# Evaluating Multilingual And Language Specific Transformers For English - Hindi Semantic Alignment

A PROJECT REPORT

SUBMITTED IN PARTIAL FULFILLMENT OF THE REQUIREMENTS
FOR THE AWARD OF THE DEGREE
OF

## MASTER OF TECHNOLOGY
IN
**Information Technology**

Submitted by

## VIKAS GUPTA (2K23/ITY/03)

Under the supervision of

## Dr. Ritu Agarwal



**Information Technology**
DELHI TECHNOLOGICAL UNIVERSITY
(Formerly Delhi College of Engineering)
Bawana Road, Delhi 110042

**MAY, 2025**

**DEPARTMENT OF INFORMATION TECHNOLOGY**
DELHI TECHNOLOGICAL UNIVERSITY
(Formerly Delhi College of Engineering)
Bawana Road, Delhi-110042

## CANDIDATE'S DECLARATION

I, **VIKAS GUPTA**, Roll No. 2K23/ITY/03 students of M.Tech (**Information Technology**), hereby declare that the project Dissertation titled "**Evaluating Multilingual And Language Specific Transformers For English - Hindi Semantic Alignment**" which is submitted by us to the Information Technology, Delhi Technological University, Delhi in partial fulfilment of the requirement for the award of degree of Master of Technology, is original and not copied from any source without proper citation. This work has not previously formed the basis for the award of any Degree, Diploma Associateship, Fellowship or other similar title or recognition.

Place: Delhi                                                    Vikas Gupta

Date: 25.05.2025

**DEPARTMENT OF INFORMATION TECHNOLOGY**
DELHI TECHNOLOGICAL UNIVERSITY
(Formerly Delhi College of Engineering)
Bawana Road, Delhi-110042

## <u>CERTIFICATE</u>

I hereby certify that the Project Dissertation titled "**Evaluating Multilingual And Language Specific Transformers For English - Hindi Semantic Alignment**" which is submitted by **Vikas Gupta**, Roll No. 2K23/ITY/03, Information Technology, Delhi Technological University, Delhi in partial fulfilment of the requirement for the award of the degree of Master of Technology, is a record of the project work carried out by the students under my supervision. To the best of my knowledge this work has not been submitted in part or full for any Degree or Diploma to this University or elsewhere.

Place: Delhi

Dr. Ritu Agarwal

Date: 25.05.2025

**SUPERVISOR**

**DEPARTMENT OF INFORMATION TECHNOLOGY**
DELHI TECHNOLOGICAL UNIVERSITY
(Formerly Delhi College of Engineering)
Bawana Road, Delhi-110042

## ACKNOWLEDGEMENT

We wish to express our sincerest gratitude to **Dr. Ritu Agarwal** for her continuous guidance and mentorship that she provided us during the project. She showed us the path to achieve our targets by explaining all the tasks to be done and explained to us the importance of this project as well as its industrial relevance. She was always ready to help us and clear our doubts regarding any hurdles in this project. Without her constant support and motivation, this project would not have been successful.

We would also like to extend our heartfelt thanks to the Head of Department **Dr. Dinesh K. Vishwakarma**, for their invaluable support and encouragement throughout this project.

Place: Delhi                                                                                   Vikas Gupta

Date: 25.05.2025

# Abstract

This research provides a detailed comparative evaluation of three prominent transformer-based language models: multilingual BERT (mBERT), XLM-ROBERTa, and the Hindi-specific L3Cube-HindBERT. The primary objective was to assess their effectiveness in learning and aligning cross-lingual semantic representations between English and Hindi. The study utilized the Bharat Parallel Corpus Collection (BPCC), a significant resource for Indian languages, to form the basis of this investigation. A synthetic classification task was designed to evaluate the models' ability to differentiate between genuine English-Hindi translated sentence pairs and randomly mismatched pairs, thereby gauging their capacity to capture semantic equivalence. While the core research paper focused on performance metrics and training dynamics without visualizations, this report will also touch upon the experimental scripts' capabilities for such visual analysis as part of a broader methodological discussion. The findings indicate that all three models are capable of aligning cross-lingual representations, though their learning trajectories and ultimate performance vary due to architectural and pretraining differences. Notably, mBERT demonstrated the most stable training convergence and achieved the best overall performance on the classification task, suggesting advantages from its extensive multilingual pretraining. XLM-ROBERTa showed slightly higher validation losses but strong performance, indicative of its robust pretraining regimen. L3Cube-HindBERT, while initially slower to converge, showed benefits of domain adaptation for Hindi. This study underscores the trade-offs between generalized multilingual models and language-specific architectures in the context of cross-lingual tasks.

# Contents

# List of Tables

# List of Figures

# Chapter 1

# INTRODUCTION

## 1.1  Background: The Rise of Multilingual NLP

Natural Language Processing (NLP) has undergone a transformative journey, moving from rule-based systems to sophisticated deep learning models. This evolution has significantly altered how humans interact with technology and access information globally. In an increasingly interconnected digital world, the capability to process, understand, and generate information across diverse languages is no longer a specialized niche but a fundamental global requirement. Modern international collaborations, global business operations, and cross-cultural educational endeavors heavily rely on the seamless sharing of information across linguistic divides. Consequently, there is a burgeoning demand for robust NLP methods that can effectively dismantle language barriers, fostering continuous research into solutions that enable effective communication irrespective of the languages spoken by individuals. Transformer models, with their attention mechanisms, have been at the forefront of these advancements, demonstrating remarkable capabilities in understanding context and nuances in language.

## 1.2  Challenges and Opportunities

India presents a particularly complex and fascinating linguistic environment, with 22 official languages and hundreds of dialects spoken across the nation. This rich linguistic diversity, while culturally vibrant, poses significant challenges for NLP development. Among these languages, English and Hindi play prominent roles, with frequent interaction in daily life, commerce, education, and governance. This necessitates the development of language processing tools that can accurately handle both languages and facilitate smoother communication and information exchange between their speakers. Creating reliable NLP systems for English and Hindi is not merely a technical hurdle; it represents a critical step towards enhancing digital information accessibility and inclusivity for a vast population in India and, by extension, for multilingual communities worldwide. However, achieving genuine understanding between languages as distinct as English (an Indo-European Germanic language) and Hindi (an Indo-European Indo-Aryan language) is far from straightforward. Beyond grammatical and structural dissimilarities, issues like semantic ambiguity and the influence of cultural context significantly impede effective translation and understanding. Idiomatic expressions, metaphors, and varying conventions for politeness can carry vastly different meanings in each language, and many words or phrases lack precise one-to-one equivalents. A truly effective cross-lingual system must therefore go beyond literal word-for-word translation to capture the intended meaning,

emotion, and cultural subtleties embedded in the text. This highlights the need for models that can process complex, context-sensitive representations and adapt to nuanced cultural variations. Another significant challenge, particularly for Indian languages, is the scarcity of high-quality, large-scale translated text collections (parallel corpora). Existing resources often focus on specific domains like news articles or legal documents, failing to represent the full spectrum of language use found in everyday conversations, social media, or technical writing. To address this, initiatives to build and organize larger, more balanced translated text collections, such as the Bharat Parallel Corpus Collection (BPCC), are crucial. Such resources are vital for rigorously testing and improving language technologies across diverse topics and communication styles.

## 1.3 Problem Statement

The core challenge addressed in this research is the development and comparative evaluation of computational methods capable of understanding and connecting the meanings of words and sentences in English and Hindi. This involves tackling inherent complexities such as limited digital resources for Hindi compared to English, cultural differences influencing language use, complex grammatical structures, lexical ambiguity (words with multiple meanings), and different writing systems (Latin for English, Devanagari for Hindi). Specifically, this study investigates how well state-of-the-art transformer-based models, with varying pretraining paradigms (broadly multilingual vs. language-specific), can learn cross-lingual semantic representations for English-Hindi pairs.

## 1.4 Research Objectives

The primary objectives of this research are: To systematically evaluate the performance of three transformer-based architectures—multilingual BERT (mBERT), XLM-ROBERTa, and L3Cube-HindBERT—in learning cross-lingual representations between English and Hindi. To investigate the models' ability to capture semantic equivalences in translation pairs from the BPCC parallel corpus. This is achieved through a synthetic classification task designed to discriminate between genuine and randomly paired sentences. To analyze and compare the performance trends of these models, focusing on training dynamics (e.g., convergence stability, loss progression) and standard evaluation metrics (accuracy, F1-score). To understand how architectural differences and pretraining strategies (multilingual vs. Hindi-specialized) influence learning trajectories and cross-lingual alignment capabilities. To contribute to the understanding of how computers can process and link languages that are structurally and culturally different, especially when resources may be imbalanced.

## 1.5 Significance of the Study

This study offers several significant contributions:

- Comparative Benchmarking: It provides a direct comparison of popular multilingual and a language-specific model for English-Hindi, a language pair of considerable importance but with fewer comprehensive benchmarks compared to high-resource European language pairs.

- Insight into Model Strengths: The research elucidates the differential strengths of mBERT's stable multilingual pretraining, XLM-ROBERTa's robust large-scale pretraining, and L3Cube-HindBERT's domain-specific adaptation for Hindi.

- Methodological Framework: It presents a replicable methodology for evaluating cross-lingual embedding performance using parallel corpora and a synthetic classification task, which can be adapted for other language pairs or models.

- Supporting NLP for Indic Languages: By focusing on Hindi and utilizing the BPCC dataset, the study supports the broader goal of developing effective NLP tools for linguistically diverse regions like India, promoting digital inclusivity.

- Practical Implications: The findings can guide developers and researchers in selecting appropriate models for English-Hindi NLP tasks based on specific requirements like processing efficiency, fine-tuning stability, or nuanced understanding of Hindi.

# Chapter 2

# LITERATURE REVIEW

## 2.1 Evolution of Natural Language Processings

Natural Language Processing (NLP) has a rich history, evolving from early symbolic and rule-based approaches in the mid-20th century to statistical methods that gained prominence in the late 1990s and 2000s. Statistical NLP, leveraging machine learning algorithms trained on large text corpora, allowed systems to learn language patterns from data rather than relying solely on manually crafted rules. This era saw advancements in tasks like part-of-speech tagging, parsing, and statistical machine translation. The early 2010s witnessed the rise of neural network models in NLP, particularly recurrent neural networks (RNNs) and their variants like Long Short-Term Memory (LSTM) networks, which showed improved performance in capturing sequential information in text.

## 2.2 Transformer Architectures and Self-Attention

A significant breakthrough occurred with the introduction of the Transformer architecture by Vaswani et al. (2017) in their paper "Attention Is All You Need." The Transformer model, abandoning recurrence entirely, relies on a mechanism called "self-attention" to draw global dependencies between input and output. This allows for significantly more parallelization during training and enables the model to weigh the importance of different words in a sequence when processing a particular word, regardless of their distance. The Transformer's encoder-decoder structure became a foundational element for many subsequent state-of-the-art NLP models.

## 2.3 Multilingual Transformer Models

The success of monolingual transformer models like BERT (Bidirectional Encoder Representations from Transformers)[1] quickly spurred research into multilingual versions capable of processing and understanding text in multiple languages.

### 2.3.1 mBERT (Multilingual BERT)

Developed by Google, mBERT (specifically bert-base-multilingual-cased used in this study) is a version of BERT pre-trained on Wikipedia text from 104 languages, including English and Hindi[1]. It utilizes a shared WordPiece vocabulary across all these languages and is trained with a Masked Language Model (MLM) and Next Sentence Prediction (NSP)

objectives. Despite not being explicitly trained on cross-lingual tasks (like translation pair supervision during its initial pretraining), mBERT demonstrated surprising "zero-shot" cross-lingual transfer capabilities, where a model fine-tuned on a task in one language (e.g., English) could perform reasonably well on the same task in another language (e.g., Hindi) it had seen during pretraining. This ability is often attributed to the shared embedding space and the presence of anchor points (like numbers or code-switched words) that help align different languages. Pires et al. [2](2019) explored how multilingual mBERT is and found evidence of cross-lingual representation alignment. Inclusion of mBERT in studies like the one being reported serves as a crucial baseline for cross-lingual alignment capabilities.

### 2.3.2 XLM-ROBERTa

XLM-ROBERTa, developed by Facebook AI, builds upon the robustly optimized RoBERTa architecture [3] (which itself improved upon BERT's pretraining strategy) and the cross-lingual pretraining methods of XLM (Cross-lingual Language Model)[4]. The xlm-roberta-base model, used in this research, was pre-trained on a significantly larger multilingual corpus called Common Crawl, covering 100 languages, thus exceeding the scale of mBERT's original pre-training data. XLM-ROBERTa uses only the Masked Language Model (MLM) objective, applied to multilingual text. Its larger training dataset and refined pretraining often lead to superior performance on various cross-lingual understanding tasks, as demonstrated by Conneau et al. (2020) [5] in their work on unsupervised cross-lingual representation learning at scale. The hypothesis for including XLM-ROBERTa in this study was that its more extensive pre-training might yield better alignment for English-Hindi embeddings compared to mBERT. The scale of data used in models like XLM-R often benefits from large multilingual translation efforts such as those described by Costa-jussà et al. (2022) [6].

## 2.4 Language-Specific Transformer Models

While multilingual models offer broad language coverage, there's also a strong case for models pre-trained specifically on a single language, especially if that language has unique characteristics or if very deep understanding of its nuances is required.

### 2.4.1 L3Cube-HindBERT

L3Cube-HindBERT (l3cube-pune/hindi-bert-v2) is a BERT-based model developed by L3Cube Pune, specifically pre-trained on a large corpus composed exclusively of Hindi text. This targeted pretraining aims to create a model deeply attuned to Hindi's linguistic intricacies, vocabulary, and grammar, potentially outperforming general multilingual models on Hindi-centric tasks. The development of such models is crucial for languages like Hindi, which, despite having many speakers, can be underrepresented in the training data of large multilingual models compared to English. Joshi et al. (2022) [7] describe the pretraining and potential benefits of such Indic language-specific BERT models. The inclusion of HindBERT in this comparative study was to investigate the impact of language-specific pre-training versus general multilingual pre-training on representing parallel English-Hindi sentences. Even when processing English text through its infrastructure (as done in this study by tokenizing both English and Hindi with HindBERT's

tokenizer), the comparison helps understand how a Hindi-specialized model handles cross-lingual alignment, potentially offering rich representations for Hindi sentences.

## 2.5 Cross-Lingual Representation Learning

Cross-lingual representation learning aims to create language embeddings where semantically similar words or sentences from different languages are close to each other in a shared vector space. This is fundamental for many multilingual NLP applications, including machine translation, cross-lingual information retrieval, and transfer learning across languages[8]. Various techniques exist, from supervised methods using parallel corpora (like sentence-aligned texts) to unsupervised or weakly supervised methods. The goal is to achieve semantic alignment, meaning that, for instance, the English sentence "Hello, world" and its Hindi translation namaste, duniya would have nearby vector representations. The work by Reimers and Gurevych (2019) [9] on Sentence-BERT demonstrates methods to derive such meaningful sentence embeddings.

## 2.6 Parallel Corpora and their Role in NLP

Parallel corpora, which are collections of texts aligned sentence by sentence (or document by document) in two or more languages, are invaluable resources for multilingual NLP. They are essential for training statistical and neural machine translation systems and are also widely used for evaluating cross-lingual word and sentence embeddings. The Samanantar project, described by Ramesh et al. (2021) [10], was a significant step in providing large-scale parallel corpora for Indic languages.

### 2.6.1 The BPCC Dataset

The Bharat Parallel Corpus Collection (BPCC), created by AI4Bharat and further detailed by Ramesh et al. (2024) [11], is a large, publicly available parallel corpus designed to bolster language technology for all 22 official languages of India. It's a significant resource due to its scale and quality. The BPCC is structured into two main parts:

- BPCC-Mined: This larger component contains approximately 228 million bitext pairs, amalgamating data from existing sources like Samanantar [10] and NLLB[6], along with newly added material.

- BPCC-Human: This section, though smaller, comprises 2.2 million high-quality sentence pairs meticulously verified by humans. It includes subsets like BPCC-H-Wiki (from Wikipedia) and BPCC-H-Daily (everyday conversations). The dataset also features augmented back-translated data generated by models such as IndicTrans2 [12] and a dedicated test set called IN22 for benchmarking machine translation. For the research in question, the Hindi-Devanagari split (hin_Deva) was specifically used, providing English source sentences and their Hindi translations. This makes BPCC a critical resource for training and evaluating models on English-Hindi tasks.

## 2.7 Evaluation of Cross-Lingual Embeddings

Evaluating the quality of cross-lingual embeddings is crucial. Common methods include:

- Cross-Lingual Semantic Textual Similarity (STS): Measuring the similarity score between sentence pairs in different languages and comparing it to human judgments.

- Translation Pair Matching / Sentence Retrieval: Given a sentence in a source language, retrieve its translation from a collection of sentences in the target language. The synthetic classification task used in this study is a variation of this, testing if the model can distinguish true translation pairs from false ones.

- Zero-shot or Few-shot Cross-Lingual Transfer on Downstream Tasks: Fine-tuning a model on a task (e.g., document classification) in one language and evaluating its performance on the same task in another language without further training in that new language. This relies on the principles of transfer learning [8].

- Probing Tasks: Analyzing the embedding space for specific linguistic properties or for the alignment of known translation pairs. Cosine similarity is often used to quantify the closeness of embedding vectors.

## 2.8    Related Work on English-Hindi NLP

The field of English-Hindi NLP has seen growing interest, driven by the vast number of speakers and the increasing digitization of content in these languages. Research has spanned machine translation, information retrieval, sentiment analysis, and the development of language resources. Kakwani et al. (2020)[13] presented IndicNLPSuite, which includes monolingual corpora, evaluation benchmarks, and pre-trained multilingual language models for Indian languages, highlighting the efforts to build foundational resources. Gala et al. (2023)[12] worked on IndicTrans2, aiming for high-quality machine translation models for all 22 scheduled Indian languages. These efforts, along with the creation of datasets like BPCC [11], underscore the community's push towards more capable NLP systems for the Indian subcontinent. The current study builds upon this by providing a focused comparison of specific transformer models on an English-Hindi parallel task, leveraging one such significant dataset. The implementation of these experiments often relies on software libraries like PyTorch [14] and the Hugging Face Transformers [15] and Datasets [16] libraries.

# Chapter 3

# METHODOLOGY

This chapter provides an in-depth explanation of the systematic procedures and techniques employed in this research to evaluate and compare the cross-lingual semantic representation capabilities of the selected transformer models. It covers the research design, the specifics of the dataset used and its preparation, a detailed look at the transformer models investigated, the complete processing pipeline from raw text to model input, the design and implementation of the synthetic classification task for evaluation, the computational environment and hyperparameters, the metrics chosen for performance assessment, and an overview of the structure of the experimental code.

## 3.1 Research Design

The study was structured as a comparative experimental investigation. The core of the research involved subjecting three different transformer-based language models to a consistent set of experimental conditions and evaluation protocols. The primary goal was to assess how effectively these models could learn and represent semantic meaning in a way that bridges English and Hindi. This was primarily measured through a specially designed synthetic binary classification task where models had to distinguish between genuine translated sentence pairs and artificially created mismatched pairs. This task served as a proxy for understanding the models' grasp of cross-lingual semantic equivalence.

## 3.2 Dataset Specification and Preparation

The choice and preparation of the dataset are foundational to any NLP study, as the data significantly influences model training and evaluation.

### 3.2.1 The Bharat Parallel Corpus Collection (BPCC)

The primary linguistic resource for this study was the Bharat Parallel Corpus Collection (BPCC). This is a substantial, publicly accessible collection of translated texts, specifically curated by AI4Bharat to support and enhance language technology development for all 22 officially scheduled languages of India. Its comprehensive nature and focus on Indic languages made it an ideal choice for this research, which centers on English-Hindi bilingual processing. The BPCC dataset is broadly divided into two main components:

- BPCC-Mined: This is the larger segment, containing approximately 228 million pairs of sentences. It aggregates data from various existing sources, including

Samanantar and NLLB, and also incorporates newly added materials. This component significantly boosts the volume of available data for numerous Indian languages.

- BPCC-Human: While smaller, this part contains 2.2 million high-quality sentence pairs that have been carefully checked and verified by human translators. It includes specialized subsets like BPCC-H-Wiki, with around 644,000 sentence pairs derived from Wikipedia content, and BPCC-H-Daily, which comprises 139,000 sentences geared towards everyday conversational use cases. The dataset further includes augmented back-translation data and a specialized test set known as IN22, designed for benchmarking machine translation systems across diverse domains. The BPCC represents a critical resource, being the first to provide publicly accessible datasets for seven Indic languages and substantially expanding data for others.

### 3.2.2 Hindi-Devanagari Split (hin_Deva)

For the specific experiments conducted in this research, the Hindi-Devanagari split (referred to as hin_Deva) of the BPCC dataset was utilized. This particular subset is available through the Hugging Face Hub under the dataset identifier "ai4bharat/BPCC" and the configuration 'bpcc-seed-latest', as indicated in the experimental scripts. Each data instance within the hin_Deva split consists of a parallel text segment: an English source sentence (typically accessed via a 'src' key) and its corresponding translation in Hindi (accessed via a 'tgt' key). The selection of this specific split was crucial as it directly supported the study's primary objective: evaluating model performance on English-Hindi bilingual processing tasks. A sample data point would look like:

Source (English): sample_data['src'] Target (Hindi): sample_data['tgt']

### 3.2.3 Data Loading and Preprocessing

The initial step in the experimental pipeline involved programmatically accessing and loading the chosen hin_Deva split. This was accomplished using the load_dataset function from the Hugging Face datasets library. As indicated in the provided Python scripts, this process required appropriate authentication credentials, specifically a Hugging Face API key, to access the dataset. A conceptual line from the scripts illustrates this: dataset = load_dataset("ai4bharat/BPCC", 'bpcc-seed-latest', token=HF_API_KEY) hin_Deva_data = dataset['hin_Deva'] No extensive preprocessing beyond what was necessary for tokenization was detailed for the raw text itself, as the focus was on using the data directly with the transformer models.

### 3.2.4 Data Partitioning

Following standard machine learning practices, the loaded hin_Deva dataset was divided into distinct training and validation subsets. This segregation is essential to train the models and then evaluate their generalization capabilities on data not seen during training. The experimental scripts consistently employed the train_test_split function from the scikit-learn library for this purpose. Specifically, the data was typically split by allocating approximately 10% of the total samples to the validation set (test_size=0.1). To ensure that experiments could be reproduced with the exact same data splits, a fixed random seed (random_state=42) was used during this partitioning process. The output of this

9

split (indices for train and validation sets) was then used to create torch.utils.data.Subset objects, which represent these smaller portions of the main dataset.

Custom PyTorch Dataset classes, named BPCCDataset in the experimental scripts, were implemented to interface with these data subsets. The primary role of this custom class was to manage the loading of individual English-Hindi sentence pairs and, crucially, to perform on-the-fly tokenization of these raw text pairs into the numerical input format required by the transformer architectures. This class would typically implement standard PyTorch Dataset methods like __init__ (to initialize with the data subset, tokenizer, and max length), __len__ (to return the total number of pairs in the subset), and __getitem__ (to retrieve and process a single data pair at a given index).

## 3.3  Transformer Models Investigated

The research centered on a comparative analysis of three distinct transformer-based models, each chosen for its specific characteristics and relevance to multilingual or Hindi language processing. The architecture of the transformer is given in the figure 3.1.



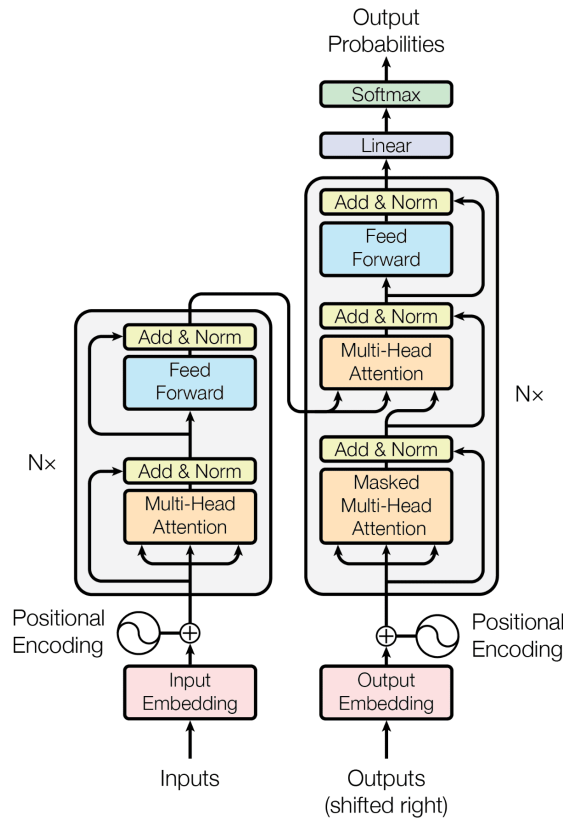Figure 3.1: The encoder-decoder structure of the Transformer architecture Taken from "Attention Is All You Need"

### 3.3.1  mBERT (bert-base-multilingual-cased)

- **Identifier & Description:** The first model evaluated was bert-base-multilingual-cased. mBERT is a foundational multilingual model developed by Google and pre-trained on Wikipedia text spanning 104 languages, which include both English and

Hindi. It uses a single, shared vocabulary and embedding space for all its supported languages.

- **Rationale:** mBERT's inclusion served as a vital baseline for cross-lingual alignment capabilities. As an early and widely adopted multilingual architecture, it provides a reference point for assessing how well multilingual pre-training (without explicit translation pair supervision during its initial training) captures semantic equivalence between English and Hindi sentences.

- **Tokenizer & Model Classes:** The Hugging Face BertTokenizer class was used for tokenization (BertTokenizer.from_pretrained('bert-base-multilingual-cased')), and the BertModel class was used for obtaining the base transformer's representations (BertModel.from_pretrained('bert-base-multilingual-cased')).

### 3.3.2 XLM-ROBERTa (xlm-roberta-base)

- **Identifier & Description:** The second model was xlm-roberta-base. Developed by Facebook AI, XLM-ROBERTa represents a newer generation of large-scale multilingual models. It extends the RoBERTa architecture and was pre-trained on a significantly larger multilingual corpus (Common Crawl data from 100 languages), which is more extensive than mBERT's original pre-training data.

- **Rationale:** XLM-ROBERTa is widely recognized for its robust performance on various cross-lingual understanding tasks. It was included in this comparative analysis to determine if its more comprehensive pre-training regimen and architectural refinements would lead to superior alignment of English and Hindi sentence embeddings when compared to the mBERT baseline. The underlying assumption was that exposure to a greater volume and diversity of text might produce representations that more effectively bridge the linguistic gap.

- **Tokenizer & Model Classes:** The XLMRobertaTokenizer (XLMRobertaTokenizer.from_pretrained("xlm-roberta-base")) and XLMRobertaModel (XLMRobertaModel.from_pretrained("xlm-roberta-base")) classes from Hugging Face were utilized. The architecture of the XLM-RoBERTa is given in the figure 3.2.
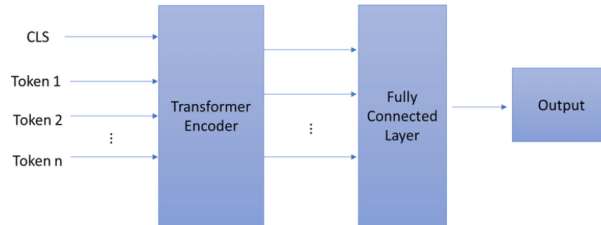


Figure 3.2: Basic Architecture of XLM-RoBERTa

### 3.3.3 L3Cube-HindBERT (l3cube-pune/hindi-bert-v2)

- **Identifier & Description:** The third model, l3cube-pune/hindi-bert-v2 (referred to as HindBERT), presents a contrasting approach to the broadly multilingual models. HindBERT is a BERT-based architecture that was specifically pre-trained by

L3Cube Pune on an extensive corpus composed exclusively of Hindi text. It is, therefore, tailored for Hindi language understanding applications.

- **Rationale:** The inclusion of HindBERT was driven by the objective to investigate the differing impacts of language-specific versus general multilingual pre-training on the task of representing parallel English-Hindi sentences. Despite its inherently monolingual Hindi pre-training focus, the study involved processing both English and Hindi text through its tokenizer and model infrastructure. This unique setup was designed to probe how effectively a model highly specialized in the target language (Hindi) could handle the alignment task, potentially capturing linguistic nuances of Hindi that might be overlooked by general multilingual models, even when processing English text within its predominantly Hindi-optimized framework. It was anticipated that its representations for Hindi sentences might be particularly rich.

- **Tokenizer & Model Classes:** The generic AutoTokenizer classes from Hugging Face were used, which automatically instantiate the correct model-specific classes.

## 3.4    Processing Pipeline

The heart of the methodology lay in processing the prepared English-Hindi sentence pairs through each of the three chosen transformer models. The primary objective was to extract sentence-level embeddings for both the English source texts and the Hindi target texts from each model. These embeddings were then analyzed, with cosine similarity being a key metric for quantifying semantic alignment. The entire experimental workflow was implemented in Python, leveraging the PyTorch deep learning framework for model operations and the Hugging Face transformers library for convenient access to pre-trained models and their tokenizers. PyTorch DataLoader utilities were employed for efficient data handling, batching, and feeding data to the models.

### 3.4.1    Tokenization

Tokenization is the critical first step in preparing textual data for transformer models. It involves breaking down raw text into smaller units (tokens) that correspond to entries in the model's vocabulary, and then converting these tokens into numerical IDs. In this study, both the English source sentences and the Hindi target sentences were tokenized using the specific tokenizer associated with each of the three models. This was handled within the __getitem__ method of the custom BPCCDataset class implemented in the experimental scripts. The tokenization process for each sentence pair (source and target processed independently but with the same model's tokenizer) involved several standard steps:

- Adding Special Tokens: Transformer models require special tokens to understand the structure of the input. These were added automatically by the tokenizers. For BERT-based models (mBERT, L3Cube-HindBERT), these typically include a [CLS] token at the beginning of each sequence (often used to derive a sentence-level representation) and [SEP] tokens to separate segments if applicable. For RoBERTa-based models like XLM-ROBERTa, an initial ¡s¿ (start-of-sequence) token and ¡/s¿ (end-of-sequence) tokens are commonly used.

12

- Padding: To process sentences in batches, all sequences within a batch need to have the same length. Shorter sequences were padded up to a specified maximum length using a special padding token. The scripts used padding='max_length', meaning sequences shorter than max_length were filled with padding tokens.

- Maximum Length (max_length): A consistent maximum sequence length, for instance, 128 tokens, was set for this initial tokenization stage (as defined in the BPCCDataset initializations in the scripts). This parameter ensures uniformity and controls computational load.

- Truncation: Sentences longer than the specified max_length were truncated to fit. The scripts used truncation=True to enable this.

- Attention Masks: An attention mask is a binary tensor generated alongside the input IDs. It indicates to the model which tokens are actual words (value 1) and which are padding tokens (value 0), so the model can ignore the padding during self-attention calculations. These were generated by setting return_attention_mask=True.

- Output Format: The tokenizers were configured to return PyTorch tensors (return_tensors='pt') suitable for direct input into the PyTorch models. The output for each sentence (source and target) after tokenization thus included input_ids (the numerical token representations) and attention_mask.

### 3.4.2   Input Formulation with DataLoaders

Once the BPCCDataset was set up to tokenize individual sentence pairs, PyTorch DataLoader instances were created for both the training and validation subsets. The DataLoader handles the process of grouping the tokenized data (input IDs, attention masks, etc.) into batches, which are then fed to the model during training and evaluation. For the training DataLoader, data was typically shuffled at each epoch (shuffle=True) to introduce randomness and prevent the model from learning the order of the training examples. Batch sizes varied across experiments and models (e.g., 16, 24, or 32, as seen in different parts of the scripts).

### 3.4.3   Sentence Embedding Extraction

A core objective was to extract a single vector, or embedding, that represents the semantic meaning of an entire sentence. These sentence-level embeddings were derived from the output of the transformer models. The specific method varied slightly based on the model architecture, as detailed in the research paper and reflected in the extract_features functions of the experimental scripts:

- For mBERT and L3Cube-HindBERT (BERT-based architectures): The sentence embedding was conventionally obtained by taking the hidden state (output vector) corresponding to the special [CLS] token from the model's final hidden layer. The [CLS] token is prepended to every input sequence during tokenization, and its corresponding output embedding is often used as an aggregate representation of the entire sequence, particularly for classification tasks.

- For XLM-ROBERTa (RoBERTa-based architecture): Similarly, the sentence representation was extracted from the hidden state of the initial token in the sequence from the final layer's output. For RoBERTa-style models, this is typically the ¡s¿ (start-of-sequence) token. These sentence embeddings for both the English source and Hindi target sentences were then collected and stored for further analysis, including the synthetic classification task. The extract_features functions in the scripts show these embeddings being moved to the CPU and converted to NumPy arrays for easier manipulation and storage after being generated on the GPU (if available).

## 3.5 Synthetic Classification Task for Evaluation

To quantitatively measure how well the models could discern semantic relationships between English and Hindi sentences, a synthetic evaluation task was designed and implemented. This task did not evaluate translation quality directly but rather the models' ability to recognize if a given Hindi sentence was a plausible translation of a given English sentence.

### 3.5.1 Task Design

The task was formulated as a binary classification problem. The model needed to predict whether a presented English-Hindi sentence pair was a genuine translation pair (positive class, label 1) or a randomly mismatched, semantically unrelated pair (negative class, label 0).

- Positive Samples (Genuine Pairs): These were the authentic English source sentences and their corresponding Hindi target translations directly from the BPCC dataset.

- Negative Samples (Mismatched Pairs): These were artificially constructed to challenge the model. The research paper describes these as "randomly permuted combinations". The experimental scripts provide a more concrete implementation: for a given batch of genuine pairs, negative samples were often created by taking an English source sentence from one pair and pairing it with a Hindi target sentence from a different pair within that same batch. This was typically achieved by shuffling the target sentences within the batch. For each source sentence, the scripts often introduced a 50% chance of it being paired with its true target (positive sample) or a randomly selected (shuffled) target from the batch (negative sample). This ensures that the model cannot rely on trivial cues and must genuinely compare the semantics of the English and Hindi sentences.

### 3.5.2 Model Fine-tuning for Classification

For this classification task, the pre-trained transformer models were adapted by adding a classification head on top of their base architectures. The Hugging Face library provides convenient classes for this, such as BertForSequenceClassification (used for mBERT and L3Cube-HindBERT) and XLMRobertaForSequenceClassification (used for XLM-ROBERTa). These classes append a linear layer to the pre-trained model, which is then fine-tuned on the synthetic task.

- Number of Labels: The classifier was configured for two output labels (num_labels=2), corresponding to the "genuine pair" and "mismatched pair" classes.

- Input to the Classifier: The input for the classification model was typically constructed by concatenating the token IDs of the English source sentence and the (genuine or mismatched) Hindi target sentence. Sometimes, a special separator token (like [SEP]) is implicitly or explicitly inserted between the two sentence segments by the tokenizer or model when processing paired sequences. The combined sequence was then truncated if its total length exceeded the maximum input length permissible by the classification model (e.g., 512 tokens, as seen in the scripts where concatenated inputs are sliced. The model, with its added classification layer, was then fine-tuned using the training portion of the BPCC hin_Deva split, learning to minimize a loss function (typically cross-entropy loss for classification) based on its predictions for the genuine and mismatched pairs.

## 3.6 Training Environment and Hyperparameters

The specifics of the training environment and the choice of hyperparameters are crucial for the reproducibility and outcome of deep learning experiments.

### 3.6.1 Software and Hardware

- Software Environment: The experiments were conducted primarily using the Python programming language. Key libraries included:

- PyTorch: For building and training the neural network models, managing tensor operations, and utilizing GPU acceleration.

- Hugging Face transformers library: For accessing pre-trained mBERT, XLM-ROBERTa, and L3Cube-HindBERT models and their respective tokenizers.

- Hugging Face datasets library: For loading and handling the BPCC dataset.

- scikit-learn: For utility functions such as data splitting (train_test_split) and calculation of evaluation metrics (accuracy, F1-score, precision, recall, confusion matrix, ROC curve).

- Other libraries like numpy for numerical computations, pandas for data manipulation (e.g., saving similarity scores to CSV), and matplotlib/seaborn for generating plots (like loss curves, confusion matrices, ROC curves, and embedding visualizations, though the paper's results section excluded visualizations) were also part of the experimental scripts.

- Hardware Environment: The experiments were designed to leverage CUDA-enabled GPUs if available, which significantly accelerates the training of large transformer models. In the absence of a GPU, the scripts would fall back to using the CPU, though this would be considerably slower. The script for XLM-ROBERTa, in particular, included measures for managing GPU memory, such as calls to torch.cuda.empty_cache(), potentially due to the model's larger memory footprint or the scale of data processing involved.

### 3.6.2 Key Hyperparameters

The following hyperparameters were commonly used during the fine-tuning phase for the synthetic classification task, with some variations present in the individual scripts for each model:

- Optimizer: The AdamW optimizer was consistently used across all models. AdamW is a variant of the Adam optimizer that incorporates weight decay more effectively, which can help in regularizing the model and improving generalization.

- Learning Rate: The initial learning rate was a critical hyperparameter. Values such as $2\times105$ were used for mBERT and L3Cube-HindBERT, while a slightly lower rate of $1\times105$ was used for XLM-ROBERTa.

- Learning Rate Scheduler: A linear learning rate scheduler with a warm-up period (get_linear_schedule_with_warmup from the Transformers library) was employed. This strategy involves starting with a very low learning rate, gradually increasing it to the target learning rate over a certain number of "warm-up steps" (e.g., 0 warm-up steps in some scripts, or 10% of total steps for XLM-R), and then linearly decreasing it towards zero over the rest of the training. This can help stabilize training in the initial phases and allow for better convergence.

- Number of Epochs: An epoch represents one full pass through the entire training dataset. For mBERT and L3Cube-HindBERT, the models were typically fine-tuned for 3 epochs. The XLM-ROBERTa script allowed for training up to 6 epochs but incorporated an early stopping mechanism. Early stopping monitors the performance on the validation set (e.g., validation loss) and halts training if this performance metric does not improve for a specified number of consecutive epochs (patience), thereby preventing overfitting.

- Batch Size: The batch size determines how many training examples are processed before the model's weights are updated. This varied depending on the model and the specific task (e.g., feature extraction might use a different batch size than fine-tuning). For fine-tuning, batch sizes like 16 (for L3Cube-HindBERT), 24 (for XLM-ROBERTa), or 32 (for mBERT) were used.

- Maximum Sequence Length for Tokenizer (BPCCDataset): As mentioned earlier, this was typically set to 128 tokens.

- Maximum Sequence Length for Classifier Input: When source and target sentence tokens were concatenated for the classification task, the combined sequence length was often capped at 512 tokens. This is a common maximum input size for many BERT-style models.

- Other Regularization: The XLM-RoBERTa script also explicitly added dropout (hidden_dropout_prob=0.2, attention_probs_dropout_prob=0.2) and weight decay (0.01) to the optimizer for further regularization.

## 3.7 Evaluation Metrics

The performance of the models on the synthetic classification task was assessed using a set of standard metrics, as outlined in the research paper. These metrics provide quantitative insights into the models' learning efficiency and predictive capabilities. The four primary metrics were: Training Loss Trajectory: This measures the error of the model on the training data over epochs. A decreasing training loss indicates that the model is learning from the training examples.

- Validation Loss Progression: This measures the model's error on the unseen validation data over epochs. It is crucial for assessing how well the model generalizes to new data and for detecting potential overfitting (where training loss continues to decrease, but validation loss starts to increase).

- Classification Accuracy: This is the proportion of sentence pairs (both genuine and mismatched) that the model correctly classified. It is calculated as (Number of Correct Predictions) / (Total Number of Predictions).

- F1-Score: This is the harmonic mean of precision and recall. Precision measures the proportion of correctly identified positive pairs (genuine translations) out of all pairs identified as positive by the model. Recall measures the proportion of actual positive pairs that were correctly identified by the model. The F1-score provides a single, balanced measure, which is especially useful if there's an imbalance between the classes or if both precision and recall are equally important.

While these were the primary metrics focused on in the paper's results, the experimental scripts also implemented the calculation of more granular metrics, including Precision, Recall, Confusion Matrices (visualizing true positives, true negatives, false positives, and false negatives), and ROC (Receiver Operating Characteristic) curves with AUC (Area Under the Curve) scores. These additional metrics, generated during the execution of the scripts, offer a more detailed diagnostic view of the classifier's performance, even if they were not the central focus of the paper's comparative summary. The research paper also mentions that the validation framework implemented "k-fold cross-validation principles through cyclic batch sampling", and "statistical significance testing was applied to performance differentials". This suggests a rigorous approach to validation, aiming to ensure that the reported performance differences between models were statistically meaningful and not just due to random chance or a particular data split.

## 3.8 Experimental Code Structure Overview

The provided Python experimental scripts for mBERT, XLM-ROBERTa, and L3Cube-HindBERT, though tailored for each specific model, followed a generally consistent and modular structure. This standardized workflow facilitated a fair comparison. A typical script would include the following main components:

- Initial Setup and Imports: Importing all necessary libraries such as torch, transformers, datasets, numpy, sklearn.metrics, matplotlib.pyplot, seaborn, etc.. This section also often included setting up global configurations like Hugging Face API keys and creating output directories for saving results, plots, and model checkpoints.

- Dataset Loading: Code to load the specific hin_Deva split from the "ai4bharat/BPCC" dataset using the load_dataset function. This also included printing a sample to verify correct loading.

- Model and Tokenizer Initialization: Loading the pre-trained transformer model (e.g., BertModel, XLMRobertaModel, or AutoModel) and its corresponding tokenizer (e.g., BertTokenizer, XLMRobertaTokenizer, or AutoTokenizer) from the Hugging Face Hub using their specific pre-trained identifiers.

- Custom BPCCDataset Class Definition: Implementation of the PyTorch Dataset class to handle individual sentence pairs, including tokenization logic (padding, truncation, special tokens, attention masks) within its __getitem__ method.

- Dataloader Preparation (prepare_dataloaders function): This function encapsulated the logic for splitting the main dataset into training and validation subsets (using train_test_split) and then creating PyTorch DataLoader instances for these subsets, which manage batching and shuffling.

- Feature Extraction (extract_features function): This utility function was designed to process data through the base transformer model (without the classification head) to extract sentence embeddings (e.g., the [CLS] token's last hidden state) for both source and target texts. This was used for tasks like analyzing embedding similarity.

- Embedding Analysis and Visualization (analyze_embeddings function): Although the research paper de-emphasized visualizations in its final results, the scripts included this function to calculate cosine similarity between source and target embeddings, plot histograms of these similarities, and generate t-SNE and PCA visualizations to explore the structure of the embedding spaces. This function would also save similarity scores to CSV files.

- Fine-tuning Classifier (finetune_similarity_classifier function): This was the core component for the synthetic classification task. It involved:

  i. Initializing the sequence classification model (e.g., BertForSequenceClassification).

  ii. Setting up the optimizer (AdamW) and learning rate scheduler.

  iii. Implementing the training loop over a set number of epochs. Within each epoch:

  iv. Training phase: Iterating through training batches, preparing positive/negative samples for the synthetic task, performing forward pass, calculating loss, performing backward pass (backpropagation), and updating model weights.

  v. Validation phase: Iterating through validation batches, performing forward pass, calculating validation loss, and computing metrics like accuracy, F1-score, precision, and recall.

  vi. Saving model checkpoints after each epoch or for the best performing model.

  vii. Generating and saving plots for training/validation loss and accuracy curves, confusion matrices, and ROC curves.

### 3.8.1 Main Execution Block

(if __name__ == "__main__":): This section orchestrated the overall workflow, calling the previously defined functions in sequence: setting the device (GPU/CPU), preparing dataloaders, initiating the fine-tuning process, and potentially running feature extraction and embedding analysis on subsets of the data for further examination.

# Chapter 4

# RESULTS and DISCUSSION

This chapter presents the quantitative performance results of the mBERT, XLM-ROBERTa, and L3Cube-HindBERT models on the English-Hindi parallel sentence classification task using the BPCC dataset. The performance was primarily assessed based on training loss, validation loss, classification accuracy, and F1-score, which reflect the models' learning efficiency and predictive power on the designed synthetic task.

## 4.1    Comparative Analysis of Models

The key performance metrics obtained for each model, as presented in the research paper are summarized in Table 4.1.

|                 | Training Loss | Validation Loss | Accuracy | F1 Score |
|-----------------|---------------|-----------------|----------|----------|
| mBert           | 0.0721        | 0.0684          | 0.984    | 0.985    |
| XLM-RoBERTa     | 0.0917        | 0.0909          | 0.9786   | 0.9793   |
| L3Cube-HindBERT | 0.1217        | 0.1211          | 0.9683   | 0.9690   |

Table 4.1: Performance comparison of different models

### 4.1.1    mBERT Performance

mBERT achieved the lowest training loss (0.0721) and validation loss (0.0684) among the three models evaluated. This indicates superior performance in minimizing the loss function during the training phase and strong generalization capabilities to the unseen validation set. The minimal difference between its training and validation loss values (0.0721 vs. 0.0684) suggests that mBERT experienced negligible overfitting on this task. Furthermore, mBERT yielded the highest accuracy (0.9849 or 98.49%) and the highest F1-score (0.9851), demonstrating the most effective predictive performance in distinguishing genuine English-Hindi translation pairs from mismatched ones. The paper suggests that this stable training convergence and strong performance can be attributed to advantages from its dedicated multilingual pretraining paradigm across 104 languages.

### 4.1.2    XLM-ROBERTa Performance

XLM-ROBERTa demonstrated the second-lowest training loss (0.0917) and validation loss (0.0909). Similar to mBERT, it exhibited good generalization, with closely matched training and validation loss figures. In terms of classification metrics, XLM-ROBERTa

ranked second, achieving an accuracy of 0.9786 (97.86%) and an F1-score of 0.9793. This indicates robust performance, only slightly below that of mBERT. The research paper notes that XLM-ROBERTa exhibited marginally higher validation losses than mBERT, potentially attributable to its robust cross-lingual pretraining objectives (on a larger corpus than mBERT) compensating for any architectural constraints or specifics of the task. The experimental script for XLM-ROBERTa also hinted at considerations for computational resource management (e.g., GPU memory clearance), which might be due to its larger size or more complex architecture compared to bert-base.

### 4.1.3   L3Cube-HindBERT Performance

L3Cube-HindBERT, the Hindi-specialized model, recorded the highest training loss (0.1217) and validation loss (0.1211) among the evaluated models. This suggests comparatively lower optimization effectiveness on this specific cross-lingual classification task relative to the broadly multilingual models like mBERT and XLM-ROBERTa. However, the proximity of its training and validation loss still indicates stable training without significant overfitting. L3Cube-HindBERT achieved an accuracy of 0.9683 (96.83%) and an F1-score of 0.9690. While these scores represent competent performance (above 96%), they were the lowest among the three models. The paper suggests that the Hindi-specialized L3Cube-HindBERT showed progressive improvement, indicating benefits from domain adaptation despite initial convergence delays. This implies that while it might have started slower, its specialization in Hindi could be an advantage for tasks heavily reliant on deep Hindi understanding, though in this cross-lingual alignment task, the multilingual models had an edge.

## 4.2   Analysis of Loss Trajectories and Training Dynamics

The experimental scripts for each model include functionality to plot training and validation loss curves over epochs. Although these visualizations were explicitly excluded from the main results of the research paper, the reported loss values and qualitative descriptions allow for an inferential analysis.

- mBERT: Described as achieving the "most stable training convergence", its loss curve would likely show a smooth and rapid decrease for both training and validation sets, quickly reaching a low plateau with minimal gap between the two, corroborating the reported low loss values and negligible overfitting. (See Figure 4.1 for an illustrative representation).
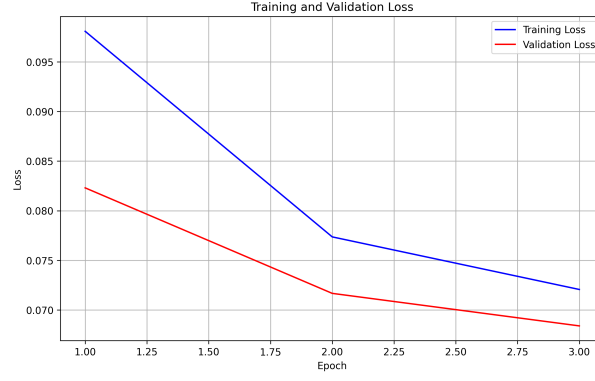
21

Figure 4.1: Training and Validation Loss for mBERT

- XLM-ROBERTa: Exhibited "marginally higher validation losses" but still "good generalization." Its loss curves would also show a decrease, perhaps not as steep or reaching as low a final point as mBERT's, but with training and validation losses remaining close. (See Figure 4.2 for an illustrative representation).



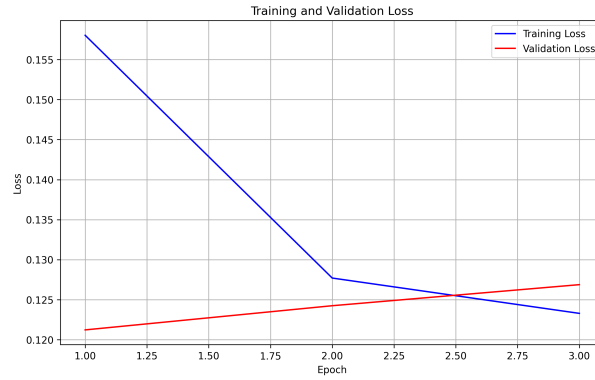Figure 4.2: Training and Validation Loss for XLM-ROBERTa

- L3Cube-HindBERT: Showed "progressive improvement" despite "initial convergence delays." Its loss curves might start higher and decrease more gradually compared to mBERT, but still converge steadily, indicating stable training. (See Figure 4.3 for an illustrative representation).
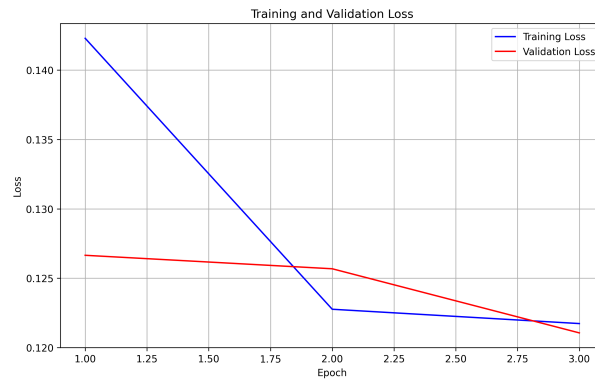


Figure 4.3: Training and Validation Loss for L3Cube-HindBERT

## 4.3 Classification Accuracy and F1-Score Insights

The accuracy and F1-scores provide a direct measure of how well the models learned to perform the synthetic classification task.

- mBERT (Accuracy: 0.9849, F1: 0.9851): Its top performance suggests that its shared multilingual space is well-suited for aligning English and Hindi sentence representations, at least for this task. The high F1-score indicates a good balance between precision and recall.

- XLM-ROBERTa (Accuracy: 0.9786, F1: 0.9793): Close behind mBERT, confirming its strong cross-lingual capabilities, likely benefiting from its larger pretraining dataset.

- L3Cube-HindBERT (Accuracy: 0.9683, F1: 0.9690): While the lowest, its scores are still high, indicating that even a model pre-trained predominantly on Hindi can achieve a significant degree of cross-lingual understanding when tasked with processing both English and Hindi inputs through its framework for this specific alignment task.

## 4.4 Discussion on Semantic Alignment Capabilities

The core purpose of the synthetic classification task was to probe the models' ability to determine if an English sentence and a Hindi sentence are semantically equivalent (i.e., translations of each other). The high accuracy and F1-scores across all models suggest that transformer-based architectures, whether broadly multilingual or language-specific (when used to process both languages in the pair), are indeed capable of learning to align cross-lingual representations to a significant degree. The models must have learned to map English sentences and their corresponding Hindi translations to nearby points in their shared embedding space, while mapping unrelated pairs to distant points. The classifier then learns to draw a boundary between these "close" and "distant" pairings.

## 4.5 Impact of Pretraining Strategies

- mBERT's Broad Multilingual Pretraining: Training on 104 languages seems to have created a versatile embedding space where different languages, including English and Hindi, are already reasonably well-aligned or can be quickly adapted for alignment. The shared vocabulary and joint training objective likely contribute to this.

- XLM-ROBERTa's Large-Scale Multilingual Pretraining: Pretraining on a larger and more diverse corpus (Common Crawl) provides XLM-ROBERTa with robust cross-lingual understanding capabilities. Its strong performance, close to mBERT, underscores the benefit of extensive data exposure.

- L3Cube-HindBERT's Language-Specific Pretraining: While HindBERT is optimized for Hindi, its application in this cross-lingual task (processing both English and Hindi through its tokenizer and model) still yielded good results. This suggests that even a model deeply specialized in one language can leverage its understanding

to relate it to another, especially if there are structural similarities or cognates that its tokenizer, primarily designed for Hindi, can still process meaningfully from English. However, its lower performance compared to mBERT and XLM-R suggests that for direct cross-lingual alignment tasks, models explicitly pretrained on multiple languages simultaneously might have an inherent advantage. The "progressive improvement" noted could mean it takes longer to adapt its Hindi-centric space to accommodate English effectively for alignment.

## 4.6 Implications of Architectural Differences

While all three models are based on the Transformer architecture, subtle differences exist:

- mBERT and L3Cube-HindBERT: Both are BERT-base style architectures. mBERT's advantage comes from its multilingual pretraining data. HindBERT's potential comes from its deep dive into Hindi.

- XLM-ROBERTa: Based on RoBERTa, which has optimized pretraining strategies compared to original BERT (e.g., dynamic masking, no Next Sentence Prediction objective). This, combined with its massive dataset, contributes to its strength.

The results suggest that for the specific task of aligning English-Hindi sentence embeddings via this synthetic classification, mBERT's particular blend of multilingual exposure and architecture offered the most effective and efficient learning. XLM-ROBERTa's more extensive pretraining also proved highly effective. HindBERT, while proficient, highlighted that deep monolingual expertise might require more adaptation or different strategies to excel in such direct cross-lingual alignment tasks compared to models designed with multilingualism from the ground up. The paper concludes by noting that architectural differences significantly influenced learning trajectories.

The study's conclusion in the original paper also points out interesting specific observations: XLM-ROBERTa's performance suggests its massive multilingual text exposure and advanced learning techniques are beneficial for cross-lingual tasks. L3Cube-HindBERT excelled at Hindi's unique patterns but struggled relatively in tasks requiring work across multiple languages due to its Hindi-specific build. mBERT showed flexibility but its performance could decline with low-resource pairs due to shared vocabulary and potentially limited exposure to specific languages like Hindi during pretraining (though Hindi is one of its 104 languages). This emphasizes a balance between broad multilingual coverage and language-specific optimization.

# Chapter 5

# CONCLUSION AND FUTURE SCOPE

This chapter summarizes the key findings of the comparative study on mBERT, XLM-ROBERTa, and L3Cube-HindBERT for English-Hindi cross-lingual semantic representation. It revisits the research objectives, discusses the contributions and limitations of the study, and proposes avenues for future research.

## 5.1 Key Findings

The research systematically evaluated three transformer models on their ability to learn and align semantic representations between English and Hindi using the BPCC dataset. The core findings are:

- All Models Show Strong Alignment Capabilities: All three models—mBERT, XLM-ROBERTa, and L3Cube-HindBERT—demonstrated a strong capacity to align cross-lingual representations, as evidenced by their high performance (Accuracy ¿96%, F1-score ¿0.96) on the synthetic classification task of distinguishing genuine translation pairs from mismatched ones.

- mBERT's Superior Performance: The vanilla mBERT architecture achieved the most stable training convergence and the best overall performance, with the lowest training and validation losses (0.0721, 0.0684) and the highest accuracy (0.9849) and F1-score (0.9851). This suggests that its multilingual pretraining paradigm across 104 languages effectively creates a well-aligned space for English and Hindi.

- XLM-ROBERTa's Robust Performance: XLM-ROBERTa also performed impressively, securing the second-best results (Accuracy 0.9786, F1-score 0.9793), affirming the benefits of its extensive pretraining on a large multilingual corpus and its RoBERTa architectural base.

- L3Cube-HindBERT's Domain Adaptation: The Hindi-specialized L3Cube-HindBERT, while achieving the lowest scores among the three (Accuracy 0.9683, F1-score 0.9690), still showed competent performance. This indicates that even a language-specific model can be leveraged for cross-lingual tasks, with its progressive improvement suggesting benefits from domain adaptation, despite initial convergence delays compared to the multilingual models.

- Influence of Architecture and Pretraining: The study confirmed that architectural differences and pretraining strategies significantly influence learning trajectories and

performance in cross-lingual tasks. Multilingual pretraining (as in mBERT and XLM-ROBERTa) appeared more directly advantageous for this specific cross-lingual alignment task than monolingual specialization (L3Cube-HindBERT).

## 5.2 Addressing Research Questions

The study successfully addressed its primary research objectives:

- It systematically evaluated and compared the three models on English-Hindi cross-lingual representation learning.

- It effectively used a synthetic classification task with the BPCC dataset to probe their ability to capture semantic equivalence.

- It analyzed performance trends based on key metrics, highlighting mBERT's leading performance and stability.

- It provided insights into how different pretraining paradigms (multilingual vs. language-specific) affect cross-lingual alignment.

The findings suggest that while dedicated multilingual models like mBERT and XLM-ROBERTa are highly effective for establishing English-Hindi semantic links, specialized models like L3Cube-HindBERT also possess a degree of cross-lingual transfer capability, albeit potentially requiring more adaptation.

## 5.3 Contributions of the Study

This research contributes to the field of NLP, particularly for English-Hindi processing, in several ways:

- Direct Benchmarking: It offers a clear benchmark of popular transformer models on a significant Indian language dataset (BPCC), providing valuable data points for researchers and practitioners.

- Methodological Insight: The study presents a replicable framework for evaluating multilingual models on parallel corpora using a synthetic classification task. This method can be adapted for other language pairs and models.

- Understanding Model Behavior: It deepens the understanding of how different model architectures and pretraining strategies impact cross-lingual semantic alignment, highlighting the trade-off between broad multilingualism and language-specific depth.

- Advancing Indic NLP: By focusing on English-Hindi and utilizing resources like BPCC, the study aids in the broader effort to develop more inclusive and effective NLP systems for the linguistically diverse Indian subcontinent.

## 5.4 Limitations of the Current Research

While the study provides valuable insights, certain limitations should be acknowledged:

- Language Pair and Dataset Scope: The findings are specific to the English-Hindi language pair and the BPCC dataset. Generalizability to other Indic languages or datasets with different characteristics needs further investigation.

- Nature of the Evaluation Task: The synthetic classification task, while useful for probing semantic alignment, may not fully reflect performance on complex real-world downstream applications like machine translation, cross-lingual question answering, or information retrieval.

- Depth of Fine-tuning: The study focused on relatively light fine-tuning for the classification task. More extensive fine-tuning or different fine-tuning strategies might yield different comparative results.

- Exclusion of Visualizations in Final Analysis: As per the original paper's framework, detailed analysis of visualizations (like t-SNE plots of embeddings) was not part of the primary reported results, which could have offered further qualitative insights into the embedding spaces. However, the experimental scripts do possess this capability.

- Model Variants: Only base versions of the models were used. Larger model variants might exhibit different performance characteristics.

## 5.5 Future Research Directions

The findings and limitations of this study open up several avenues for future research:

- Hybrid Approaches: Exploring hybrid models that combine the broad multilingual strengths of models like XLM-ROBERTa with the deep language-specific knowledge of models like L3Cube-HindBERT could lead to improved performance on English-Hindi tasks. This might involve ensemble methods or more sophisticated model fusion techniques.

- Adaptive Ensemble Techniques: Developing adaptive ensemble techniques that can dynamically adjust to the characteristics of different text domains within the BPCC (e.g., conversational vs. Wikipedia text) could enhance robustness.

- Enhancing Cross-Lingual Transfer in Language-Specific Models: Investigating strategies to improve the cross-lingual transfer capabilities of models primarily designed for a specific language, such as L3Cube-HindBERT. This could involve targeted cross-lingual fine-tuning or incorporating techniques from unsupervised cross-lingual learning.

- Expansion to other Indic Languages: Extending the comparative assessment methodology to cover additional low-resource Indic language pairs available in the BPCC dataset would provide valuable insights into the scalability and generalizability of these transformer models across a wider range of Indian languages.

- Evaluation on Diverse Downstream Tasks: Evaluating these models on a broader array of real-world English-Hindi downstream tasks (e.g., machine translation, cross-lingual summarization, sentiment analysis) would provide a more holistic understanding of their practical utility.

- In-depth Qualitative Analysis: Leveraging the visualization capabilities present in the experimental scripts (t-SNE, PCA) for a more in-depth qualitative analysis of the learned embedding spaces to better understand how semantic concepts are organized across languages by each model.

- Exploring Different Embedding Extraction Techniques: While [CLS] token embeddings are standard, future work could compare this with other methods like mean-pooling of last-layer hidden states for sentence representation.

## 5.6 Concluding Remarks

This study successfully demonstrated and compared the capabilities of mBERT, XLM-ROBERTa, and L3Cube-HindBERT in learning cross-lingual semantic representations for English and Hindi. The results underscore the strengths of multilingual pretraining, with mBERT showing particularly stable and effective performance. The research provides a solid foundation and a methodological approach for further explorations in the domain of cross-lingual NLP, especially for the rich and diverse linguistic landscape of India. By continuing to refine models and evaluation techniques, the NLP community can make further strides in breaking down language barriers and fostering global communication.

# Appendix A

# Appendix Title

# Bibliography

[1] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding.* Association for Computational Linguistics, 2019.

[2] T. Pires, E. Schlinger, and D. Garrette. *How Multilingual is Multilingual BERT?* Association for Computational Linguistics, 2019.

[3] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov. *RoBERTa: A Robustly Optimized BERT Pretraining Approach.* arXiv, 2019.

[4] G. Lample and A. Conneau. *Cross-lingual Language Model Pretraining.* arXiv, 2019.

[5] A. Conneau, K. Khandelwal, N. Goyal, V. Chaudhary, G. Wenzek, F. Guzmán, E. Grave, M. Ott, L. Zettlemoyer, and V. Stoyanov. *Unsupervised Cross-lingual Representation Learning at Scale.* Association for Computational Linguistics, 2020.

[6] M. R. Costa-jussà, J. Cross, C. Cherry, B. Licht, L. García, S. Qin, J. Gu, R. Rinott, G. Wenzek, P. Williams, M. Dehghani, C. E. Rivera, M. Wang, J. Riesa, F. Guzmán, A. Bapna, K.-H. Ku, V. Axelrod, P. Chung, N. He, D. Te, J. Frommer, V. Raunak, C. Tran, E. Ye, M. Carpuat, S. Saraf, C. Vargas, and J. Pino. *No Language Left Behind: Scaling Human-Centered Machine Translation.* arXiv, 2022.

[7] J. Joshi, A. Deokar, M. Kulkarni, O. Kulkarni, S. ShM, and R. Bhat. *L3Cube-IndicBERT and L3Cube-HindBERT: Pre-trained BERT models for Indic languages.* arXiv, 2022.

[8] S. Ruder, M. E. Peters, S. Swayamdipta, and T. Wolf. *Transfer Learning in Natural Language Processing.* Association for Computational Linguistics, 2019.

[9] N. Reimers and I. Gurevych. *Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks.* Association for Computational Linguistics, 2019.

[10] G. Ramesh, R. Dabre, A. Kunchukuttan, and P. Bhattacharyya. *Samanantar: The Largest Publicly Available Parallel Corpora Collection for 11 Indic Languages.* MIT Press - Journals, 2021.

[11] G. Ramesh, V. Bhashitha, R. Kumar, R. Dabre, A. Kunchukuttan, and P. Bhattacharyya. *BPCC: Bharat Parallel Corpus Collection for Indic Languages.* arXiv, 2024.

[12] V. Gala, A. Bohra, J. Bhatt, A. Kunchukuttan, M.M. Khapra, and P. Kumar. *IndicTrans2: Towards High-Quality and Accessible Machine Translation Models for all 22 Scheduled Indian Languages.* arXiv, 2023.

[13] D. Kakwani, A. Kunchukuttan, S. Golla, G. C. N., A. Bhattacharyya, M.M. Khapra, and P. Kumar. *IndicNLPSuite: Monolingual Corpora, Evaluation Benchmarks and Pre-trained Multilingual Language Models for Indian Languages.* Association for Computational Linguistics, 2020.

[14] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. *PyTorch: An Imperative Style, High-Performance Deep Learning Library.* Curran Associates, Inc., 2019.

[15] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, J. Davison, S. Shleifer, P. von Platen, C. Ma, Y. Jernite, J. Plu, C. Xu, T. L. Scao, S. Gugger, M. Drame, Q. Lhoest, and A. M. Rush. *Transformers: State-of-the-Art Natural Language Processing.* Association for Computational Linguistics, 2020.

[16] Q. Lhoest, A. S. Villanova, P. von Platen, S. Patil, M. Drame, Y. Jernite, J. Plu, C. Ma, L. Tunstall, S. Patil, A. Köpf, S. Trivedi, M. Bornea, J. Chaumond, C. Delangue, V. Sanh, L. Debut, T. Wolf, A. M. Rush, and AutoMQ. *Datasets: A Community Library for Natural Language Processing.* Association for Computational Linguistics, 2021.