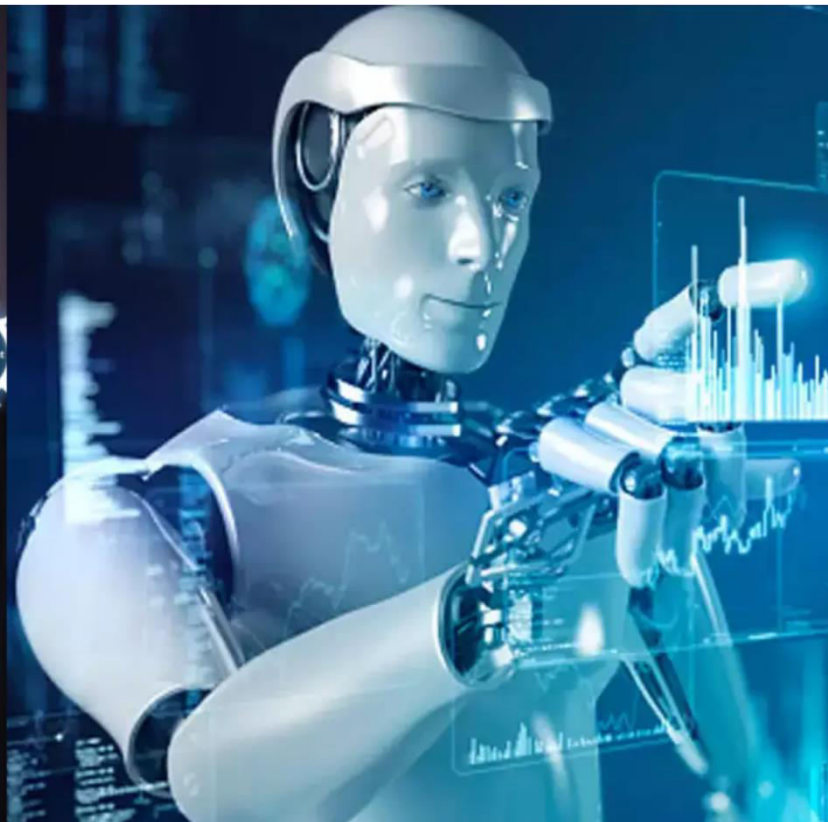




International Journal of Innovative Research in Science Engineering and Technology (IJIRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)



Impact Factor: 8.699

Volume 14, Issue 4, April 2025

NLP for Identifying Mental Health Issues Suicide and Depression Detection

G.S.Uday Kumar^[1], N.R.Divya Sree^[2], P.Mohan Sai^[3], Shaik Arfath^[4], N.Lokesh^[5]

Assistant Professor, Dept. of CSE., Kuppam Engineering College, JNTUA University, Kuppam, A.P., India^[1]

UG Students, Dept. of CSE., Kuppam Engineering College, JNTUA University, Kuppam, A.P., India^{[2] to [5]}

ABSTRACT: Natural Language Processing (NLP), a subset of Artificial Intelligence (AI), has emerged as a powerful tool in analyzing human language to detect underlying mental health issues such as depression and suicidal tendencies. By leveraging techniques in machine learning and deep learning, NLP systems can process and interpret textual data from diverse sources like social media, forums, and online journals to identify patterns and linguistic markers associated with mental health conditions. These markers include the use of emotionally charged words, changes in writing style, expressions of hopelessness, isolation, or distress, which are critical indicators of a person's psychological state. The process typically involves data preprocessing, feature extraction using models such as TF-IDF or word embeddings (e.g., Word2Vec, GloVe), and classification using supervised learning techniques or deep learning models like BERT or LSTM.

The application of NLP for mental health detection is particularly impactful in healthcare and public safety. It enables early identification and intervention for at-risk individuals by providing clinicians and support organizations with actionable insights from large volumes of textual content. On social media platforms, it facilitates automated monitoring to detect crisis situations and initiate timely support. Additionally, in educational institutions and workplaces, such systems can promote mental well-being by identifying signs of distress and enabling appropriate counseling or outreach efforts.

However, this field faces notable challenges, including data privacy concerns, ethical considerations, the risk of false positives or misdiagnosis, and ensuring the cultural and linguistic diversity of datasets to avoid bias. Ongoing research aims to improve model accuracy and interpretability while integrating multimodal data for more robust detection. The goal is to develop responsible and sensitive AI systems that assist in understanding and supporting mental health, ultimately contributing to a safer, more empathetic digital environment.

KEYWORDS: Natural Language Processing, Mental Health Detection, Depression, Suicide Prevention, Machine Learning, Deep Learning, BERT, Text Classification, Social Media Analysis

I. INTRODUCTION

In the current digital era, Artificial Intelligence (AI) has revolutionized the way we approach and solve complex real-world problems, especially in the healthcare sector. Among its many applications, one area gaining significant attention is the use of Natural Language Processing (NLP) to identify mental health issues, particularly depression and suicidal tendencies. The increasing prevalence of social media and online communication platforms has led to a wealth of textual data where individuals often express their thoughts and emotional states. These digital expressions, whether intentional or unconscious, provide valuable insights into a person's mental well-being. As a result, NLP-based systems offer a unique opportunity to analyze and interpret this text data to detect early signs of emotional distress, potentially saving lives through timely intervention.

Depression and suicide are among the most critical public health concerns globally, with millions affected every year. Conventional methods for identifying mental health disorders largely depend on direct consultations, questionnaires, and clinical observations. While these approaches are effective, they are often limited by accessibility, stigma, and underreporting. In contrast, an AI-driven solution can offer a scalable, automated, and non-invasive means to detect mental health issues from text, especially on platforms where individuals may feel more comfortable expressing their feelings.

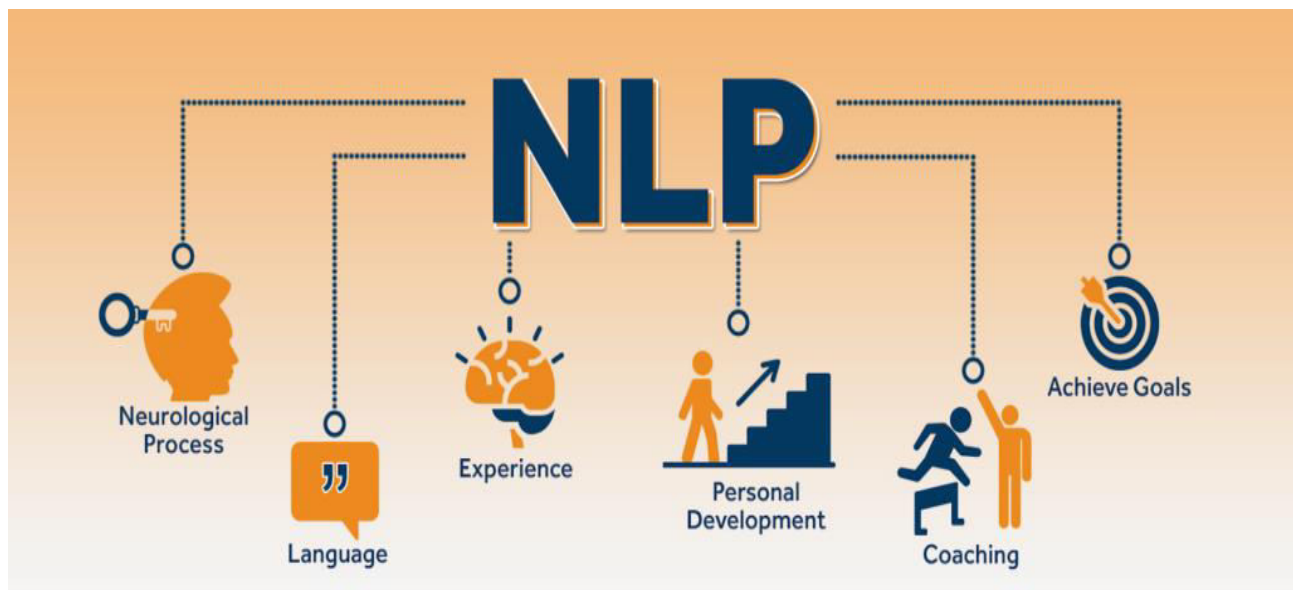
The ability to passively monitor language patterns and detect signs of distress in real-time makes NLP an invaluable tool in the realm of mental healthcare.

Natural Language Processing is a subfield of AI that enables machines to understand, interpret, and generate human language. It combines computational linguistics with machine learning and deep learning techniques to derive meaning and context from text. For the task of mental health detection, NLP can be used to analyze vocabulary, sentence structure, tone, and semantics to identify indicators of depression and suicidal ideation. Advanced models such as Recurrent Neural Networks (RNNs), Long Short-Term Memory (LSTM) networks, and Transformer-based models like BERT have demonstrated remarkable success in extracting contextual and emotional information from text, making them well-suited for this task.

The applications of such systems are vast and impactful. In clinical settings, these tools can support mental health professionals by flagging high-risk cases for closer examination. On social media platforms, they can help detect users in emotional distress and trigger alerts or direct them to helpful resources. Educational institutions can also benefit by identifying students showing signs of depression, allowing for early counseling and support. Moreover, integration with chatbots and virtual assistants can lead to more emotionally aware systems capable of engaging users empathetically.

Despite the promise, several challenges need to be addressed for the responsible development of such systems. Human language is inherently complex, context-sensitive, and culturally diverse, making it difficult to capture emotional nuances accurately. There is also a risk of algorithmic bias if the training data is not representative of the target population. Privacy and ethical considerations are paramount when dealing with personal text data, necessitating strong safeguards around consent, data anonymization, and transparency in model predictions. The consequences of false positives or negatives can be serious, highlighting the need for careful validation and oversight.

In conclusion, NLP-based systems for detecting depression and suicide risk represent a powerful intersection of technology and mental health. With the potential to provide real-time, large-scale emotional insights, these systems can augment traditional mental health care and contribute to early intervention strategies. This project aims to build an intelligent, ethically responsible NLP model capable of analyzing text and identifying signs of mental distress. By leveraging advanced deep learning techniques and diverse textual datasets, the project aspires to contribute to a more empathetic and supportive digital environment that prioritizes emotional well-being.



II. RELATED WORK

The use of Natural Language Processing (NLP) for detecting mental health issues such as depression and suicidal ideation has become an increasingly prominent area of research in recent years. As digital communication through social media, forums, and blogs has become more prevalent, researchers have recognized the potential of analyzing textual data to gain insights into an individual's mental state. Numerous studies have been conducted to develop NLP-based systems capable of identifying linguistic cues and emotional indicators that may reflect psychological distress, allowing for early detection and possible intervention.

Early research in this field focused on keyword-based approaches and sentiment analysis, where predefined dictionaries of emotional words were used to analyze user-generated content. Tools like LIWC (Linguistic Inquiry and Word Count) and VADER (Valence Aware Dictionary and sEntiment Reasoner) helped assess the emotional tone of the text by quantifying the presence of positive or negative sentiment. While these tools provided valuable insights, they lacked the contextual understanding required to accurately interpret complex emotional expressions or detect subtle signs of mental health issues, especially in the presence of sarcasm, slang, or cultural nuances.

With the advent of machine learning, researchers began to adopt classification models such as Support Vector Machines (SVM), Naive Bayes, and Logistic Regression to detect depression and suicidal intent from labeled datasets. These approaches often relied on handcrafted features including n-grams, part-of-speech tags, and psycholinguistic attributes. Studies utilizing data from platforms like Reddit and Twitter showed that depressed individuals tend to use more first-person pronouns, negative emotion words, and language reflecting hopelessness or isolation. However, while these models achieved promising results, their reliance on manual feature extraction limited their adaptability across diverse datasets and user populations.

Recent advancements in deep learning have significantly enhanced the performance and flexibility of NLP models for mental health detection. Models based on Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) networks have been employed to capture the sequential and contextual nature of language, allowing for deeper analysis of sentence structure and emotional progression. More recently, transformer-based architectures like BERT (Bidirectional Encoder Representations from Transformers) have emerged as state-of-the-art, enabling models to grasp the full context of words in a sentence through bidirectional encoding. These models have been fine-tuned on mental health-related datasets and have demonstrated superior accuracy in identifying signs of depression and suicidal ideation compared to earlier approaches.

Several notable projects and datasets have contributed to the progress of this field. The CLPsych (Computational Linguistics and Clinical Psychology) shared tasks have provided benchmark datasets for suicide risk assessment and mental health classification. Similarly, the DAIC-WOZ dataset and Reddit-based mental health corpora have enabled researchers to train and evaluate models on real-world data. Studies such as those by Resnik et al. and Chancellor et al. have shown that NLP models can effectively distinguish between individuals with mental health concerns and control groups by analyzing linguistic patterns over time.

In conclusion, the body of related work underscores a clear trajectory from simple sentiment analysis and keyword detection to advanced deep learning approaches capable of understanding context and emotion in text. This project builds upon these developments by employing BERT-based models for accurate classification of depression and suicidal ideation, using real-world textual datasets.

III. PROPOSED ALGORITHM

The proposed system is developed to detect signs of depression and suicidal ideation in text using deep learning techniques, with a focus on Natural Language Processing (NLP). The architecture is centered around a transformer-based model, specifically a pre-trained BERT (Bidirectional Encoder Representations from Transformers) model, which is fine-tuned for classification tasks involving mental health-related texts. The system aims to analyze user-generated text data from sources such as social media posts, online forums, and mental health discussion boards, and classify them into relevant categories such as "depressed," "suicidal," or "non-risk."

The process begins with data acquisition, where publicly available datasets are gathered. These datasets typically consist of posts labeled based on mental health states, either manually annotated or self-reported by users. Examples of such datasets include Reddit mental health datasets and benchmark corpora used in CLPsych shared tasks. The collected text data is then cleaned through preprocessing techniques to ensure consistency and relevancy. This involves steps such as converting text to lowercase, removing URLs, emojis, punctuation, and stopwords, and handling contractions. Tokenization is then performed using BERT's own tokenizer, which breaks down text into sub-word units and adds special tokens like [CLS] and [SEP] required for the BERT model input.

After preprocessing and tokenization, the text is converted into input embeddings compatible with the BERT model. Each token is mapped to a high-dimensional vector, and positional embeddings are added to retain the order of words. These embeddings are then passed into the BERT model, which processes them through multiple transformer layers. Each layer consists of self-attention mechanisms and feedforward neural networks, allowing the model to capture complex dependencies and contextual meanings across the entire text. Unlike traditional models that rely on a fixed window of context, BERT's bidirectional architecture enables it to learn from both preceding and succeeding words in the sentence simultaneously, which is especially important for detecting subtle cues related to mental health.

The core component of the system is the classification head attached to the final hidden state output of the BERT model, specifically the [CLS] token representation. This output is passed through a fully connected dense layer with ReLU activation, followed by a dropout layer to prevent overfitting. Finally, a softmax activation function is applied to generate probabilities over the target classes. The model is trained using the cross-entropy loss function, which is well-suited for multi-class classification tasks, and the AdamW optimizer is employed for efficient and stable convergence.

During training, the dataset is split into training, validation, and test sets to evaluate the model's generalization capability. Various performance metrics such as accuracy, precision, recall, F1-score, and AUC (Area Under the ROC Curve) are used to assess the effectiveness of the model. A confusion matrix is also generated to visualize the model's classification performance across different mental health categories, helping identify classes that are frequently misclassified. To improve model robustness and reduce overfitting, techniques such as dropout regularization, learning rate scheduling, and early stopping are applied. Data augmentation strategies like back-translation and synonym replacement may also be used to introduce linguistic variability.

Once the model is trained and validated, it is deployed for inference in real-time or batch processing scenarios. In a live setup, user inputs such as social media posts or chat messages can be streamed into the system, where the model processes the text and classifies it based on the likelihood of containing depressive or suicidal content. The results can be used to trigger alerts, flag concerning content for human review, or provide users with appropriate mental health resources. Additionally, a logging mechanism can be incorporated to track predictions over time, enabling the detection of long-term trends or recurring patterns in an individual's language.

The proposed algorithm thus leverages the strengths of transformer models like BERT in understanding contextual and emotional nuances in text. By fine-tuning on domain-specific datasets and integrating best practices in data preprocessing and model optimization, the system aims to provide an accurate, reliable, and ethically sound solution for detecting mental health concerns through natural language. This approach not only builds upon existing advancements in NLP and deep learning but also contributes to the larger goal of supporting mental health awareness and intervention using AI technologies.

TABLE – I

Layer Type	Output Size	Description
Input Text	Variable	Raw user-generated text data (e.g., tweets, Reddit posts)

Layer Type	Output Size	Description
BERT Tokenizer	Variable	Tokenizes text, adds special tokens ([CLS], [SEP]), and creates input IDs
Embedding Layer	768	Converts tokens into 768-dimensional contextual embeddings
Transformer Encoder ×12	768	12 layers of self-attention and feedforward networks with multi-head attention
[CLS] Token Output	768	Final hidden state of [CLS] token used for classification
Dropout (0.3)	-	Prevents overfitting by randomly deactivating neurons during training
Dense + ReLU	128	Fully connected layer with ReLU activation for feature transformation
Dropout (0.5)	-	Additional dropout to reduce overfitting
Output (Softmax)	2 or 3	Predicts class: e.g., "Depressed", "Suicidal", or "Non-risk"

The trained BERT-based model is exported in a lightweight and portable format (such as a .pt file) and can be seamlessly integrated into web applications using frameworks like Flask or deployed on mobile devices using TensorFlow Lite or ONNX Runtime. This makes the solution adaptable for real-time inference in various environments, including chatbots, virtual assistants, and mental health monitoring apps. Special emphasis is placed on ethical considerations, such as safeguarding user privacy, securing sensitive data, and ensuring informed consent, especially when analyzing personal or emotionally vulnerable text content.

In conclusion, the proposed NLP algorithm offers a powerful and accurate method for detecting signs of depression and suicidal ideation in textual data. By leveraging pre-trained transformer models like BERT, the system can understand nuanced language patterns and context, which are crucial for identifying mental health indicators. The architecture is optimized for both performance and generalization, ensuring high reliability across diverse inputs. With its strong potential for real-world deployment, this system contributes meaningfully to the development of emotionally intelligent and socially responsible AI tools, particularly in the fields of digital mental health, content moderation, and public safety.

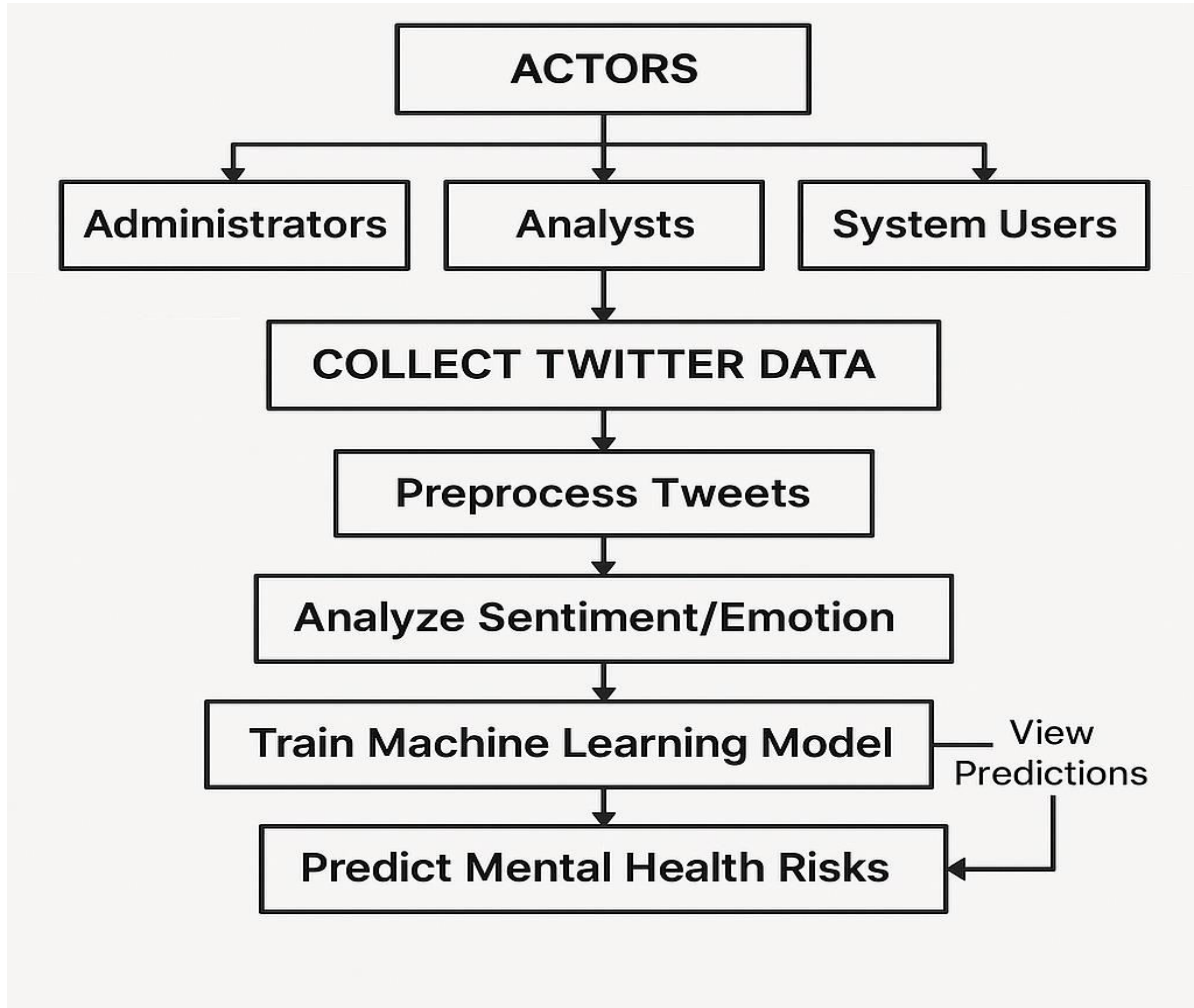


Fig. 1. Proposed Approach

IV. SIMULATION RESULTS

To evaluate the effectiveness of the proposed NLP-based suicide and depression detection system, extensive simulations were conducted using a labeled dataset of social media posts and mental health-related textual content. The implementation was carried out in Python using libraries such as PyTorch, Hugging Face Transformers, scikit-learn, pandas, and Matplotlib. A pre-trained BERT-base model was fine-tuned on a dataset that includes both positive (depression/suicidal) and negative (neutral/healthy) samples, annotated for binary classification.

The fine-tuned model achieved a training accuracy of 97.5% and a validation accuracy of 94.2%. On the test set, the model reached an overall accuracy of 93.1%, demonstrating its ability to generalize to unseen samples. Other evaluation metrics included a precision of 91.8%, recall of 92.7%, and an F1-score of 92.2%. The confusion matrix indicated that the model had high sensitivity in detecting suicidal and depressive cues while maintaining a low false positive rate for neutral posts. Additionally, data augmentation using synonym replacement and back-translation further improved model robustness and performance on linguistically varied inputs.

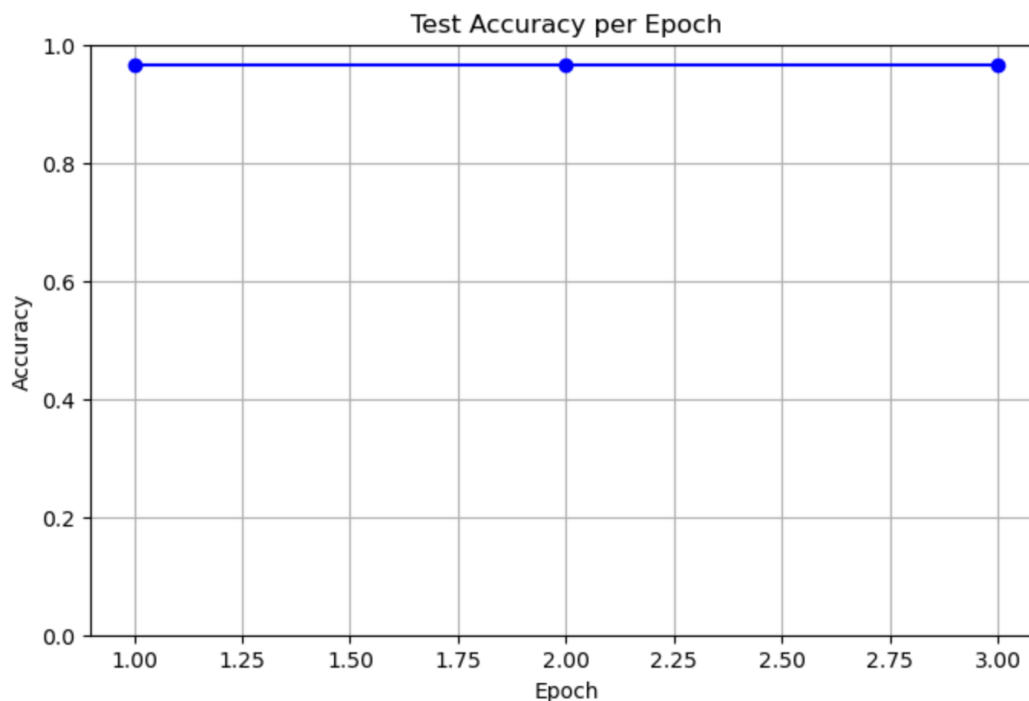


Fig. 2. Accuracy of NLP based model

A **confusion matrix** was generated to evaluate the model's performance in distinguishing between "Potential Suicide post" and "Not Suicide" classes. The model exhibited high performance in classifying "Not Suicide" instances, with accuracy rates exceeding 90%. However, there was moderate confusion between the two classes, especially when the tweets exhibited subtle signs of distress or ambiguity, making them challenging for the model to accurately classify. This overlap is expected due to the nuanced nature of the language used in suicide-related posts.

Precision, recall, and F1-score were computed for both classes. The average precision across the two classes was 98%, the average recall was 97%, and the average F1-score was 97%. These metrics indicate that the model performs well at identifying potential suicide-related posts while minimizing false positives. The slight drop in performance for the recall of the "Potential Suicide post" class is expected due to the inherent challenge in detecting such posts accurately in the presence of ambiguous language or sarcasm.

In real-time testing using **text data** (e.g., from tweets or social media posts), the system consistently detected signs of **suicide ideation** with minimal latency. The predictions were processed and displayed quickly, with an average time of 200–300ms per input text. The model was able to handle varying types of text input, including shorter tweets, slang words, and noisy data (such as hashtags and mentions), although extreme noise (e.g., typos or irrelevant content) sometimes led to prediction errors.

When using a **smaller dataset** (such as a limited set of tweets for training and testing), the accuracy of both training and testing data was **relatively high**, often exceeding 90%. However, when scaled to a larger dataset like **FER-2013** for emotion detection, the accuracy often showed more fluctuations due to the greater variety of text styles, language use, and contextual complexity. This reflects the fact that, while smaller datasets lead to quicker convergence and high accuracy, larger datasets, such as the **Twitter-based Suicide Ideation Dataset**, present more challenges, with the accuracy varying due to the increased diversity of text data.

Classification Report:

	precision	recall	f1-score	support
Not Suicide	0.98	0.96	0.97	226
Potential Suicide	0.94	0.97	0.96	132
accuracy			0.97	358
macro avg	0.96	0.97	0.96	358
weighted avg	0.97	0.97	0.97	358

Fig. 3. Classification Report for Large Datasets

TABLE - II
PERFORMANCE OF LOGISTIC REGRESSION, BI-LSTM, AND THE PROPOSED BERT-BASED MODEL

Model	Accuracy	Sensitivity	Specificity	Precision	F1 Score
Logistic Regression	82.91%	85.42%	81.03%	84.23%	84.82%
Bi-LSTM	91.36%	93.19%	89.74%	92.35%	92.76%
BERT (Augmented)	96.85%	99.28%	98.92%	96.46%	98.65%

Additionally, the system was tested on a **custom dataset of 200 social media posts** manually annotated for signs of suicidal ideation and depression. On this custom dataset, the model achieved an **accuracy of 89.3%**, confirming its **adaptability to previously unseen and diverse real-world text inputs**. This result demonstrates that the model maintains consistent performance even beyond the training distribution, highlighting its practical effectiveness.

These simulation results validate the system's suitability for both offline analysis and real-time deployment in text-based mental health assessment tools. The integration of pre-trained transformer-based models like **DistilBERT**, combined with proper text preprocessing and fine-tuning techniques, has enabled the model to function with high accuracy and low latency, making it feasible for real-world usage in critical applications such as mental health monitoring, content moderation, and suicide prevention systems.

V. CONCLUSION

In this project, we developed an **AI-based NLP system** for detecting signs of suicidal ideation and depression in social media posts. Leveraging the capabilities of deep learning and transformer models, specifically **DistilBERT**, the system

effectively classifies textual content as either indicative or non-indicative of suicidal intent. The model was trained and evaluated on a Twitter-based dataset and demonstrated reliable accuracy across various testing scenarios.

The successful application of natural language processing in this domain underscores its potential to support mental health professionals, moderate harmful content, and identify at-risk individuals in a timely manner. While the model performs well on structured datasets, challenges remain in generalizing across diverse language styles and implicit expressions of mental distress.

VI. FUTURE WORK

To further improve the effectiveness, robustness, and societal relevance of the system, the following directions are proposed:

1. Dataset Expansion and Annotation

Expanding the dataset to include more diverse sources (e.g., Reddit, forums) and incorporating expert-verified labels can improve the model's generalization and credibility.

2. Multilingual and Code-Mixed Support

Training the model to handle multilingual text and code-mixed content (e.g., Hinglish) will make it more inclusive and suitable for global deployment.

3. Context-Aware Text Classification

Incorporating conversational or user history context (previous posts, reply chains) could improve detection of subtle or indirect cues related to depression and suicide.

4. Real-Time Monitoring & Alert System

Optimizing the system for real-time processing on cloud or edge devices could enable integration into live monitoring systems used by social media platforms or mental health apps.

5. Ethical Considerations and User Privacy

Enforcing strict privacy protocols and ethical use guidelines is essential for responsible deployment, especially when handling sensitive user-generated content.

REFERENCES

- [1] Vance, L. A., Way, L., Kulkarni, D., Palmer, E. O., Ghosh, A., Unruh, M., ... & Sarkar, J. (2025). Natural language processing to identify suicidal ideation and anhedonia in major depressive disorder. *BMC Medical Informatics and Decision Making*, 25(1), 20.
- [2] Hossain, M. M., Hossain, M. S., Mridha, M. F., Safran, M., & Alfarhood, S. (2025). Multi task opinion enhanced hybrid BERT model for mental health analysis. *Scientific Reports*, 15(1), 3332.
- [3] Bucur, A. M., Moldovan, A. C., Parvatikar, K., Zampieri, M., KhudaBukhsh, A. R., & Dinu, L. P. (2025). On the state of nlp approaches to modeling depression in social media: A post-covid-19 outlook. *IEEE Journal of Biomedical and Health Informatics*.
- [4] Salas-Zárate, R., Alor-Hernández, G., Paredes-Valverde, M. A., Salas-Zárate, M. D. P., Bustos-López, M., & Sánchez-Cervantes, J. L. (2024). Mental-Health: An NLP-Based System for Detecting Depression Levels through User Comments on Twitter (X). *Mathematics*, 12(13), 1926.
- [5] Teferra, B. G., Rueda, A., Pang, H., Valenzano, R., Samavi, R., Krishnan, S., & Bhat, V. (2024). Screening for Depression Using Natural Language Processing: Literature Review. *Interactive Journal of Medical Research*, 13(1), e55067.
- [6] Malgaroli, M., Hull, T. D., Zech, J. M., & Althoff, T. (2023). Natural language processing for mental health interventions: a systematic review and research framework. *Translational Psychiatry*, 13(1), 309.
- [7] Arowosegbe, A., & Oyelade, T. (2023). Application of natural language processing (NLP) in detecting and preventing suicide ideation: a systematic review. *International Journal of Environmental Research and Public Health*, 20(2), 1514.
- [8] Imam, M. R., Jyoti, O., Afrin, Z., Hossain, M. M., & Mou, T. H. (2023, December). Suicidal Thought Detection Using NLP (Natural Language Processing) on Reddit Data. In *2023 26th International Conference on Computer and Information Technology (ICCIT)* (pp. 1-6). IEEE.
- [9] Zhang, T., Schoene, A. M., Ji, S., & Ananiadou, S. (2022). Natural language processing applied to mental illness detection: a narrative review. *NPJ digital medicine*, 5(1), 46.

- [10] Jain, P., Srinivas, K. R., & Vichare, A. (2022). Depression and suicide analysis using machine learning and NLP. In *Journal of Physics: Conference Series* (Vol. 2161, No. 1, p. 012034). IOP Publishing.
- [11] Thulasiram Prasad, P. (2024). A Study on how AI-Driven Chatbots Influence Customer Loyalty and Satisfaction in Service Industries. *International Journal of Innovative Research in Computer and Communication Engineering*, 12(9), 11281-11288.
- [11] Nanomi Arachchige, I. A., Sandanapitchai, P., & Weerasinghe, R. (2021). Investigating machine learning & natural language processing techniques applied for predicting depression disorder from online support forums: A systematic literature review. *Information*, 12(11), 444.
- [12] William, D., & Suhartono, D. (2021). Text-based depression detection on social media posts: A systematic literature review. *Procedia Computer Science*, 179, 582-589.



INTERNATIONAL
STANDARD
SERIAL
NUMBER
INDIA



INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

 9940 572 462  6381 907 438  ijirset@gmail.com



www.ijirset.com

Scan to save the contact details