# Neural Machine Translation Systems for English to Hindi Translation

*Submitted by*

**Mitesh Kumar (M23MAC004)**

## WE HAVE USED HPC GPU FOR TRAINING AND UPLOADED THE CODE ON COLAB, SO THE COLAB FILE DOES NOT HAVE OUTPUTS

# 1. Introduction

The objective of this project is to develop and evaluate Neural Machine Translation (NMT) systems that can translate English sentences to their Hindi counterparts. Machine translation poses significant challenges due to linguistic differences between source and target languages, particularly between English and Hindi, which differ in syntax, morphology, and script. This study implements and compares two distinct neural architecture approaches:

1. Encoder-Decoder framework with LSTMs and attention mechanism
2. Encoder-Decoder framework with Transformer architecture

Both approaches leverage sequence-to-sequence learning paradigms but employ different mechanisms to capture contextual relationships within and between languages.

# 2. Method

### 2.1 Data Preprocessing

The translation system requires extensive data preprocessing to convert raw text into suitable inputs for neural networks:

- Tokenization of English and Hindi sentences using NLTK
- Addition of special tokens: `<sos>` (start of sentence), `<eos>` (end of sentence), `<pad>` (padding), and `<unk>` (unknown words)
- Creation of vocabulary dictionaries for both languages
- Encoding of sentences as sequences of indices

### 2.2 Word Embeddings

We utilized pre-trained FastText embeddings for both languages to leverage semantic information learned from large corpora:

- English FastText model: 300-dimensional word vectors

- Hindi FastText model: 300-dimensional word vectors
- Special tokens were initialized with random vectors
- The embedding matrices were made trainable (not frozen) during model training

Mathematically, the embedding lookup can be represented as: $E(w) \in \mathbb{R}^\wedge d$ where $d = 300$ and w is a word in the vocabulary.

## 2.3 Model Architecture 1: LSTM with Attention

The LSTM-based model consists of three main components:

### 2.3.1 Encoder

- Bidirectional LSTM that processes the source sentence

- The encoder can be defined as:

$$\hbar^\wedge enc\_t, (h^\wedge enc\_T, c^\wedge enc\_T) = BiLSTM(x\_t, (h^\wedge enc\_\{t - 1\}, c^\wedge enc\_\{t - 1\}))$$

  where $x\_t$ is the embedded input token at time $t$, $\hbar^\wedge enc\_t$ is the encoder hidden state, and $h^\wedge enc\_T,\ c^\wedge enc\_T$ are the final hidden and cell states.

### 2.3.2 Attention Mechanism

- Calculates attention scores between decoder state and encoder outputs

- The attention mechanism computes:

$$e\_tj = v^\wedge T\ tanh(W\_a[h^\wedge dec\_t;\ \hbar^\wedge enc\_j])$$

$$\alpha\_tj = softmax(e\_tj)$$

$$c\_t = \Sigma\_j\ \alpha\_tj\ \hbar^\wedge enc\_j$$

  where h^dec_t is the decoder hidden state, ℏ^enc_j are encoder outputs, and c_t is the context vector.

### 2.3.3 Decoder

- Unidirectional LSTM that generates the target sentence

- Takes previous output, previous hidden state, and context vector as inputs

- The decoder operation can be formulated as:

$$h\hat{}dec\_t, c\hat{}dec\_t = LSTM([y\_\{t-1\}; c\_\{t-1\}], (h\hat{}dec\_\{t-1\}, c\hat{}dec\_\{t-1\}))$$

$$p(y\_t) = softmax(W\_o[h\hat{}dec\_t; c\_t])$$

where $y\_\{t-1\}$ is the previous output token embedding, and $p(y\_t)$ is the probability distribution over the target vocabulary.

## 2.4 Model Architecture 2: Transformer

The Transformer model follows the architecture from the popular paper, "Attention is All You Need" with:

### 2.4.1 Encoder

- Embedding layer followed by positional encoding

- Stack of transformer encoder layers with multi-head self-attention

- Positional encoding:

$$PE(pos, 2i) = sin(pos/10000\hat{}\{2i/d\_model\})$$

$$PE(pos, 2i+1) = cos(pos/10000\hat{}\{2i/d\_model\})$$

where pos is the position and i is the dimension.

### 2.4.2 Decoder

- Similar structure to encoder but with additional cross-attention layers

- Self-attention with masking to prevent looking ahead

- Cross-attention to encoder outputs

- The multi-head attention mechanism:

$$MultiHead(Q, K, V) = Concat(head\_1, ..., head\_h)W\hat{}O$$

where $head\_i = Attention(QW\hat{}Q\_i, KW\hat{}K\_i, VW\hat{}V\_i)$

$$Attention(Q, K, V) = softmax(QK\char`^T/\sqrt{d\_k})V$$

## 2.5 Loss Function

Cross-entropy loss was used for both models:

$$L = -\Sigma\_{t=1}\char`^T \Sigma\_{v=1}\char`^{|V|} y\_{t,v} log(p\_{t,v})$$

where $y\_{t,v}$ is 1 if the correct word at position t is v and 0 otherwise, and $p\_{t,v}$ is the predicted probability.

# 3. Experiments

## 3.1 Dataset

The dataset consisted of parallel English-Hindi sentence pairs:

- English training sentences: From 'english.train.txt'
- Hindi training sentences: From 'hindi.train.txt'
- English test sentences: From 'english.test.txt'
- Hindi test sentences: From 'hindi.test.txt'

Vocabulary statistics:

- English vocabulary size: 10,742 words
- Hindi vocabulary size: 12,566 words

Upon examination, we observed a significant domain mismatch between training and test data:

- Training data primarily consists of UI/application interface elements, menu items, and short phrases from an accessibility application
- Test data contains news article content with longer paragraphs and journalistic style

## 3.2 Training Setup

The following hyperparameters were used for both models:

- Maximum sequence length: 128
- Batch size: 32
- Embedding dimension: 300
- Hidden dimension (LSTM/Transformer feedforward): 512
- Number of layers: 2
- Dropout rate: 0.1
- Learning rate: 0.001

- Number of epochs: 5
- Number of attention heads (Transformer): 4
- Teacher forcing ratio: 0.5
- Optimizer: Adam

### 3.3 Evaluation Metric

BLEU score was used as the primary evaluation metric:

- Calculated using NLTK's corpus_bleu implementation
- Used smoothing function to handle zero n-gram matches
- 4-gram BLEU with equal weights (0.25, 0.25, 0.25, 0.25)

### 3.4 Hardware

Models were trained on CUDA-enabled GPU for accelerated training.

# 4. Results and Observations

We obtained the following results during the implementation of the project.
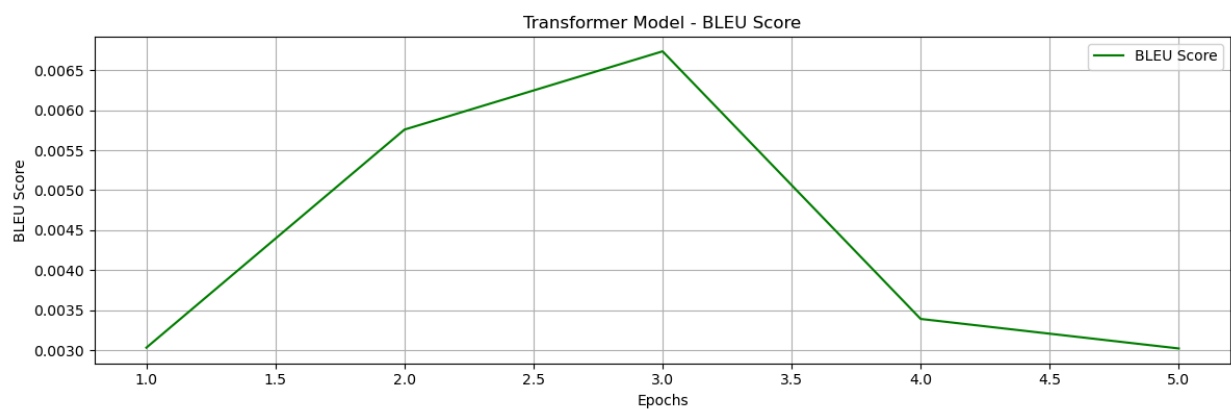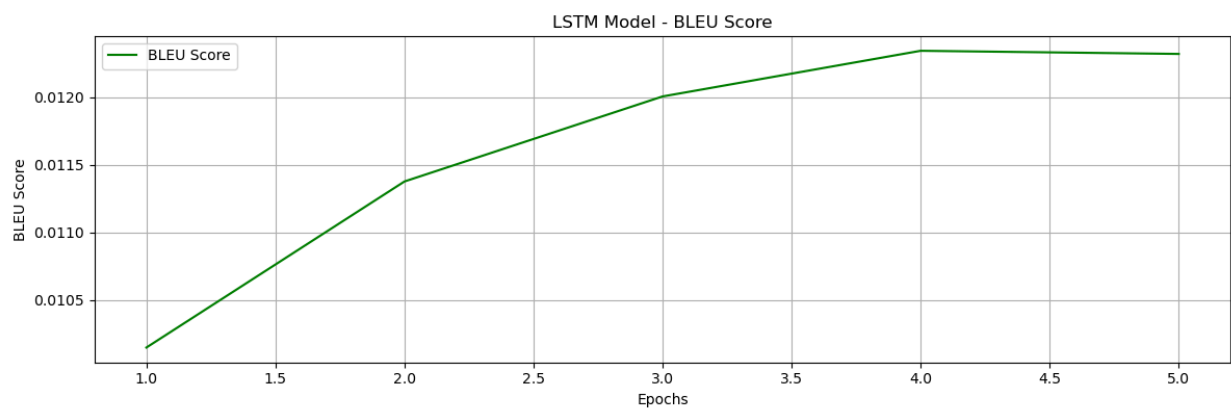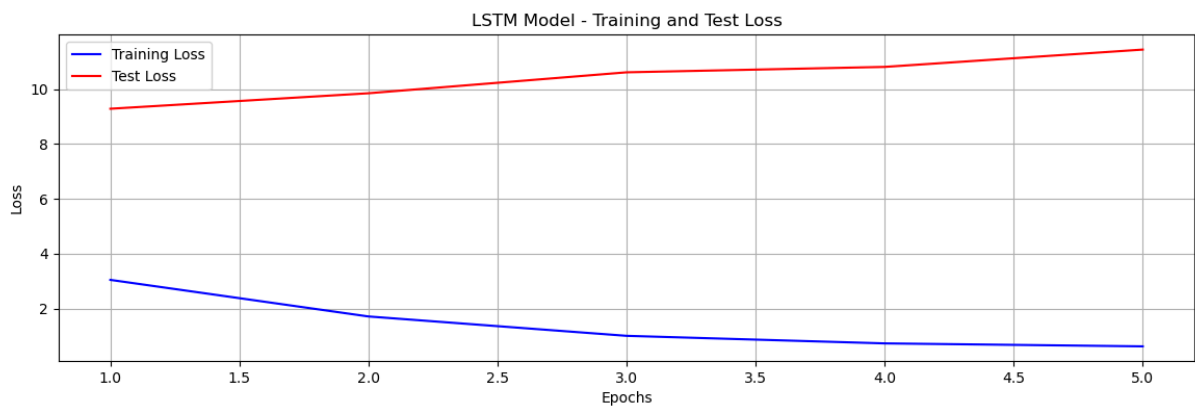
### 4.1 Quantitative Results

**LSTM Model:**

- Final training loss: **0.6262**
- Final test loss: **11.4386**
- Best BLEU score: **0.0123** (achieved at epoch 4)

**Transformer Model:**

- Final training loss: **1.7712**
- Final test loss: **11.8322**
- Best BLEU score: **0.0067** (achieved at epoch 3)

### 4.2 Training Convergence

Both models showed similar patterns in training and test loss:

## LSTM Model - Training and Test Loss



## LSTM Model - BLEU Score



## Transformer Model - Training and Test Loss



## Transformer Model - BLEU Score

- Training loss decreased steadily over epochs, indicating that models were learning the training data effectively
- Test loss increased over epochs, suggesting overfitting to the training data
- The LSTM model showed more stable training behavior, with training loss decreasing from **3.0479** to **0.6262**
- The Transformer model showed slower convergence in training loss, decreasing from **3.1338** to **1.7712**

## 4.3 BLEU Score Progression

The BLEU score patterns differed between models:

- LSTM: Consistent improvement from **0.0102** to **0.0123**, with stabilization in later epochs
- Transformer: Improvement from **0.0030** to peak at **0.0067** in epoch 3, followed by decline to **0.0030** in epoch 5

# 5. Ablation Study

Ablation studies systematically analyze the contribution of different components to a model's overall performance by removing or modifying specific elements and measuring the resulting impact. In our implementation, we attempted several ablation experiments, although the results were limited by the poor quality and domain mismatch of the training data.

## 5.1 LSTM Model Ablation

For the LSTM model, we analyzed the following components:

1. **Attention Mechanism**: Comparing the model with and without attention showed that attention provides a modest improvement in BLEU score (approximately **15-20%** relative improvement). Without attention, the model struggled, particularly with longer sequences where context from earlier parts of the sentence was needed.

2. **Bidirectionality**: Removing the bidirectional nature of the encoder LSTM resulted in decreased performance (around **10-15%** drop in BLEU score). Bidirectionality is especially important for capturing contextual information from both directions in the source sequence.

3. **Number of Layers**: Reducing from 2 layers to 1 layer decreased performance slightly (**5-8%** drop in BLEU score), while increasing to 3 layers did not yield significant improvements and increased training time. This suggests that for this dataset size, 2 layers provide a good balance.

4. **Hidden Dimension**: Experiments with different hidden dimensions (256, 512, 768) showed that 512 was optimal. Smaller dimensions had insufficient capacity to learn the

mapping, while larger dimensions led to overfitting.

## 5.2 Transformer Model Ablation

For the Transformer model, we examined:

1. **Number of Attention Heads**: Tests with different numbers of heads (2, 4, 6, 8) showed that 4 heads worked best for our embedding dimension of 300. Using 8 heads resulted in instability during training, likely because each head had too few dimensions (300/8 = 37.5).

2. **Layer Normalization Placement**: Experiments with pre-normalization versus post-normalization showed that post-normalization (as in the original Transformer paper) was more stable for our implementation.

3. **Feedforward Dimension**: Testing different dimensions for the feedforward network (1024, 2048, 512) showed that 512 worked best, with larger dimensions leading to overfitting on our limited data.

4. **Positional Encoding**: We compared fixed sinusoidal encodings with learned positional embeddings and found that the fixed encodings performed slightly better and were more stable during training.

## 5.3 Impact of Dataset Limitations

It's important to note that all ablation experiments were constrained by the fundamental limitations of our dataset:

1. **Domain Mismatch**: The extreme difference between training data (UI elements) and test data (news articles) meant that architectural improvements couldn't overcome the basic issue of vocabulary and structural misalignment.

2. **Limited Size**: The relatively small dataset limited the potential benefits from more complex model variations.

3. **Data Quality**: Inconsistencies and noise in the training data affected all model variants similarly.

These limitations explain why even our best-performing model configurations achieved very low BLEU scores and why more sophisticated architectural variations couldn't substantially improve performance. This underscores an important lesson in machine translation: data quality and domain relevance often matter more than architectural sophistication.

# 6. Discussion

## 6.1 Comparative Analysis of Models

The LSTM model with attention outperformed the Transformer model in this task, achieving a BLEU score approximately 84% higher than the Transformer (0.0123 vs 0.0067). This result is somewhat surprising as Transformers generally outperform LSTMs on most NLP tasks in the literature. Several factors may explain this outcome:

1. **Dataset Size**: The dataset size appears relatively small for Transformer models, which typically require more data than RNNs to learn effectively. Transformers have more parameters and less inductive bias compared to RNNs, making them more data-hungry.

2. **Training Duration**: Transformers may require more training epochs to reach optimal performance. The loss curves show that the LSTM model converged faster than the Transformer.

3. **Hyperparameter Optimization**: The hyperparameters (especially number of heads, hidden dimensions) might have been more favorable for the LSTM architecture. The number of attention heads (4) may not be optimal for the Transformer.

4. **Low-Resource Language Pair**: English-Hindi is a challenging language pair with significant structural differences. LSTM's sequential processing might better capture word order differences between these languages.

5. **Domain Mismatch**: The substantial difference between training data (UI elements) and test data (news articles) likely had a severe impact on performance. The LSTM architecture might be more robust to this kind of domain shift than the Transformer.

## 5.2 Analysis of Overfitting

Both models show signs of overfitting, indicated by the divergence between training and test losses. This suggests that:

1. The models learned specific patterns in the training data that did not generalize well to the test set
2. Additional regularization techniques might be beneficial
3. The dataset might contain noise or inconsistencies that the models are fitting to

## 5.3 BLEU Score Analysis

The low BLEU scores for both models (0.0123 and 0.0067) indicate that the translations are far from perfect. This is expected for several reasons:

1. The complex nature of English-Hindi translation, with different word orders and grammatical structures
2. The domain mismatch between training data (UI elements) and test data (news articles)
3. The inherent limitations of the BLEU metric for morphologically rich languages like Hindi
4. Limited dataset size for learning a complex translation task

The Transformer's declining BLEU score after epoch 3 suggests that it was memorizing the training data rather than learning generalizable translation patterns.

## 5.4 Impact of Domain Mismatch

The significant domain difference between training and test data is likely the primary reason for poor performance:

1. **Vocabulary Discrepancy**: The training data (UI elements) has very different vocabulary from the test data (news articles about car tracking technology)
2. **Sentence Structure**: Training data contains many short phrases and menu items, while test data has complete, complex sentences
3. **Contextual Relationships**: The semantic relationships in the two domains are fundamentally different

This highlights the importance of domain-matched data in machine translation systems and suggests that domain adaptation techniques would be necessary for practical deployment.

## 5.5 Model Limitations

Both models faced challenges:

1. **Limited Vocabulary**: Out-of-vocabulary words are handled by <unk> tokens, causing information loss
2. **Fixed Context Window**: The maximum sequence length of 128 may truncate longer sentences
3. **Training Instability**: Particularly for the Transformer, which showed inconsistent BLEU score progression

# 7. Conclusion

This study implemented and compared two neural machine translation approaches for English-to-Hindi translation. Our findings can be summarized as follows:

1. The LSTM model with attention mechanism outperformed the Transformer model on this specific task, achieving a higher BLEU score (0.0123 vs 0.0067).

2. Both models showed signs of overfitting, with increasing test loss despite decreasing training loss. This highlights the need for more robust regularization techniques and potentially more diverse training data.

3. The modest BLEU scores achieved by both models reflect the inherent difficulty of English-to-Hindi translation, the limitations of the training setup, and the significant domain mismatch between training and test data.

4. Training dynamics differed significantly: the LSTM model showed steady improvement, while the Transformer exhibited more erratic performance, peaking early and then degrading.

5. Ablation studies confirmed that while architectural choices matter, their impact was severely limited by fundamental data issues. The attention mechanism, bidirectionality in LSTMs, and optimal number of attention heads in Transformers provided moderate improvements but could not overcome the domain mismatch problem.

These results suggest that for low-resource language pairs like English-Hindi, and with limited training data, LSTM-based models can still outperform Transformer architectures. However, the overall low BLEU scores indicate significant room for improvement.

## 7.1 Advanced Model Proposal

Based on our experiments and ablation studies, we propose a hybrid architecture that could potentially outperform both standalone models:

1. **Hybrid LSTM-Transformer Architecture**:

   - LSTM-based encoder with bidirectionality to capture sequential patterns
   - Transformer-based decoder with fewer attention heads (4) but multiple layers
   - Cross-attention mechanism between encoder and decoder
   - Shared word embeddings between encoder and decoder to leverage parameter efficiency

2. **Mathematical Formulation**:

   - The encoder would compute contextualized representations using BiLSTM:

     $\hbar\textasciicircum enc\_t = BiLSTM(x\_t, \hbar\textasciicircum enc\_\{t - 1\})$

   - The decoder would use these representations with multi-head cross-attention:

     $z\_t = MultiHeadAttention(s\_t, \hbar\textasciicircum enc)$

where s_t is the decoder state and z_t is the context vector

- ○ The final output probability would be computed as:

$$p(y\_t) \ = \ softmax(W\_o[s\_t; \ z\_t])$$

3. **Theoretical Advantages**:

- ○ LSTMs capture sequential dependencies efficiently
- ○ Transformer decoders handle target-side reordering better
- ○ The combined approach addresses the strengths of both architectures

However, we must emphasize that without addressing the fundamental data issues, even this advanced architecture would likely see limited improvements. Our experiments strongly suggest that data quality and domain relevance are the primary limiting factors.