# LINGUALINK — LINKS LANGUAGES AND PEOPLE

**Swati Rathod[1], Mayuri Ambore[2], Krishna Kanade[3]**

*1Computer Science &Engineering, Babasaheb Naik College Of Engineering, Pusad*
*Sant Gadgebaba University, Amravati, Maharastra, Indi-445215*
*Swatirathod682@gmail.com*
*2Computer Science &Engineering, Babasaheb Naik College Of Engineering, Pusad*
*Sant Gadgebaba University, Amravati, Maharastra, India-445215*
*amboremayuri38@gmail.com*
*3Computer Science &Engineering, Babasaheb Naik College Of Engineering, Pusad*
*Sant Gadgebaba University, Amravati, Maharastra, India-445215 [3]*
*Krushnakanade2003@gmail.com*

**ABSTRACT**

*The rapid growth of digital platforms in India has increased the demand for vernacular language support, especially in a region with high linguistic diversity. This paper introduces LinguaLink, a multilingual chatbot designed to facilitate communication and information access in Hindi, English, and Marathi. Leveraging advanced NLP techniques and a fine-tuned transformer model (like MuRIL BERT), LinguaLink provides accurate, context-aware responses to user queries in their chosen language. The system's architecture, methodology for corpus creation, and performance evaluation metrics are detailed. The results demonstrate that LinguaLink effectively removes the need for expensive machine translation layers, offering a seamless, real-time conversational experience that is vital for connecting people across linguistic barriers.*

## I. INTRODUCTION

India is one of the most linguistically diverse countries in the world, with over 22 officially recognized languages and hundreds of regional dialects spoken across its vast geography. As digital adoption surges, particularly through mobile-first platforms, there is a growing need to make online services accessible in multiple regional languages. However, most existing conversational AI systems in India are limited to English or rely on costly and error-prone machine translation pipelines, which often fail to capture the nuances of vernacular languages

In this context, the development of multilingual conversational agents that can understand and respond natively in regional languages is essential for fostering digital inclusivity. This paper presents **LinguaLink**, a multilingual chatbot designed to support Hindi, English, and Marathi — three widely spoken languages in India. By leveraging transformer-based architectures, specifically a fine-tuned version of MuRIL BERT, LinguaLink offers real-time, context-aware interactions without the need for intermediary translation systems. This work details the system's overall architecture, the process of multilingual corpus creation, and the evaluation metrics used to assess its performance. Through comprehensive testing, LinguaLinkS demonstrates its capability to deliver a seamless, high-quality conversational experience across languages, thereby contributing to the broader goal of democratizing information access in multilingual societies.

## II. METHODS AND MATERIAL [ Page Layout ]

### System Architecture

LinguaLink is built upon a transformer-based architecture, utilizing a fine-tuned MuRIL BERT model specifically designed for Indian languages. The system supports multilingual input by detecting the language of the query and processing it directly without reliance on machine translation. The architecture comprises modules for language identification, intent recognition, response generation, and context management to ensure coherent and context-aware conversations.

### Corpus Creation

To train and fine-tune the chatbot, a multilingual dataset was curated comprising user queries and responses in Hindi, English, and Marathi. Data sources included publicly available conversational datasets, crowdsourced dialogues, and domain-specific question-answer pairs. The corpus was preprocessed through normalization, tokenization, and language-specific text cleaning to handle unique linguistic features such as script variations and idiomatic expressions.

### Model Training and Fine-Tuning

The MuRIL BERT model was fine-tuned on the prepared corpus using transfer learning techniques to enhance its understanding of the target languages' syntax and semantics. Hyperparameters such as learning rate, batch size, and number of epochs were optimized through cross-validation to maximize performance. Training was conducted on GPUs to handle computational demands efficiently.

### Evaluation Metrics

Performance was evaluated using standard NLP metrics including accuracy, precision, recall, and F1-score for intent classification and response relevance. Additionally, human evaluation was conducted to assess conversational fluency, context retention, and user satisfaction. Latency measurements were also recorded to ensure the system supports real-time interaction.

## III. RESULTS AND DISCUSSION [Page Style ]

IV. **Performance Metrics:** Evaluate the chatbot's effectiveness using metrics such as:

    A. **Accuracy:** Percentage of correct responses.

    B. **Response Time:** Latency of the chatbot's replies.

    C. **F1 Score:** Measures the model's accuracy, considering both precision and recall.

V. **Qualitative Analysis:** Discuss user feedback and observations on the conversational flow and usefulness of the chatbot.

VI. *Comparative Analysis: Compare LinguaLink's performance with a baseline model using traditional machine translation.*

Incorporate sentiment analysis to provide more emotionally intelligent responses

## . II. Literature review

- **Evolution of chatbots:** From rule-based systems to modern AI-powered conversational agents.
- **Multilingual NLP challenges:** Discussing the scarcity of resources for low-resource Indian languages, dialectal variations, and code-switching.
- **Existing multilingual chatbot models:** Reviewing recent approaches, including those using multilingual Pre-trained Language Models (PLMs) like mT5 and mBERT.
- **Multilingual datasets for Indian languages:** Referencing existing work like the Hindi and Marathi QA datasets

## III. Methodology
## A. System architecture
The LinguaLink architecture would include:

- **User Interface (UI):** A simple chat window where users can type their queries.
- **Language Detection Module:** This initial module auto-detects the input language (English, Hindi, or Marathi).
- **Natural Language Understanding (NLU) Engine:**
  - **Intent Recognition:** Classifies the user's query into predefined categories.
  - **Entity Extraction:** Identifies and extracts key information from the query.
- **Multilingual Transformer Model (e.g., MuRIL BERT):** This is the core of the system. It is fine-tuned for question-answering on a custom dataset of fixed-response queries in Hindi, English, and Marathi.
- **Dialogue Management:** Manages the flow of the conversation, maintains context, and retrieves the most relevant response from the knowledge base.
- **Knowledge Base/Database:** Stores the fixed-response questions and their corresponding answers in all three supported languages.
- **Natural Language Generation (NLG):** Generates the final, human-like response in the detected language.

## B. Dataset creation
- **Corpus Sourcing:** Gathering a parallel dataset of questions and answers in English, Hindi, and Marathi. This could involve translating an existing corpus (like SQuAD) or building one from scratch.
- **Data Cleaning and Pre-processing:** Tokenization, handling of code-switching, and normalization of text.

## C. Model training
- **Pre-trained Model Selection:** The MuRIL BERT model is a strong candidate, as it is pre-trained on Indian languages and has shown robust performance.
- **Fine-tuning:** The MuRIL model is fine-tuned on the custom-created, fixed-response QA dataset for the downstream task of question-answering.

Conclusion

LiguaLink successfully addresses the linguistic challenges posed by India's diverse language landscape by offering an efficient, multilingual

chatbot solution. Its use of advanced NLP and transformer models ensures accurate and contextually relevant communication without relying on costly translation services. This innovation not only enhances user engagement across Hindi, English, and Marathi speakers but also sets a precedent for scalable, real-time multilingual applications in the digital ecosystem.

- **Summary of Findings: Reiterate the effectiveness of LinguaLink in providing an efficient multilingual conversational experience.**
- **Future Work:**
  - Expand language support to include more Indian languages.
  - Integrate the chatbot with other platforms and messaging apps.
  - Explore advanced dialogue management for multi-turn conversations and context switching.

Incorporate sentiment analysis to provide more emotionally intelligent responses

REFERENCE

- Ayyalasomayajula, R., et al. (2023). Multilingual Chatbot for Indian Languages. ResearchGate.
- Haptik. (2022). Multilingual Chatbots Making Conversational AI Vernacular.
- Rasa Forum. (2021). Multilingual Chatbot for Indian Languages.
- A Question-Answering Dataset for Hindi and Marathi. (2024). arXiv.
- Analytics Vidhya. (2024). How to Build a Multilingual Chatbot using Large Language Models?
- Patidar, P., et al. (2023). Hindi Chatbot for Supporting Maternal and Child Health. ACL Anthology.
- Patidar, P., et al. (2023). Hindi Chatbot for Supporting Maternal and Child Health Related Queries. ResearchGate.
- Bhosale, S., et al. (2022). A Novel Marathi Speech-Based Question and Answer Chatbot for the Educational Domain. ResearchGate.
- M. S. Wani, R. S. Kulkarni. (2023). Study and evaluation of "Seq-2-Seq" model competency in AI-based educational chatbot for the Marathi language. ResearchGate.
- ijrpr. Agriculture Chatbot Marathi. ijrpr.com.