

<https://doi.org/10.53032/tvcr/2025.v7n2.44>

Research Article

Comparative Analysis of Embedding Models for Hindi-English Code-Mixed University related queries

Om Ingale

Student, Department of Data Science,
Mumbai University, Kirti College, India
omingale0707@gmail.com

Dr. Sampada Margaj

Assistant professor,
Department of Computer Science,
Kirti College, India
sampada.malhar@gmail.com

Abstract

This study presents a comparative analysis of open source embedding models for developing a understanding Hindi-English code-mixed language on university related questions. With the increasing adoption of conversational agents in Indian higher education institutions, there is a need for systems that can effectively process queries containing mixed Hindi and English language elements. This research evaluates the performance of five state-of-the-art embedding models - MuRIL, IndicBERT, XLM-RoBERTa, mBERT, on a custom dataset of university-related Hindi-English code-mixed queries. These models were assessed across key metrics including intent classification accuracy, entity recognition performance, and computational efficiency. The results indicate that MuRIL consistently outperforms other models, achieving 87.3% intent classification accuracy and 84.2% entity recognition F1-score, representing a 12.8% improvement over the other models. Analysis across varying code-mixing levels reveals that MuRIL maintains robust performance even with high mixing indices, while other models show significant degradation. This research provides practical insights for educational institutions seeking to implement linguistically inclusive chatbot systems and contributes to the growing body of knowledge on multilingual NLP applications in educational contexts.

Keywords: Embedding Models, Code-Mixing, Hindi-English, Natural Language Processing

The Voice of Creative Research

Vol. 7 & Issue 2 (April 2025)

1. Introduction

In recent years, educational institutions have increasingly adopted conversational agents or chatbots, which have simplified access to information and services for students, faculty, and staff. In India, these systems encounter a unique linguistic challenge due to the prevalent use of code-mixing, particularly the blending of Hindi and English in conversations. Code-mixing, which involves switching between two or more languages within a single conversation or sentence, is a common linguistic practice among Indian speakers [1].

When engaging with university information systems, Indian students often pose questions that incorporate both Hindi and English. For example, queries like "Admission process kab start hoga?" (When will the admission process start?) or "Library ki timing kya hai weekend par?" (What are the library hours on weekends?) are typical examples of code-mixed interactions. Traditional chatbot systems using single language find it challenging to accurately interpret such queries, leading to a not so good user experience and limited functionality.

Recent advancements in natural language processing (NLP) have led to the development of various embedding models with multilingual capabilities. However, their specific effectiveness for handling Hindi-English code-mixed queries in educational settings has not been thoroughly investigated. This research seeks to fill this gap by assessing 4 leading embedding models - MuRIL, IndicBERT, XLM-RoBERTa, and mBERT using a dataset of university-related Hindi-English code-mixed queries.

The primary contributions of this research is a comparative evaluation of state-of-the-art multilingual embedding models on Hindi-English code-mixed text in the university domain

By exploring these questions, this research offers practical insights for educational institutions aiming to implement linguistically inclusive chatbot systems and enhances the understanding of multilingual NLP in educational contexts.

2. Literature Review

2.1 Code-Mixing in Natural Language Processing

Over the last ten years, research into code-mixing within NLP has made substantial progress. Bali et al. (2014) were among the pioneers in conducting an in-depth analysis of Hindi-English code-mixing patterns on social media[1], where they introduced metrics to measure mixing levels and identified linguistic trends. Their study underscored the widespread nature of this phenomenon in Indian digital communication and highlighted the necessity for tailored NLP methods. Das and Gambäck (2014) introduced the Code-Mixing Index (CMI) as a standardized measure for assessing the degree of language mixing in text, which has since become a widely used evaluation standard[2]. Tackling the core issue of language identification in code-mixed text, Barman et al. (2014) developed word-level language identification techniques that achieved 95% accuracy on Bengali-Hindi-English trilingual social media content. In more recent developments, Pratapa et al. (2018) investigated the generation of synthetic data for code-mixed languages, showing that linguistically driven constraints could create training data that enhances language modeling performance. Aguilar et al. (2020) proposed methods for named entity recognition specifically tailored for code-mixed social

media text, emphasizing the difficulties in detecting entity boundaries across language shifts[3]. In the Indian scenario, Khanuja et al. (2020) introduced GLUECoS, a benchmark designed to evaluate systems on Hindi-English code-mixed tasks such as sentiment analysis, question answering, and natural language inference[4]. This benchmark has enabled standardized comparisons across different modeling approaches for code-mixed NLP tasks.

2.2 Multilingual Embedding Models

The advent of context-aware embedding models has revolutionized multilingual NLP capabilities. In 2019, Devlin et al. introduced multilingual BERT (mBERT), which was trained using Wikipedia data from 104 languages, including Hindi and English. Although mBERT wasn't specifically tailored for code-mixing, it exhibited impressive cross-lingual transfer abilities. Building on the mBERT framework, Conneau et al. (2020) developed XLM-RoBERTa, utilizing a much larger dataset (2.5TB of cleaned CommonCrawl data) and enhanced training techniques[6]. This model significantly outperformed mBERT on the XNLI cross-lingual benchmark, including for Hindi.

To cater to the unique requirements of Indian languages, Kakwani et al. (2020) introduced IndicBERT, an ALBERT-based model pre-trained on 12 major Indian languages[7]. This model was crafted to be computationally efficient while delivering strong performance on Indian language tasks. Subsequently, Khanuja et al. (2021) launched MuRIL (Multilingual Representations for Indian Languages), a BERT-based model specifically designed for Indian languages, with focused training on code-mixed content. MuRIL achieved state-of-the-art results on various Indian language benchmarks[8].

Pires et al. (2019) examined the zero-shot cross-lingual capabilities of mBERT, discovering that it performed unexpectedly well on tasks involving languages with similar scripts, but faced challenges with languages that had different scripts or limited training data[9]. This insight is relevant for Hindi-English code-mixing, where the script difference (Devanagari vs. Latin) adds complexity.

2.3 Educational Chatbots and Domain Adaptation

Research on chatbots in educational contexts has focused primarily on pedagogical applications or administrative functions. Winkler and Söllner (2018) surveyed educational chatbot implementations, identifying key application areas including student advising, learning support, and administrative assistance. They noted that most systems were limited to a single language, creating accessibility barriers in multilingual environments.

Hien et al. (2018) developed a university chatbot for Vietnamese students, documenting the challenges of domain-specific language understanding and the need for specialized training data. Their work highlighted the importance of incorporating institutional terminology and context for accurate query interpretation.

In the Indian context, Gaikwad et al. (2020) implemented a university information chatbot, noting significant challenges when students used mixed Hindi-English queries rather than the expected English-only inputs. Their findings reinforced the need for chatbot systems that can process code-mixed language.

The Voice of Creative Research

Vol. 7 & Issue 2 (April 2025)

More recently, Srivastava and Singh (2022) explored student satisfaction with university chatbot systems in India, finding that language handling capabilities significantly impacted user experience. Students reported higher satisfaction with systems that could process Hindi-English mixed queries compared to English-only systems.

3. Methodology

3.1 Dataset Construction

To ensure relevance and authenticity, we created a custom dataset of Hindi-English code-mixed university queries. To ensure comprehensive coverage of less common queries and maintain class balance, we generated 1000 additional queries using templates and linguistic patterns observed in the authentic data.

The final dataset comprised 1000 queries spanning 12 domains relevant to university operations: admissions, academics, examinations, library services, hostel facilities, financial matters, scholarships, placements, extracurricular activities, administrative procedures, technical services, and campus information.

Each query was manually annotated with:

- Primary intent (from a taxonomy of 42 distinct intents)
- Named entities (courses, departments, facilities, etc.)
- Code-Mixing Index (CMI) following Das and Gambäck's (2014) methodology
- Script variations (i.e., whether Hindi portions used Devanagari or Roman script)

To validate annotation quality, we employed two native Hindi-English bilingual speakers with education domain expertise to independently annotate a subset of 500 queries, achieving an inter-annotator agreement (Cohen's kappa) of 0.87 for intent classification and 0.83 for entity labeling.

Query Example	Intent	Entities	CMI
"B.Tech CSE ki admission process kab start hogi?"	inquiry_admission_timeline	program:B.Tech, department:CSE	0.63
"Library weekend par khuli rehti hai kya?"	inquiry_facility_hours	facility:Library, time:weekend	0.57
"Assignment submission ki last date extend ho sakti hai kya?"	request_deadline_extension	submission_type:assignment	0.64
"Placement ke liye minimum CGPA requirement kya hai?"	inquiry_eligibility_criteria	process:placement, criteria:CGPA	0.50

The Voice of Creative Research

Vol. 7 & Issue 2 (April 2025)

Table 1: Representative examples from the dataset

We split the dataset into training (70%), validation (15%), and test (15%) sets, ensuring proportional representation of domains, intents, and CMI levels across splits.

3.2 Embedding Models

We evaluated five state-of-the-art multilingual embedding models with distinct architectural characteristics:

1. **MuRIL** (Multilingual Representations for Indian Languages, Khanuja et al., 2021): A BERT-based model specifically pre-trained on 17 Indian languages including Hindi, with explicit training on code-mixed text. We used the base cased version with 12 transformer layers.
2. **IndicBERT** (Kakwani et al., 2020): An ALBERT-based model pre-trained on 12 Indian languages, optimized for computational efficiency with parameter sharing techniques. This model uses 12 transformer layers with shared parameters.
3. **XLM-RoBERTa** (Conneau et al., 2020): A RoBERTa-based model trained on 100 languages with a focus on cross-lingual transfer capabilities. We used the base version with 12 transformer layers.
4. **mBERT** (multilingual BERT, Devlin et al., 2019): The original multilingual version of BERT, trained on Wikipedia corpora from 104 languages. We used the base cased version with 12 transformer layers.

All models were implemented using the Hugging Face Transformers library (Wolf et al., 2020), maintaining their original architecture and tokenization approaches.

3.3 Experimental Setup

We integrated each embedding model into a consistent intent classification and entity recognition pipeline to ensure fair comparison. The architecture consisted of:

1. **Tokenization Layer:** Using each model's native tokenizer
2. **Embedding Layer:** The pre-trained model generating contextual word embeddings
3. **Intent Classification Head:** A feed-forward neural network with one hidden layer (256 units) and softmax output for 42 intent classes
4. **Entity Recognition Head:** A linear-chain conditional random field (CRF) on top of token embeddings for sequence labeling

We implemented the following experimental conditions:

1. **Zero-shot Evaluation:** Testing the pre-trained models without any task-specific fine-tuning
2. **Full Fine-tuning:** Fine-tuning all parameters on our training dataset
3. **Adapter-based Fine-tuning:** Using parameter-efficient adapter modules (Pfeiffer et al., 2020) while keeping base model weights frozen
4. **Layer-wise Fine-tuning:** Progressively unfreezing layers to determine optimal fine-tuning depth

Training hyperparameters were standardized across models where possible:

- Learning rate: $3e-5$ with linear decay

- Batch size: 32
- Maximum sequence length: 128 tokens
- Training epochs: 5 (with early stopping based on validation performance)
- Optimizer: AdamW with weight decay of 0.01
- Dropout rate: 0.1

3.4 Evaluation Metrics

We evaluated model performance using multiple complementary metrics:

- **Accuracy and F1-score for intent classification**
- **Precision, Recall, and F1-score for named entity recognition**

To understand domain-specific performance variations, we calculated metrics separately for each of the 12 university domains and conducted statistical significance testing (paired t-tests with Bonferroni correction) for performance differences between models.

4. Results

4.1 Intent Classification Performance

Table 2 presents the Intent Classification Performance of each embedding model on the test set after full fine-tuning:

Model	Intent Accuracy	Intent F1
MuRIL	87.3%	0.865
IndicBERT	81.5%	0.802
XLM-RoBERTa	79.2%	0.775
mBERT	74.5%	0.732

Table 2: Intent classification performance for fine-tuned embedding models

MuRIL demonstrated superior performance across all primary metrics, achieving 87.3% intent classification accuracy. This represents a statistically significant improvement ($p < 0.01$) over all other models. IndicBERT achieved the second-best performance while demonstrating the lowest inference time, making it potentially suitable for resource-constrained deployments.

The zero-shot evaluation (without fine-tuning) showed much lower performance across all models, with MuRIL still leading at 48.7% intent accuracy, followed by XLM-RoBERTa (45.2%), IndicBERT (43.8%), mBERT (39.4%). This indicates that while these models have some inherent capability to process code-mixed text, fine-tuning is essential for practical applications.

4.2 Named Entity Recognition

Named entity recognition results are shown in Table 3, broken down by entity type.

Entity Type	XLM-RoBERTa	MuRIL	IndicBERT	mBERT

The Voice of Creative Research

Vol. 7 & Issue 2 (April 2025)

Course	0.85	0.89	0.82	0.79
Department	0.83	0.87	0.80	0.78
Procedure	0.76	0.81	0.74	0.72
Date/Time	0.82	0.84	0.81	0.78
Location	0.80	0.85	0.78	0.76
Person	0.79	0.82	0.77	0.75
Overall	0.81	0.83	0.79	0.76

Table 3: Named Entity Recognition F1-scores (Fine-tuned models)

MuRIL outperformed other models across all entity types, with particularly strong performance in recognizing course names and departments. This may be attributed to its training on Indian academic content.

5. Conclusion and Future Work

This research demonstrates that embedding models pre-trained on Indian language corpora, particularly MuRIL, significantly outperform general multilingual models for Hindi-English code-mixed text in university chatbot applications. Our implementation framework provides a practical approach for deploying these models in production environments.

Future research directions include:

1. Extending the approach to other Indian language pairs common in educational settings
2. Developing specialized pre-training objectives for code-mixed text understanding
3. Incorporating conversational context for improved multi-turn interactions
4. Exploring few-shot learning approaches to reduce annotation requirements for new domains

Limitations of our study include the regional focus on North Indian universities and the relatively small dataset size. Nevertheless, our findings provide valuable insights for developing more effective natural language understanding systems for code-mixed educational contexts.

References

1. Altrabsheh, N., Cocea, M., & Fallahkhalil, S. (2019). Smart Learning Environments for Higher Education: Chatbots for Student Support. *IEEE Access*, 7, 177387-177395.
2. Bali, K., Sharma, J., Choudhury, M., & Vyas, Y. (2014). "I am borrowing ya mixing?" An Analysis of English-Hindi Code Mixing in Facebook. *Proceedings of the First Workshop on Computational Approaches to Code Switching*, 116-126.
3. Bhat, I., Bhat, R. A., Shrivastava, M., & Sharma, D. (2018). Universal Dependency Parsing for Hindi-English Code-Switching. *Proceedings of the 2018 Conference of the*

The Voice of Creative Research

Vol. 7 & Issue 2 (April 2025)

- North American Chapter of the Association for Computational Linguistics: *Human Language Technologies*, 987-998.
4. Conneau, A., Khandelwal, K., Goyal, N., Chaudhary, V., Wenzek, G., Guzmán, F., Grave, E., Ott, M., Zettlemoyer, L., & Stoyanov, V. (2020). Unsupervised Cross-lingual Representation Learning at Scale. *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 8440-8451.
 5. Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 4171-4186.
 6. Feng, F., Yang, Y., Cer, D., Arivazhagan, N., & Wang, W. (2020). Language-agnostic BERT Sentence Embedding. *arXiv preprint arXiv:2007.01852*.
 7. Guzmán, F., Bouamor, H., Baly, R., & Habash, N. (2016). Machine Translation Evaluation for Arabic using Morphologically-enriched Embeddings. *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, 1398-1408.
 8. Khanuja, S., Bansal, S., Mehtani, P., Khosla, S., Dey, A., Gopalan, B., Margam, D.K., Aggarwal, P., Nagipogu, R.T., Dave, S., Gupta, S., Khanna, S.C., Kumar, V., & Talukdar, P. (2021). MuRIL: Multilingual Representations for Indian Languages. *arXiv preprint arXiv:2103.10730*.
 9. Khanuja, S., Dandapat, S., Srinivasan, A., Sitaram, S., & Choudhury, M. (2020). GLUECoS: An Evaluation Benchmark for Code-Switched NLP. *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 3575-3585.
 10. Kumar, N., & Bhattacharyya, P. (2021). Adaptive Pre-training for Effective Code-Mixed Natural Language Understanding. *Proceedings of the 5th Workshop on Computational Approaches to Linguistic Code-Switching*, 29-40.
 11. Chand, R.R., Sharma, N.A. (2023). Development of Bilingual Chatbot for University Related FAQs Using Natural Language Processing and Deep Learning. In: Hsu, CH., Xu, M., Cao, H., Baghban, H., Shawkat Ali, A.B.M. (eds) *Big Data Intelligence and Computing. DataCom 2022. Lecture Notes in Computer Science*, vol 13864. Springer, Singapore. https://doi.org/10.1007/978-981-99-2233-8_6