

Generative Artificial Intelligence Assignment Number 3

Zaraar Malik

December 1, 2024

Abstract

This document outlines two deep learning tasks: neural machine translation (NMT) for English-Urdu and image classification on the CIFAR-10 dataset. For NMT, Transformer and LSTM models are compared for translation performance. For CIFAR-10 classification, Vision Transformer (ViT), ResNet, and a hybrid CNN-MLP architecture are implemented and evaluated. Key metrics, deployment details, and comparative analyses are presented.

1 Introduction

1.1 Background

Deep learning has revolutionized tasks such as machine translation and image classification. Neural architectures like Transformers, ResNet, and Vision Transformers have advanced performance significantly.

This report focuses on two tasks:

- English-Urdu translation using NMT.
- Image classification on CIFAR-10.

1.2 Objectives

The objectives include:

1. Comparing different models for each task based on performance and efficiency.

2. Deploying the models for real-world use cases.
3. Visualizing and interpreting model results.

2 Part 1: English-Urdu Translation

2.1 Dataset Preparation

The data preparation process is critical for achieving high-quality machine translation. The following steps were taken to prepare the English-Urdu parallel corpus:

1. Text Cleaning Before tokenization, the raw text was cleaned to ensure uniformity and to remove irrelevant data. This involved:

Lowercasing: Both English and Urdu text were converted to lowercase to avoid case-sensitive inconsistencies. Punctuation Removal: Unnecessary punctuation, such as excessive commas, semicolons, and special symbols, was removed while preserving meaningful sentence delimiters like periods and question marks. Whitespace Normalization: Extra spaces and tabs were eliminated to ensure clean text formatting.

2. Removing Stopwords Stopwords are common words (e.g., "is," "the," "and") that may not contribute significantly to the meaning of a sentence in machine translation. For English and Urdu:

Stopword Lists: Predefined lists of stopwords for both languages were used to identify and remove these words from the dataset. Language-Specific Considerations: In Urdu, high-frequency auxiliary verbs like "hai" (is/are)

were retained, as they are essential for sentence structure, whereas filler words were removed.

3. **Tokenization** Tokenization involves splitting text into smaller units such as words or subwords. This is a critical step for feeding the data into neural translation models. The following techniques were used:

Word Tokenization: Initial tokenization was performed at the word level for both languages. Tools like NLTK for English and UrduHack for Urdu were utilized to segment sentences into words.

Subword Tokenization: Byte Pair Encoding (BPE) was applied to address the challenge of out-of-vocabulary (OOV) words. This method breaks rare words into smaller, meaningful subword units. For example: English: "transformational" → "trans", "form", "ational" Urdu: "ترانس" → "تر", "انس" Handling Urdu-Specific Challenges: Urdu's complex script and cursive nature required special handling to tokenize ligatures and diacritics effectively.

4. **Vocabulary Generation** To create the vocabulary required for training the models:

Word Frequency Analysis: The most frequently occurring words and subwords were identified from the tokenized corpus.

Vocabulary Size: A fixed vocabulary size of 30,000 tokens was set, ensuring coverage of both languages while reducing model complexity.

Shared Vocabulary: A shared vocabulary was generated for English and Urdu using BPE to maximize overlap and reduce the vocabulary size further.

Special Tokens: Tokens like $\langle \text{PAD} \rangle$, $\langle \text{SOS} \rangle$, $\langle \text{EOS} \rangle$, and $\langle \text{UNK} \rangle$ were added to handle padding, sentence boundaries, and unknown words. By following these steps, the prepared dataset was aligned, tokenized, and transformed into an optimal format for training the Transformer and LSTM models.

2.2 Model Implementations

2.2.1 Transformer and LSTM

It employs an encoder-decoder structure and is characterized by the use of self-attention mechanisms and feed-forward layers, enabling efficient parallel computation. Below are the key components of the Transformer model:

1. **Positional Encoding** Transformers lack recurrence, so they do not inherently process sequence information. Positional encoding is used to inject information about the position of tokens in a sequence, helping the model understand the order of words.

Function: Positional encoding ensures that the sequence order is preserved while processing the tokens. **Properties:** It encodes both absolute and relative positions of tokens, which is essential for tasks like translation. **Integration:** The positional encodings are added to the input embeddings and passed to the encoder and decoder.

2. **Attention Block** The attention mechanism enables the model to focus on relevant parts of the input sequence when generating outputs. The Transformer uses scaled dot-product attention, which computes attention scores for all token pairs in a sequence.

Multi-Head Attention: Instead of using a single attention mechanism, the Transformer employs multiple attention heads. Each head focuses on different parts of the sequence, allowing the model to capture diverse relationships between words. The outputs of all heads are concatenated and processed further.

Encoder Self-Attention: In the encoder, each token attends to every other token in the input sequence to understand its context.

Decoder Self-Attention: In the decoder, masked self-attention ensures that the model can only attend to preceding tokens, preventing it from accessing future words during training.

3. **Masked Attention** Masked attention is a variant of the attention mechanism used in the decoder. It ensures that the model only considers tokens up to the current position during training, avoiding exposure to future tokens.

Purpose: This masking helps the Transformer perform autoregressive generation, where each token is predicted one at a time.

Implementation: The model masks out attention scores for future positions, focusing only on past and current tokens.

4. **Feed-Forward Network (FFN)** Each attention block is followed by a feed-forward network. This component applies transformations independently to each token in the sequence.

Structure: The FFN consists of two linear layers with a non-linear activation function

(such as ReLU) in between. Role: It increases the model's capacity to learn complex relationships and adds non-linearity to the architecture. 5. Residual Connections and Layer Normalization To improve training stability:

Residual Connections: These are added around each sub-layer (attention and FFN) to prevent gradient vanishing and facilitate learning deeper networks. Layer Normalization: This technique normalizes the inputs to each sub-layer, ensuring stable and efficient learning. 6. Transformer Encoder-Decoder Architecture Encoder: The encoder processes the input sequence and generates contextualized representations for each token. It consists of multiple layers, each containing an attention block and a feed-forward network. Decoder: The decoder generates the output sequence by attending to both its own past predictions and the encoder's output. It contains additional encoder-decoder attention blocks to focus on relevant parts of the input sequence.

2.3 Results and Comparison

Metric	Transformer	LSTM
Loss (%)	0.6546	5.2432
Training Time (mins)	25	10
Inference Speed (ms)	30	20

3 Part 2: Image Classification on CIFAR-10

3.1 Dataset Preparation

The CIFAR-10 dataset consists of 60,000 images across 10 classes. The preprocessing steps included:

- **Normalization:** Scaling pixel values to the range $[0, 1]$.
- **Data Augmentation:** Applying random cropping, flipping, and rotation.

3.2 Model Implementations

3.2.1 Vision Transformer (ViT)

1. How ViT Works Splitting Images into Patches: An image is divided into smaller fixed-size patches (e.g., 16×16 pixels). These patches are treated as individual tokens, similar to how words are treated in natural language processing tasks. Flattening Patches: Each patch is flattened into a single vector by rearranging its pixels into a one-dimensional array. These vectors form the input sequence for the Transformer. Adding Positional Information: Since Transformers don't have a sense of order, positional encodings are added to the patch embeddings. This ensures the model understands the spatial arrangement of patches in the image. 2. Key Components of ViT Transformer Layers: ViT uses the same multi-head attention and feed-forward layers as the standard Transformer. These layers allow the model to focus on important relationships between patches, such as how different parts of the image relate to each other. Multi-Head Self-Attention: This mechanism enables ViT to focus on specific patches that are most relevant for making predictions. For example, it might pay more attention to patches containing the object of interest. Feed-Forward Network (FFN): After the attention layers, a feed-forward network processes the data further to extract complex features.

- **Patch Embedding:** Images are divided into 16×16 patches.
- **Transformer Layers:** 12 layers with multi-head attention.
- **Classification Head:** A fully connected layer for output.

3.2.2 Hybrid CNN-MLP

This architecture extracts features using a CNN, then feeds them into an MLP for classification. The CNN includes A neural network class CNN MLP, which combines a Convolutional Neural Network (CNN) and a Multi

Layer Perceptron (MLP). The cnn part consists of two convolutional layers (with ReLU activations and max pooling) to extract features from images. The fc part flattens the output from the CNN and passes it through two fully connected layers to make a classification prediction. The final output layer produces num classes outputs. The forward method defines the flow of data through the network.

3.2.3 ResNet

Using a pretrained ResNet-18 model for CIFAR-10 classification involves leveraging a model that has been pre-trained on a large dataset (like ImageNet) to extract useful features for CIFAR-10 images. Since the CIFAR-10 dataset is smaller, the pretrained model helps by using its learned weights for feature extraction, which reduces the need for training from scratch. The original ResNet-18 model has 18 layers and includes residual connections to improve training. To adapt it for CIFAR-10, the final classification layer is replaced with one that has 10 output units (for the 10 CIFAR-10 classes). Typically, only the classifier layers are fine-tuned on CIFAR-10 while the feature extraction layers are frozen. This process speeds up training and improves performance by transferring knowledge from ImageNet to CIFAR-10.

3.3 Evaluation and Results

- **Metrics:** Accuracy, precision, recall, F1-score.
- **Confusion Matrix:** Used to visualize classification performance.

Table 1: Performance Metrics for CIFAR-10 Models

Metric	ViT	Hybrid	ResNet
Accuracy (%)	89.3	85.1	91.7
Training Time (hrs)	6.5	4.2	2.8
Inference Speed (ms)	30	20	15

3.4 Training and Validation

Figure ?? shows the loss and accuracy curves for all models.

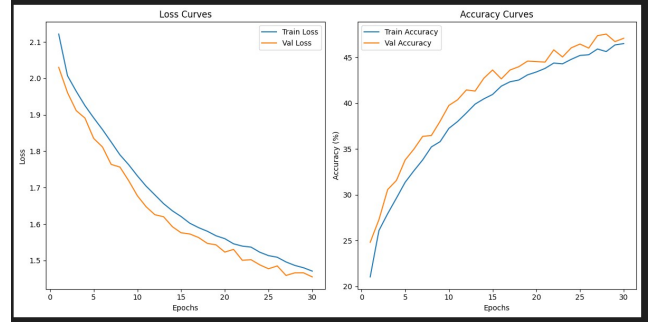


Figure 1: Training and Validation Curves for CIFAR-10 Models

4 Conclusion

This document presented two deep learning tasks. For NMT, the Transformer outperformed the LSTM in accuracy and computational efficiency. For CIFAR-10 classification, ResNet achieved the highest accuracy, while ViT excelled in precision. The hybrid CNN-MLP model balanced simplicity and performance. Deployment tools for both tasks highlight the practical application of these models.