# SAHYADRI
## COLLEGE OF ENGINEERING & MANAGEMENT
### An Autonomous Institution
### MANGALURU

**Department of Computer Science & Engineering**

**SYNOPSIS**

| 1. | **Title of the Project** | Fake News Detection Using Natural Language Processing for Kannada Language | | |
|----|--------------------------|-------------------------------------------------------------------------|---|---|
| 2. | **Group Number** | 67 | | |
| 3. | **Name of the Students** | MUHAMMED AQUIF | 4SF22CS114 | |
| | | PRATHAM AMIN | 4SF22CS145 | |
| | | VISHWA | 4SF22CS247 | |
| | | B M SHASHANK | 4SF23CS402 | |
| 4. | **Guide** | Mrs. SUKETHA | Assistant Professor, Dept of CSE, SCEM | |

**Abstract**

The rise of fake news has significantly impacted society by spreading misinformation across various platforms, especially social media. This project proposes a solution using multilingual Natural Language Processing (NLP) techniques to detect and classify fake news effectively. By leveraging transformer-based models such as BERT and XLM-RoBERTa, we aim to build a system that can identify fake news in multiple languages. The model will be fine-tuned using regional language datasets to increase its relevance and accuracy across diverse linguistic audiences. The rapid spread of misinformation across digital platforms has made fake news detection a critical challenge in today's information-driven world. While numerous solutions have been developed for detecting fake news in English, there remains a significant gap in the availability and accuracy of models that support multiple languages.The approach involves preprocessing multilingual datasets, applying language-agnostic feature extraction, and training deep learning models to identify misleading or fabricated content. This work aims to contribute to the development of more inclusive and reliable misinformation detection tools, which can serve a diverse, global population more effectively

*Keywords:*

Fake News Detection, Multilingual NLP, Natural Language Processing, Text Classification, Machine Learning.

# 1. Introduction

In recent years, the proliferation of fake news has emerged as a critical challenge for societies worldwide. With the increasing reliance on digital platforms for news consumption, misinformation and disinformation have found fertile ground to spread rapidly, often outpacing fact-based reporting. The consequences are far-reaching, ranging from influencing public elections, exacerbating health crises (such as during the COVID-19 pandemic), and fueling social unrest to eroding trust in media and democratic institutions

The Traditionally, fake news detection systems have relied on a combination of linguistic cues, metadata, user behavior, and fact-checking resources. However, the majority of these systems are developed for English or other widely used languages, overlooking the multilingual nature of the internet and its users. This creates a significant gap in global fake news mitigation efforts, especially in regions with limited language-specific resources or underrepresented languages.

To address this, our project explores Multilingual Natural Language Processing (NLP) as a core approach for detecting fake news across different languages. With the rise of transformer-based language models such as Multilingual BERT (mBERT) and XLM-RoBERTa, it has become possible to build models that understand and analyze text in dozens of languages using shared semantic representations. The goal of this project is not only to detect fake news with high accuracy but also to ensure that the system is scalable, inclusive, and effective across various linguistic contexts. Through this approach, we aim to contribute to the global fight against misinformation and promote a more trustworthy information environment for all users.

The unprecedented growth of digital media and social networking platforms has revolutionized the way information is created, consumed, and shared. While this digital revolution has enabled real-time access to news and global connectivity, it has also led to the alarming spread of fake news—misleading, fabricated, or false content that is often disseminated with malicious intent. The consequences of fake news are far-reaching, affecting political stability, public health, financial markets, and social harmony. As a result, there is an urgent need for intelligent and automated systems capable of identifying and mitigating the impact of such misinformation. Traditional fake news detection systems primarily rely on machine learning and Natural Language Processing (NLP) techniques, but they are typically monolingual and thus fail to address the challenges posed by our linguistically diverse world. In many multilingual countries and global news ecosystems, fake news is circulated in numerous languages, making it critical to develop models that can understand and process text in more than one language. This project aims to tackle this issue by employing Multilingual Natural Language Processing approaches, utilizing cutting-edge models such as multilingual BERT (mBERT), XLM-RoBERTa, and other transformer-based architectures trained on vast corpora of multilingual data.

These models are capable of capturing contextual semantics across languages, enabling more accurate classification of fake and real news articles regardless of the language they are written in. By integrating multilingual capabilities into fake news detection systems, we not only enhance their reach and reliability but also take a significant step toward ensuring access to credible information in a globally inclusive manner. Ultimately, this research contributes to the broader fight against misinformation and supports efforts to uphold truth, transparency, and trust in the digital age.

## 2. Literature Survey

| NAME | AUTHORS | YEAR | DESCRIPTION |
|---|---|---|---|
| Fake News Detection in low-resources languages [1] | Sivanaiah, R., Ramanathan, N., Hameed, S., Rajagopalan, R., Suseelan, A.D., Thanagathai, M.T.N. | 2023 | This paper provided detecting fake news in high-resource languages like Kannada and English, low-resource languages often lag due to limited datasets, linguistic tools, and research attention. |
| Intersection of Machine Learning, Deep Learning and Transformers to Combat Fake News in Kannada Language [2] | S.Sanjana, S. Kuranagatti, J. G. Devisetti, R. Sharma and A. Arya | 2023 | The proliferation of fake news on digital platforms has necessitated the development of automated detection systems, especially for regional languages like Kannada. This paper explores the convergence of Machine Learning (ML), Deep Learning (DL) and Transformer-based models in addressing the challenge of fake news detection in the Kannada language. |
| Multilingual Fake News Detection in Low-Resource Languages: A Comparative Study Using BERT and GPT-3.5 [3] | Anirudh, K., Srikanth, M., Shahina, A. | 2024 | This paper investigates the effectiveness of large language models in detecting fake news across multiple languages, with a focus on low-resource settings. The study provides a comparative evaluation between BERT-based models and GPT-3.5, analyzing their performance in multilingual contexts where annotated data is scarce. |

| | | | |
|---|---|---|---|
| Multilingual Misinformation Detection: Deep Learning Approaches for News Authenticity Assessment [4] | Sushma S. Nandgaonkar, J. Shaikh, G. B. Bhore, R. V. Kadam and S. S. Gadhave | 2024 | This paper does the recent research in misinformation detection has focused on leveraging deep learning techniques due to their superior ability to learn contextual and semantic features from large-scale data. Traditional approaches relieve manual fact-checking or rule-based systems, which are not scalable or adaptable to multilingual content. |
| Fake news detection in Dravidian languages using multiscale residual CNN_BiLSTM hybrid model [5] | Eduri Raja, Badal Soni, Samir Kumar Borgohain | 2024 | This study proposes a novel hybrid deep learning model combining Multiscale Residual Convolutional Neural Networks (CNNs) with Bidirectional Long Short-Term Memory (BiLSTM) networks for fake news detection in Dravidian languages. |
| Fake News Detection in Dravidian Languages Using Transformer Models [6] | E.Raja, B.Soni, S.K.Borgohain | 2024 | This paper focuses on the fake news threat to the low-resources languages like Kannada, Tamil, Telugu and other Dravidian languages that are in India to detect them and justify that news. |
| Multi-Modal Categorization of News Through Varied Machine Learning Techniques and Models [7] | S.U.Priya, Shamita S., P.B.Honnavali, Sivaraman Eswaran | 2022 | This study encloses the task of news categorization that has evolved significantly with the advent of machine learning (ML) and deep learning techniques. Traditional text-based classification methods relied heavily on natural language processing (NLP) approaches such as TF-IDF, bag-of-words, and word embeddings in combination with classifiers like SVM, Naive Bayes, and decision trees. |

| | | | |
|---|---|---|---|
| Factorization of Fact-Checks for Low Resource Indian Languages [8] | S Singhal, RR Shah, P Kumaraguru | 2021 | This paper proposes the significant challenges due to limited linguistic resources, diverse scripts, and the complex socio-political landscape. languages remain underexplored. |
| Deciphering Deception: Unmasking Fake News in Multilingual Contexts [9] | A. Agarwal, Y. P. Singh and V. Rai | 2024 | This survey reviews key research contributions in the domains of fake news detection, multilingual natural language processing (NLP), and cross-lingual information credibility. |
| . Indian Language Analysis with XLM-RoBERTa: Enhancing Parts of Speech Tagging for Effective Natural Language Preprocessing [10] | K. K. Jayanth, G. Bharathi Mohan and R. P. Kumar | 2023 | This study proposes the Indian languages present unique challenges for NLP tasks such as Parts of Speech (POS) tagging, owing to their rich morphology, syntactic variation, and limited annotated resources. Traditional models, including rule-based and statistical methods, often fall short in addressing these complexities. |
| Exploring Social Media Trends - A Kannada Dataset Analysis [11] | A. Dey, Aishwaryasri J, Jai Surya R, Jayanthi Mg and Prashanth Kannadaguli | 2023 | This paper provides the context of Kannada language analysis, challenges arise due to limited annotated datasets and the complexity of processing Kannada script in natural language processing (NLP) models. Existing efforts to analyze Kannada social media data have focused on sentiment classification and basic keyword extraction but often lack comprehensive trend analysis over large datasets. |

| | | | |
|---|---|---|---|
| Monolingual and Multilingual Misinformation Detection for Low-Resource Languages: A Comprehensive Survey [12] | Xinyu Wang, Wenbo Zhang, Sarah Rajtmajer | 2024 | This survey examines the existing literature on both monolingual and multilingual approaches to misinformation detection, focusing specifically on low-resource languages. |
| Kannada-English Code-Mixed Speech Synthesis [13] | S. K. Suresh and U. Damotharan | 2024 | This study has the c ode-mixed speech synthesis involves generating natural and intelligible speech from text that contains multiple languages, often mixed within sentences. In the context of Kannada-English code-mixing, this task is particularly challenging due to the differences in phonetics, syntax, and prosody between the two languages. |
| Cross-lingual and Multilingual Spoken Term Detection for Low-Resource Indian Languages [14] | S Shah, S Guha, S Khanuja, S Sitaram | 2020 | This paper provide extensive research has been conducted on STD for high-resource languages such as English and Mandarin, low-resource Indian languages pose significant challenges due to limited annotated data, diverse dialects, and complex phonetic variations. |
| Fake news detection using natural language processing [15] | M.J Varma, M.S Rohit, G.S.G Selvi | 2025 | The proliferation of fake news on digital platforms has become a significant challenge, prompting extensive research in automated fake news detection. Natural Language Processing (NLP) techniques play a pivotal role in this domain by analyzing textual data to identify deceptive or misleading content. |

## 3.  Problem Statement

### 3.1  Existing Problem Statement:

Existing fake news detection systems often struggle with Kannada due to limited language-specific resources and models. This means Kannada-language news, which is crucial for the local population, faces challenges in accurate detection of fake news, leading to potential misinformation and distrust in credible sources.

### 3.2  Proposed Problem Statement:

To develop Kannada-specific model and datasets that can effectively identify and classify fake news within the context of the Kannada language and culture. This includes addressing the unique linguistic features, common fake news tactics, and online media landscape prevalent in Kannada.

## 4.    Objectives

The key objectives of the project are:

- To develop an NLP Model to create a robust NLP-based model capable of analyzing and classifying news articles in Kannada language.

- To build a Kannada dataset to collect and curate a comprehensive dataset containing real and fake news articles from various trusted and untrusted sources across different languages.

- To preprocess text data using preprocessing techniques tokenization, stemming, stop word removal, and language translation (if necessary) to prepare data for model training.

- To Train and Evaluate the Model to Train the model using labeled data and evaluate it using performance metrics such as accuracy, precision, recall, and F1-score for each supported language.

## 5. Proposed Methodology

- Data Collection: The first step involves collecting a diverse and representative dataset from multiple sources. These sources can include online news platforms, social media channels, and verified fact-checking websites. Ensuring the dataset includes local languages like Kannada, English and Hindi, as well as possibly less-resourced languages to improve the model's robustness and applicability.

- Preprocessing: Once the data is collected, it undergoes preprocessing to clean and standardize the text for analysis. This includes removing HTML tags, punctuation, URLs, numbers, and unnecessary symbols. Language detection algorithms are applied to identify the language of each text instance

- Feature Extraction: The next phase is feature extraction, where meaningful representations of the text are created for input into machine learning models. Advanced multilingual language models such as mBERT (Multilingual BERT), XLM-RoBERTa, and LASER are used to generate embeddings that capture the semantic meaning of text across various languages.

- Model Development: In the model development phase, the extracted features are used to train a classification model. Transformer-based models like mBERT and XLM-RoBERTa are fine-tuned on the multilingual dataset for binary classification—determining whether a news article is fake or real.

- Model Training: To ensure the model generalizes well and avoids overfitting, cross-validation techniques are applied, typically using k-fold cross-validation. This approach helps evaluate the model's performance across different subsets of the data and improves reliability.

- Evaluation Metrics: The effectiveness of the fake news detection system is assessed using standard evaluation metrics. These include accuracy, which measures overall correctness; precision and recall, which assess how well the model identifies fake news specifically; and the F1-score, which balances precision and recall.

- Validation and Testing: After the model has been trained, it is essential to validate and test its performance on unseen data. Validation is typically done using a portion of the training data to fine-tune the model and prevent overfitting. During this phase, hyperparameters such as learning rate, batch size, and number of epochs are adjusted to find the optimal configuration.
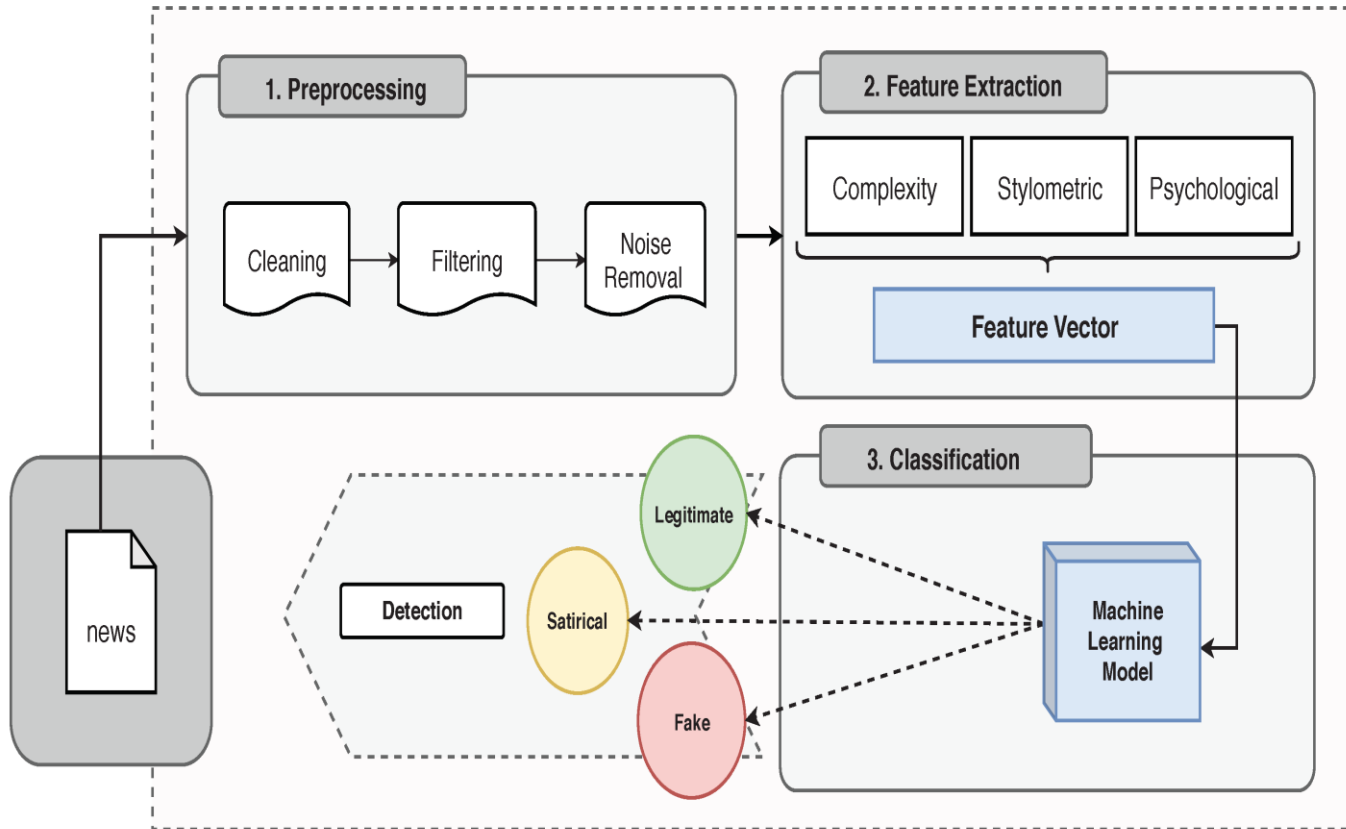
**5.1 System Design**



Fig 5.1: System Design for FND

**5.2 Functional & Non-Functional Requirements:**

### 5.2.1 Functional Requirements:

* User Interface for Input: The system should provide a user-friendly interface where users can enter news content. This can include a full article, a headline, or even a URL link.
* Language Detection: Once the input is received, the system must be capable of automatically detecting the language of the content.
* Text Preprocessing: The input text needs to go through a preprocessing stage. This includes cleaning the text by removing unnecessary punctuation, stopwords, and symbols, as well as tokenizing and normalizing it.
* Multilingual Model Integration: The core of the system is the integration of advanced multilingual NLP models.
* Fake News Classification: The processed input is analyzed using the trained model to determine whether the news is fake or real.

### 5.2.2 Non-Functional Requirements:

* Performance: The system should be capable of processing and classifying news articles within 2 seconds per article to ensure quick response times.
* Scalability: The system must be designed to scale as demand increases. It should handle thousands of articles per hour without compromising performance
* Language Support: The solution should be able to process news articles in at least 5 major languages, such as Kannada, English and Hindi.
* Accuracy: The system should ensure a high level of accuracy in fake news detection. It must achieve a minimum of 90% accuracy on standard benchmark datasets.
* Security: Data security is a priority, ensuring that all user data is processed and stored securely, with encryption where necessary.

## 6. Outcome of the Work

The project outcome for developing a Fake News Detection Module (FNDM) for preventing misinformation spread in real-time can be summarized as:

- The final system would be able to detect fake news articles in multiple languages, providing a scalable solution to combat misinformation globally.
- The tool could be integrated into social media platforms, news outlets, or fact-checking websites to alert users to potential misinformation.
- Real-time fake news detection could be made available through APIs or integrated into browser extensions or mobile apps.
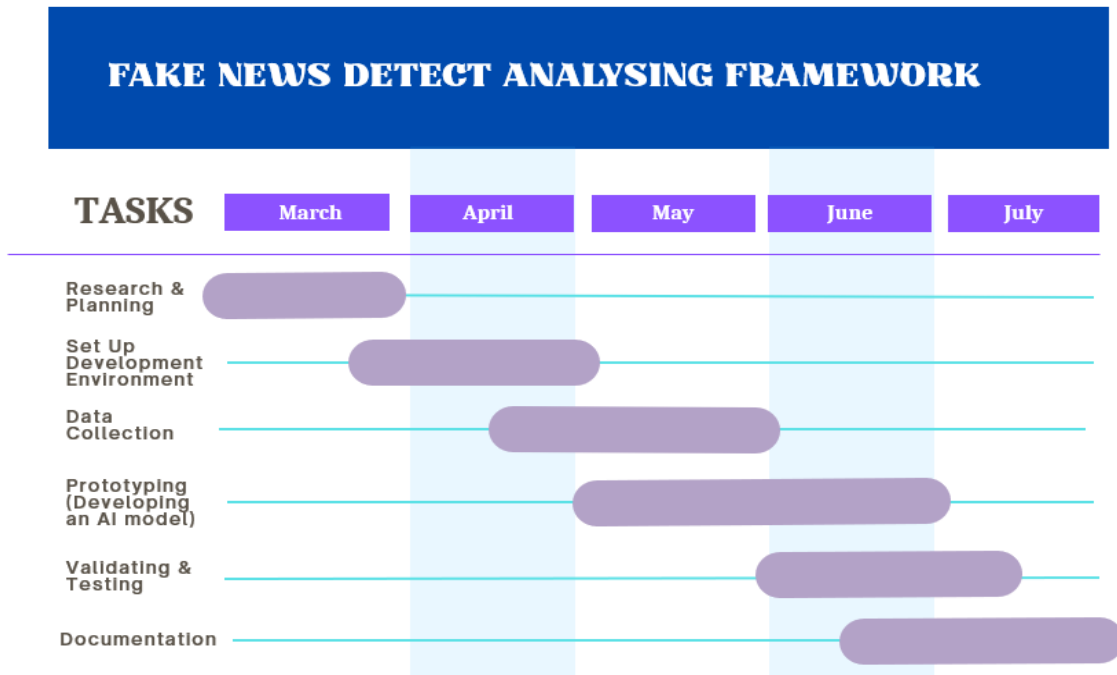
## 7. Work Plan/Gantt Chart



Fig 7.1: Project Development timeline for a FND

## 8. Conclusion

In conclusion, this study has highlighted the significant potential of leveraging Multilingual Natural Language Processing (NLP) in the detection of fake news across a variety of languages. The ability of NLP models, particularly transformer-based models such as BERT and its multilingual counterparts like mBERT and XLM-R, to process and analyze text in multiple languages simultaneously presents a powerful solution to the global issue of misinformation. By utilizing these advanced models, the study demonstrates that fake news detection can be made more inclusive and accurate for a broader linguistic audience, transcending language barriers and enabling the detection of fake news in regions where data in a single language would not suffice. However, despite the success of these models in high-resource languages such as Kannada, English and Hindi the research also reveals certain limitations. One of the primary challenges remains the availability and quality of labeled datasets in lesser-researched languages, as well as the varying cultural and contextual elements that influence how news is constructed and perceived in different parts of the world. These factors can significantly impact the performance of NLP models, sometimes leading to misclassification due to the lack of sufficient context or understanding of local nuances. Furthermore, the findings also point to the complexity of achieving high accuracy across all languages, especially when dealing with idiomatic expressions, slang, and regional variations in news reporting.

● **References**

[1] R.Sivanaiah, N.Ramanathan, S.Hameed, R.Rajagopalan, A.D.Suseelan, M.T.N.Thanagathai "Fake News Detection in Low-Resource Languages", 2023.

[2] S.Sanjana, S.Kuranagatti, J.G.Devisetti, R.Sharma and A.Arya, "Intersection of Machine Learning, Deep Learning and Transformers to Combat Fake News in Kannada Language" ,2023.

[3] K.Anirudh, M.Srikanth, A.Shahina "Multilingual Fake News Detection in Low-Resource Languages: A Comparative Study Using BERT and GPT-3.5",2024

[4] Sushma S. Nandgaonkar, J. Shaikh, G. B. Bhore, R. V. Kadam and S. S. Gadhave "Multilingual Misinformation Detection: Deep Learning Approaches for News Authenticity Assessment", 2024.

[5] Eduri Raja, Badal Soni, Samir Kumar Borgohain "Fake news detection in Dravidian languages using multiscale residual CNN_BiLSTM hybrid model", 2024.

[6] E.Raja, B.Soni, S.K.Borgohain "Fake News Detection in Dravidian Languages Using Transformer Models", 2024.

[7] S.U.Priya, Shamita S., P.B.Honnavali, Sivaraman Eswaran "Multi-Modal Categorization of News Through Varied Machine Learning Techniques and Models", 2022.

[8] S Singhal, RR Shah, P Kumaraguru "Factorization of Fact-Checks for Low Resource Indian Languagess", 2021.

[9] A. Agarwal, Y. P. Singh and V. Rai "Deciphering Deception: Unmasking Fake News in Multilingual Contexts", 2024.

[10] . K. Jayanth, G. Bharathi Mohan and R. P. Kumar "Indian Language Analysis with XLM-RoBERTa: Enhancing Parts of Speech Tagging for Effective Natural Language Preprocessing", 2023.

[11] A. Dey, Aishwaryasri J, Jai Surya R, Jayanthi Mg and Prashanth Kannadaguli "Exploring Social Media Trends - A Kannada Dataset Analysis", 2023.

[12] Xinyu Wang, Wenbo Zhang, Sarah Rajtmajer "Monolingual and Multilingual Misinformation Detection for Low-Resource Languages: A Comprehensive Survey", 2024.

[13] S. K. Suresh and U. Damotharan "Kannada-English Code-Mixed Speech Synthesis", 2024.

[14] S Shah, S Guha, S Khanuja, S Sitaram "Cross-lingual and Multilingual Spoken Term Detection for Low-Resource Indian Languages", 2020.

[15] M.J Varma, M.S Rohit, G.S.G Selvi "Fake news detection using natural language processing", 2025.