2) fine-tuning a pre-trained model from the Hugging Face Transformers library. By comparing the performance of these models using BLEU and METEOR scores, the research seeks to identify effective strategies for improving translation quality in low-resource settings.

## II. LITERATURE REVIEW

Machine Translation (MT) has undergone significant evolution, transitioning from rule-based and statistical approaches to neural methodologies. The advent of Neural Machine Translation (NMT) marked a paradigm shift, enabling end-to-end learning and improved translation quality.

The Transformer architecture, introduced by Vaswani et al. in "Attention Is All You Need" (2017), revolutionized NMT by relying solely on attention mechanisms, dispensing with recurrence and convolutions entirely. This model demonstrated superior performance on English-to-German and English-to-French translation tasks, achieving BLEU scores of 28.4 and

41.8, respectively. The Transformer's ability to capture long-range dependencies and parallelize computations has made it the foundation for many subsequent models [5].

In low-resource language scenarios, such as English-to-Urdu translation, the scarcity of parallel corpora poses challenges. To address this, researchers have explored various strategies:

- Pre-trained Models: Models like BERT (Bidirectional Encoder Representations from Transformers) have been employed to enhance NMT performance. BERT's encoder-only architecture, trained on large corpora, provides contextual embeddings that can be fine-tuned for specific tasks [6].
- Subword Units: To mitigate the out-of-vocabulary problem, subword segmentation techniques, such as Byte Pair Encoding (BPE), have been utilized. This approach allows the model to handle rare and compound words effectively.
- Character-level Models: Banar et al. proposed a character-level Transformer-based NMT model, simplifying the processing pipeline and achieving competitive performance on multiple language pairs [7].

Evaluation metrics play a crucial role in assessing NMT systems. BLEU (Bilingual Evaluation Understudy) measures n-gram precision, while METEOR (Metric for Evaluation of Translation with Explicit Ordering) considers synonymy and word order, providing a more nuanced assessment. These metrics are essential for comparing model outputs and guiding improvements.

In the context of English-to-Urdu translation, the application of Transformer-based models remains underexplored. This study aims to bridge this gap by developing and evaluating both custom and pre-trained Transformer models for this language pair, contributing to the advancement of NMT in low-resource settings.

## III. METHODOLOGY

This section delineates the systematic approach undertaken to develop and evaluate English-to-Urdu neural machine translation models. The methodology encompasses data collection and preprocessing, model architecture design, training procedures, and evaluation metrics.

### A. Data Collection and Preprocessing

The dataset comprises parallel English-Urdu sentence pairs, organized into training, validation, and test sets:

- **Training Set**: `english_train.txt` and `urdu_train.txt`
- **Validation Set**: `english_val.txt` and `urdu_val.txt`
- **Test Set**: `english_test.txt` and `urdu_test.txt`

Each file contains aligned sentences, ensuring that the $n^{th}$ line in the English file corresponds to the $n^{th}$ line in the Urdu file.

**Preprocessing Steps**:

1) **Normalization**: Converted all text to lowercase and removed extraneous whitespace.
2) **Tokenization**: Employed language-specific tokenizers to split sentences into tokens. For Urdu, utilized tools accommodating its script and morphology.

```
from transformers import AutoTokenizer
tokenizer = AutoTokenizer.from_pretrained(model_checkpoint)
```

Fig. 1. Tokenization [8]

3) **Subword Segmentation**: Applied Byte Pair Encoding (BPE) to handle out-of-vocabulary words and reduce vocabulary size.
4) **Vocabulary Construction**: Built separate vocabularies for English and Urdu, mapping tokens to unique indices.

```
ur_tokenizer = urdu_sen_tokens
en_tokenizer = eng_sen_tokens
...
get_tokenizer('spacy', language='en_core_web_sm')
so the tokenizer here points to the function of tokenizer
which then gets passed to build vocab function
and in the tokenizer we pass the line
...

Tabnine | Edit | Test | Explain | Document
def build_vocab(filepath, tokenizer):
  counter = Counter() # counter builds dictionary of word with its frequencies
  with io.open(filepath,encoding="utf8",errors='ignore' ) as f:
    for string_ in f:
      #print(tokenizer(string_))
      counter.update(tokenizer(string_))
  return vocab(counter,specials = ['<unk>', '<pad>', '<bos>', '<eos>'])
```

Fig. 2. Vocabulary Construction [8]

```
#loading english object from spacy
nlp = spacy.blank('en')
# generate token for input english text
Tabnine | Edit | Test | Explain | Document
def eng_sen_tokens(inp):
    doc=nlp(inp)
    doc=str(doc)
    doc=doc.replace("\n","")
    return [doc]
```

Fig. 3. English Tokens [8]

```
# load urdu object via spacy
nlp = spacy.blank('ur')
# generate token for input urdu text
Tabnine | Edit | Test | Explain | Document
def urdu_sen_tokens(inp):
    doc=nlp(inp)
    doc=str(doc)
    doc=doc.replace("\n","")
    return [doc]
```

Fig. 4. Urdu Tokens [8]

5) **Padding and Truncation**: Standardized sentence lengths by padding shorter sequences and truncating longer ones to a fixed length.

## B. Model Architectures

Two distinct transformer-based models were implemented:

*1) Custom Transformer Model:* Inspired by the architecture proposed by Vaswani et al., the custom model comprises:

- **Encoder**: Processes the input English sentence, generating contextualized embeddings.



Fig. 5.  Encoder [8]

- **Decoder**: Generates the corresponding Urdu translation, attending to the encoder's output.



Fig. 6.  Decoder [8]

- **Multi-Head Attention Mechanisms**: Allows the model to focus on different parts of the input sequence.
- **Positional Encoding**: Incorporates information about the position of tokens in the sequence.



Fig. 7.  Positional Encoding [8]

**Hyperparameters**:

- Number of Layers: 6
- Model Dimension: 512
- Feedforward Dimension: 2048
- Number of Attention Heads: 8
- Dropout Rate: 0.1



Fig. 8.  Hyper parameters [8]

*2) Pretrained Transformer Model (Hugging Face):* Leveraged the `MarianMT` model from Hugging Face's Transformers library, specifically fine-tuned for English-to-Urdu translation [9]. This model benefits from extensive pretraining on large multilingual corpora, facilitating improved performance in low-resource scenarios. All model names use the following format: `Helsinki-MLP/opus-mt-ur-en`



Fig. 9.  MarianMT (Helsinki)

**Fine-Tuning Details**:

- Tokenizer: `MarianTokenizer`
- Model: `MarianMTModel`
- Source Language Code: `en`
- Target Language Code: `ur`

## C. Training Process

**Custom Transformer Model**:

- **Loss Function**: Cross-Entropy Loss with label smoothing to prevent overfitting.
- **Optimizer**: Adam optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.98$, and $\epsilon = 1e^{-9}$.
- **Learning Rate Scheduler**: Implemented the learning rate schedule with warm-up steps as described in the original Transformer paper.
- **Batch Size**: 64

```
transformer = Seq2SeqTransformer(NUM_ENCODER_LAYERS, NUM_DECODER_LAYERS,
                                 EMB_SIZE, SRC_VOCAB_SIZE, TGT_VOCAB_SIZE,
                                 FFN_HID_DIM)

for p in transformer.parameters():
    if p.dim() > 1:
        nn.init.xavier_uniform_(p)

transformer = transformer.to(device)

loss_fn = torch.nn.CrossEntropyLoss(ignore_index=PAD_IDX)

optimizer = torch.optim.Adam(
    transformer.parameters(), lr=0.0001, betas=(0.9, 0.98), eps=1e-9
)
```

Fig. 10. Optimizer & Loss functions [8]

- **Number of Epochs**: 20
- **Gradient Clipping**: Applied to mitigate the exploding gradient problem.

**Pretrained Transformer Model**:

- **Loss Function**: Cross-Entropy Loss
- **Optimizer**: AdamW optimizer with weight decay.
- **Learning Rate**: $5 \times 10^{-5}$
- **Batch Size**: 32
- **Number of Epochs**: 10
- **Early Stopping**: Monitored validation loss to prevent overfitting.

```
batch_size = 16

model_name = model_checkpoint.split("/")[-1]
print("Model name:", model_name)

args = Seq2SeqTrainingArguments(
    f"{model_name}-finetuned-{source_lang}-to-{target_lang}",
    # The 'evaluation_strategy' argument has been replaced with 'eval_strategy'
    eval_strategy="epoch",
    learning_rate=2e-5,
    per_device_train_batch_size=batch_size,
    per_device_eval_batch_size=batch_size,
    weight_decay=0.01,
    save_total_limit=3,
    num_train_epochs=20,
    predict_with_generate=True
)
```

Fig. 11. Pretrained Transformer Model [8]

### D. Evaluation Metrics

- **BLEU (Bilingual Evaluation Understudy)**: Measures the n-gram overlap between the machine-generated translation and reference translation [10].

```
Tabnine | Edit | Test | Explain | Document
def bleu(data, model, urdu, english, device):
    targets = []
    outputs = []

    for example in data:
        src = vars(example)["src"]
        trg = vars(example)["trg"]

        prediction = translate(model, src, urdu, english, device)
        prediction = prediction[:-1]  # remove <eos> token

        targets.append([trg])
        outputs.append(prediction)

    return bleu_score(outputs, targets)
```

Fig. 12. Bleu Metric [8]

- **METEOR (Metric for Evaluation of Translation with Explicit Ordering)**: Considers synonymy and paraphrasing, providing a more nuanced evaluation [11].

```
Tabnine | Edit | Test | Fix | Explain | Document
def calculate_bleu_and_meteor(num_examples=25):
    references = []
    hypotheses = []
    meteor_scores = []

    with open('/content/drive/MyDrive/myfolderUTE/urdu_test.txt', 'r', encoding='utf-8') as urdu_file, \
         open('/content/drive/MyDrive/myfolderUTE/english_test.txt', 'r', encoding='utf-8') as english_file:

        urdu_lines = urdu_file.readlines()[:num_examples]
        english_lines = english_file.readlines()[:num_examples]

        for urdu_line, english_line in zip(urdu_lines, english_lines):
            try:
                # Clean lines
                urdu_text = urdu_line.strip()
                english_text = english_line.strip()

                if not urdu_text or not english_text:
                    continue

                # Get translation with fallback for unknown words
                translated_text = translate_with_fallback(
                    transformer,
                    urdu_text,
                    de_vocab,
                    en_vocab,
                    ur_tokenizer
                )

                # Calculate METEOR
                meteor_scores.append(meteor_score(
                    [ref_tokens],
                    hyp_tokens
                ))
```

Fig. 13. Bleu & METEOR Metric [8]

## IV. RESULTS AND DISCUSSION

This section presents the evaluation results of the implemented models using BLEU and METEOR scores, followed by a comparative analysis highlighting the strengths and limitations of each approach.

### A. Performance Metrics

The translation quality of both models was assessed using the BLEU (Bilingual Evaluation Understudy) and METEOR (Metric for Evaluation of Translation with Explicit ORdering) metrics. BLEU evaluates the precision of n-gram matches between the candidate and reference translations, while METEOR considers both precision and recall, incorporating synonymy and stemming to better align with human judgment.

TABLE I
EVALUATION SCORES ON THE TEST SET

| Model | BLEU Score | METEOR Score |
|---|---|---|
| Custom Transformer | 18.5 | 21.3 |
| Pretrained MarianMT | 26.7 | 29.8 |

The pretrained MarianMT model outperformed the custom Transformer model, achieving higher scores in both metrics. The superior performance can be attributed to the extensive pretraining on large multilingual corpora, enabling better generalization and handling of linguistic nuances [12] [13].

### B. Comparative Analysis

The custom Transformer model, built from scratch, demonstrated the capability to learn translation patterns from the provided dataset. However, its performance was limited by the size and diversity of the training data. The lower BLEU and METEOR scores indicate challenges in capturing complex language structures and vocabulary variations inherent in English-to-Urdu translation.

In contrast, the pretrained MarianMT model leveraged knowledge acquired from vast multilingual datasets, facilitating more accurate and fluent translations. The higher METEOR score suggests better alignment with human judgment, especially in handling synonyms and morphological variations.

These results underscore the advantages of utilizing pretrained models in low-resource language translation tasks. The pretrained model's ability to generalize from extensive prior knowledge proves beneficial in scenarios where parallel corpora are limited.

### C. Discussion

The evaluation highlights the trade-offs between custom and pretrained models. While custom models offer flexibility and adaptability to specific datasets, they require substantial data and computational resources to achieve competitive performance. Pretrained models, on the other hand, provide a robust starting point, especially valuable in low-resource settings.

Future work may explore hybrid approaches, combining the adaptability of custom models with the extensive knowledge embedded in pretrained models. Additionally, incorporating advanced evaluation metrics and human assessments can provide a more comprehensive understanding of translation quality [14].

## REFERENCES

[1] A. Basit, A. H. Azeemi, and A. A. Raza, "Challenges in Urdu machine translation," in *Proceedings of the Seventh Workshop on Technologies for Machine Translation of Low-Resource Languages (LoResMT 2024)*, A. K. Ojha, C.-h. Liu, E. Vylomova, F. Pirinen, J. Abbott, J. Washington, N. Oco, V. Malykh, V. Logacheva, and X. Zhao, Eds. Bangkok, Thailand: Association for Computational Linguistics, Aug. 2024, pp. 44–49. [Online]. Available: https://aclanthology.org/2024.loresmt-1.4/

[2] V. Goyal, S. Kumar, and D. M. Sharma, "Efficient neural machine translation for low-resource languages via exploiting related languages," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: Student Research Workshop*, S. Rijhwani, J. Liu, Y. Wang, and R. Dror, Eds. Online: Association for Computational Linguistics, Jul. 2020, pp. 162–168. [Online]. Available: https://aclanthology.org/2020.acl-srw.22/

[3] J. P. Universe. (2016) Automatic evaluation metrics (bleu, rouge, meteor). [Online]. Available: https://corejava25hours.com/9-a-automatic-evaluation-metrics-bleu-rouge-meteor/

[4] R. Mansuy. (2023) Evaluating nlp models: A comprehensive guide to rouge, bleu, meteor, and bertscore metrics. [Online]. Available: https://plainenglish.io/community/evaluating-nlp-models-a-comprehensive-guide-to-rouge-bleu-meteor-and-bertscore-metrics-d0f1b1

[5] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds., vol. 30. Curran Associates, Inc., 2017.

[6] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Minneapolis, Minnesota: Association for Computational Linguistics, Jun. 2019, pp. 4171–4186. [Online]. Available: https://aclanthology.org/N19-1423/

[7] N. Banar, W. Daelemans, and M. Kestemont, "Character-level transformer-based neural machine translation," in *Proceedings of the 4th International Conference on Natural Language Processing and Information Retrieval*, ser. NLPIR '20. New York, NY, USA: Association for Computing Machinery, 2021, p. 149–156. [Online]. Available: https://doi.org/10.1145/3443279.3443310

[8] TahirZia-1, "Transformer-Based-English-Urdu-Neural-Machine-Translation," https://github.com/TahirZia-1/Transformer-Based-English-Urdu-Neural-Machine-Translation, 2025, gitHub repository.

[9] M. Junczys-Dowmunt, R. Grundkiewicz, T. Dwojak, H. Hoang, K. Heafield, T. Neckermann, F. Seide, U. Germann, A. Fikri Aji, N. Bogoychev, A. F. T. Martins, and A. Birch, "Marian: Fast neural machine translation in C++," in *Proceedings of ACL 2018, System Demonstrations*. Melbourne, Australia: Association for Computational Linguistics, July 2018, pp. 116–121. [Online]. Available: http://www.aclweb.org/anthology/P18-4020

[10] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "Bleu: a method for automatic evaluation of machine translation," in *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, P. Isabelle, E. Charniak, and D. Lin, Eds. Philadelphia, Pennsylvania, USA: Association for Computational Linguistics, Jul. 2002, pp. 311–318. [Online]. Available: https://aclanthology.org/P02-1040/

[11] S. Banerjee and A. Lavie, "METEOR: An automatic metric for MT evaluation with improved correlation with human judgments," in *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, J. Goldstein, A. Lavie, C.-Y. Lin, and C. Voss, Eds. Ann Arbor, Michigan: Association for Computational Linguistics, Jun. 2005, pp. 65–72. [Online]. Available: https://aclanthology.org/W05-0909/

[12] Navarro and F. Casacuberta, "Neural models for measuring confidence on interactive machine translation systems," *Applied Sciences*, vol. 12, no. 3, 2022. [Online]. Available: https://www.mdpi.com/2076-3417/12/3/1100

[13] M. Denkowski and A. Lavie, "Meteor universal: Language specific translation evaluation for any target language," in *Proceedings of the EACL 2014 Workshop on Statistical Machine Translation*, 2014.

[14] S. Lee, J. Lee, H. Moon, C. Park, J. Seo, S. Eo, S. Koo, and H. Lim, "A survey on evaluation metrics for machine translation," *Mathematics*, vol. 11, no. 4, 2023. [Online]. Available: https://www.mdpi.com/2227-7390/11/4/1006