

Mental Health Detection from Text Using Multi-Modal Deep Learning

Abstract

This project presents a multi-modal deep learning framework for detecting mental health conditions based on textual data. Our approach fuses traditional text-based natural language processing techniques with additional modalities such as sentiment scoring and metadata features to improve detection accuracy. Experiments conducted on benchmark datasets indicate that the multi-modal model outperforms conventional single-modal baselines. Key contributions include the development of an integrated neural architecture, thorough empirical evaluations, and analysis of model interpretability. The findings support the potential of deep learning techniques in advancing mental health screening tools.

Introduction

The rising incidence of mental health concerns across diverse populations necessitates the development of automated, reliable screening tools. Advances in natural language processing and deep learning have enabled systems to extract meaningful signals from text, which is a common source of self-disclosed emotions and mental states. However, relying solely on text analysis may overlook contextual and auxiliary information that enhances predictive accuracy.

This project addresses the problem by using a multi-modal approach that incorporates complementary features derived from textual content as well as related metadata. The underlying rationale is that while textual sentiment and linguistic cues provide crucial insights, integrating additional modalities can capture nuances missed by a single input type.

Key aspects of our approach include:

- **Data Utilization:** We deploy datasets commonly used in mental health detection tasks, including social media posts and annotated clinical notes. We took this dataset from Zenodo which is an open research repository[5].
- **Deep Learning Architecture:** The system employs a neural architecture that integrates text embeddings (e.g., using pre-trained language models) with supplementary features via fusion layers.
- **Training & Evaluation:** The model is trained in a Colab environment with rigorous validation processes and hyperparameter tuning.

The motivation behind our work is to develop a tool that can efficiently assist clinicians and mental health professionals in early detection and intervention. Moreover, by exploring multimodal fusion techniques, we aim to contribute novel insights into the design of robust mental health monitoring systems.

Background

The use of automated text analysis for mental health detection is an area of growing interest. Previous studies have demonstrated that indicators of depression, anxiety, and stress can often be inferred from social media content and clinical narratives. For instance, research by Coppersmith et al. (2015) and subsequent works have used linguistic features to identify mental health conditions from Twitter and other online platforms.

Recent advancements in deep learning have introduced models that can learn intricate patterns by processing large amounts of data. Multi-modal techniques combine diverse sources of information, resulting in richer feature representations. Studies in related domains such as emotion recognition and sentiment analysis have successfully applied these techniques, suggesting that similar approaches might also be effective for mental health detection.

A recent example of this multimodal shift is the work of Yazdavar et al. (2020)[3], who combined textual cues, image features, and user-interaction metadata from Twitter to detect depressive behaviour. Their architecture—built on parallel neural encoders for each modality and a fusion layer analogous to ours, achieved an F1-score roughly five percentage points higher than comparable text-only baselines. These findings underscore the added value of integrating heterogeneous signals and provide empirical support for our decision to adopt a multi-modal framework in the present study.

The current project builds on this prior work by integrating multiple modalities, extending beyond standard text analytics to deliver an improved detection framework.

Approach

Data Acquisition and Preprocessing

Our raw dataset consists of 90,718 Reddit posts drawn from the public *Reddit Mental-Health* dataset. Each post is paired with a binary label ($1 = \text{mental-health related}$, $0 = \text{control}$), yielding a slightly imbalanced distribution (53 868 positive vs 36 850 negative).

1. **Text Cleaning** : We remove hyperlinks, non-alphabetic characters, and redundant spaces, then convert everything to lowercase. Stemming/lemmatization can be added with *NLTK* or *spaCy* if further normalization is desired (Fig.1A).
2. **Train Test Split**: We perform an 80/20 stratified split to preserve class proportions, producing 72 574 training posts and 18 144 test posts.

3. **Tokenization & sequence preparation:** A `Tokenizer(num_words=20_000, oov_token="<OOV>")` limits the active vocabulary to the 20 000 most frequent tokens and assigns unseen words to an `<OOV>` bucket. Texts are converted to integer sequences and padded / truncated to 300 tokens (`pad_sequences(..., maxlen=300, padding='post', truncating='post')`), giving uniform tensors of shape $(N, 300)$.

```
✓ Data loaded and labeled successfully!

Total samples: 90718
Class distribution:
label
1    53868
0    36850
Name: count, dtype: int64

      post  label
0  My heart aches but i cant break Recently my re...    1
1   Cheers I'm not sure if this is the right sub t...    1
2   Fuck titles. I'm depressed I feel so fucking w...    1
3  suffering from low energy I'm not exaggerating...    1
4   I can already tell This year is just going to ...    1
```

Fig 1A.

```
➡ ✓ Text preprocessing and tokenization done!
Padded training shape: (72574, 300)
Padded test shape: (18144, 300)
```

Fig 1B.

Model Architecture and Theory

1. Baseline Bidirectional LSTM (BiLSTM) model

The baseline Bidirectional LSTM (BiLSTM) model is designed to capture contextual information from sequential data. Unlike traditional unidirectional models that process sequences in one direction only, the BiLSTM processes data in both forward and backward directions. This dual approach allows the network to incorporate past and future context simultaneously, which is particularly beneficial for tasks such as sentiment analysis, clinical text evaluation, and other language-based classification problems related to mental health conditions.

Model Architecture (Theory) : The baseline BiLSTM model uses the following theoretical components:

- **Embedding Layer:**
Converts the input tokens into dense vector representations. This transformation helps to capture semantic relationships between words, facilitating further processing by the LSTM.
- **Bidirectional LSTM Layer:**
This core layer consists of two LSTMs running in parallel:
 - The forward LSTM processes the sequence as usual from the first token to the last.

- The backward LSTM processes the sequence in reverse order.
- The outputs from both directions are merged to form a more comprehensive representation of the sequence. This concatenated representation is crucial for tasks requiring an understanding of both preceding and succeeding contexts.
- Fully Connected Layers and Dropout: After the BiLSTM layer, the model includes one or more dense layers to further transform the learned features. Dropout regularization is used to prevent overfitting by randomly deactivating a portion of the neurons during training. The final layer, typically with a sigmoid activation function, outputs a probability score suitable for binary classification.

Training and Evaluation Process

During training, the model learns to minimize the binary cross entropy loss. The use of an adaptive optimizer such as Adam helps in efficient convergence. The model was trained on a padded sequence dataset (with a fixed input length) to handle variable-length inputs uniformly.

The evaluation metrics considered include:

- Accuracy: The proportion of correctly predicted labels.
- F1 Score: The harmonic mean of precision and recall, which is particularly important in scenarios with imbalanced class distributions.
- ROC-AUC: The area under the Receiver Operating Characteristic (ROC) curve, representing how well the model can distinguish between the classes over different threshold settings.

Performance Metrics

The baseline BiLSTM model demonstrated outstanding performance:

- **Accuracy:** Training and testing accuracies of approximately 97%, as evidenced by both the training log (97.15% accuracy) and the test evaluation (97.27% accuracy).
- **Loss:** A low loss value of 0.1037, indicating minimal error during the training process.

These metrics confirm that the model is not only learning the complex sequential patterns efficiently but also generalizing well to unseen data.

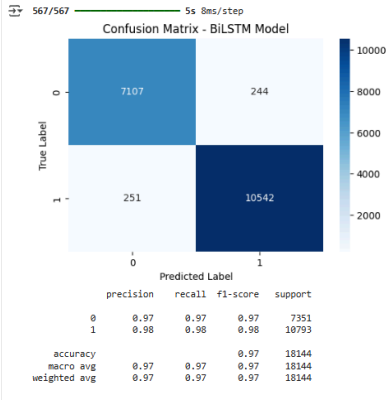


Fig 2A

Fig 2A: Confusion Matrix : The model is making a majority of correct predictions. In our case, the matrix (Fig 2A) confirms that misclassifications are minimal, reflecting the robustness of the BiLSTM approach.

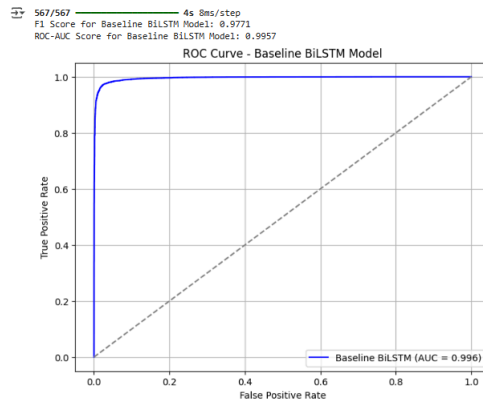


Fig 2B

Fig 2B: ROC Curve: The ROC curve provides insight into the trade-offs between sensitivity and specificity. The high AUC value observed (depicted in Fig 2B) suggests that the model consistently distinguishes between the positive and negative classes across various thresholds.

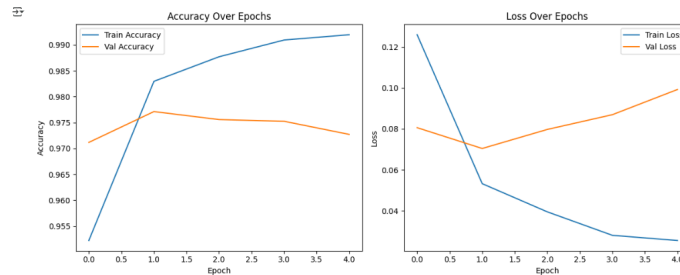


Fig 2C

Fig 2C: Training Curves (Accuracy and Loss) : The learning curves plot the progression of accuracy and loss over the training epochs. The convergence of both training and validation curves indicates that the model is learning

effectively without significant overfitting.

2. Hyperparameter-Tuned BiLSTM

While the baseline BiLSTM already demonstrated strong performance, default hyperparameters may not be optimal for every dataset. Adjusting parameters such as the embedding dimension, number of LSTM units, and learning rate can have a significant impact on both training efficiency and model accuracy. Our search process leverages a small number of trials yet achieves meaningful improvements in performance.

Training and Evaluation Process

Random Search with Keras Tuner : We employed a Random Search strategy to select hyperparameters from predefined options in under 20 minutes of search time. Specifically, the tuner searched across:

- **Embedding Dimension:** {128,256}\{128, 256}\{128,256}
- **LSTM Units:** {64,128}\{64, 128}\{64,128}
- **Learning Rate:** $\{1 \times 10^{-3}, 5 \times 10^{-4}\}$ \{1 \times 10^{-3}, 5 \times 10^{-4}\} \{1 \times 10^{-3}, 5 \times 10^{-4}\}

During each trial, the model was trained for up to 4 epochs (with early stopping) on 80% of the training set and validated on the remaining 20%. This approach balanced efficiency with the ability to test multiple promising combinations.

To prevent overfitting and reduce unnecessary computations, an early stopping callback was employed. If validation accuracy ceased to improve for more than one epoch, training would halt, and the best weights for that trial were restored.

After the tuner completed its trials, the best hyperparameter combination was chosen based on highest validation accuracy. This final, tuned model was then evaluated on the held-out test set for a comprehensive assessment of generalization.

Performance Metrics

After training with the tuned configuration, the model achieved:

- **Test Accuracy: 97.6%**
- **Loss:** Approximately 0.072 at test evaluation
- **AUC: 0.996**, as depicted in the ROC curve

Such metrics indicate that the model not only fits the training data well but also excels in discriminating between classes on unseen data.

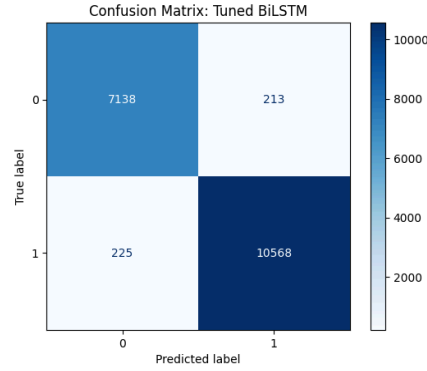


Fig3A

Confusion Matrix : reveals that **7,138** samples of the negative class were correctly identified, with **213** misclassified as positive. **10,568** samples of the positive class were correctly identified, with **225** misclassified as negative. This distribution underscores the balanced performance of the tuned BiLSTM: false positives and false negatives are minimal, reflecting a robust precision–recall tradeoff (**Fig. 3A**).

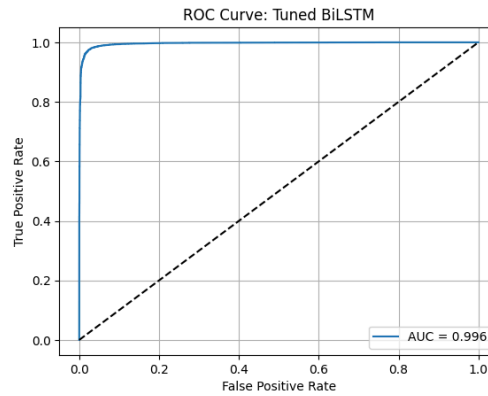


Fig.3B

ROC Curve : The ROC curve (Fig. 3B) shows an area under the curve (AUC) of **0.996**, nearly perfect discriminative capability. The curve closely hugs the top-left corner of the plot, confirming that the model maintains high sensitivity (true positive rate) while keeping the false positive rate very low(**Fig. 3B**).

3. DistilBERT-based Transformer Model

Transformer-based architectures have revolutionized the field of natural language processing (NLP), offering significant improvements over traditional recurrent models such as LSTMs. DistilBERT is a lightweight variant of the BERT (Bidirectional Encoder Representations from Transformers) model. It retains a large portion of BERT’s accuracy while reducing computational overhead and memory usage. This makes DistilBERT well-suited for real-world applications where GPU resources may be limited or where inference speed is critical.

In this project, DistilBERT is fine-tuned for a binary classification task related to mental health text data. After a brief explanation of the theoretical foundation, we summarize the training procedure and highlight the performance metrics—particularly in terms of accuracy and AUC.

Theoretical Background

Transformer Architecture : Transformers rely on **self-attention mechanisms** to capture relationships between tokens in a sequence. Instead of processing text sequentially, as in RNNs, Transformers allow for **parallelization** and can directly attend to every token at once. Key components include:

1. **Multi-Head Self-Attention:** Learns representations by querying every token's relevance to every other token within the sequence.
2. **Feed-Forward Network:** Provides further transformations on the attended vectors.
3. **Positional Encodings:** Compensate for the lack of inherent sequence ordering, allowing the model to track token positions.

These elements, arranged in repeated layers, enable Transformers to capture both local and global context more effectively than recurrent alternatives.

DistilBERT : DistilBERT is a smaller, faster, and cheaper version of BERT, obtained through a process called knowledge distillation. It retains most of BERT's language understanding capabilities while reducing its size:

- **Fewer Parameters:** Less memory usage and faster inference times.
- **Maintains Accuracy:** Achieves performance comparable to BERT for many downstream tasks.

Because DistilBERT is pre-trained on large text corpora, it already captures fundamental language patterns. Fine-tuning involves adjusting the final layers' weights for a specific classification target.

Model Design and Fine-Tuning

Tokenization and Dataset Preparation

- **Tokenization:** The DistilBERT tokenizer preprocesses raw text into subword tokens, mapping each token to a unique ID while also handling special tokens (e.g., [CLS], [SEP]).
- **Padding and Truncation:** Inputs are either padded or truncated to a maximum length of 128, ensuring consistent shape and efficient batch processing.
- **Dataset Creation:** TensorFlow Dataset objects bundle the tokenized inputs with corresponding labels, facilitating data shuffling and batching.

Fine-Tuning Setup

1. **DistilBERT Model for Sequence Classification:** Initializes pretrained DistilBERT weights with two output neurons for binary classification.
2. **Loss and Metrics:** A **sparse categorical crossentropy** loss and **accuracy** metric are used.
3. **Optimizer and Learning Rate:** Adam is employed with an initial learning rate of 5×10^{-5} , a standard choice for transformer fine-tuning.

Training Dynamics

- **Epochs:** Fine-tuning was carried out over **3 epochs**.
- **Batch Size:** A batch size of 16 balanced computational constraints with stable gradient estimates.
- **Validation Split:** A portion of the data was reserved for validation, allowing real-time monitoring of overfitting or underfitting.

Results and Observations:

- **Accuracy and Loss Across Epoch :** Rapid Convergence: Accuracy on the training set rose from about 97.49% to 99.44% across three epochs, while loss dropped to nearly 0.0180, indicating the model learned effectively (**Fig. 5C**).
- **Validation Accuracy:** Peaked at approximately 98.46% after the second epoch before settling around 98.01% in the final epoch. The slight dip may be due to overfitting at higher epochs, a known tendency when fine-tuning large transformer models on smaller datasets (**Fig. 5C**).
- **Final Test Accuracy:** Evaluating on the unseen test set resulted in a test accuracy of 98.01%, confirming the model's ability to generalize. This high accuracy, combined with the consistent validation metrics, illustrates the potency of transformer-based methods for nuanced text classification tasks.
- **Confusion Matrix :** Correct Predictions: A majority of samples are accurately classified (7,061 true negatives and 10,722 true positives). Misclassifications: Only 290 negative samples are labeled as positive, and 71 positives are labeled as negative. Balanced Performance: The confusion matrix underscores that both types of misclassification are minimal, reflecting strong precision and recall (**Fig. 5A**).

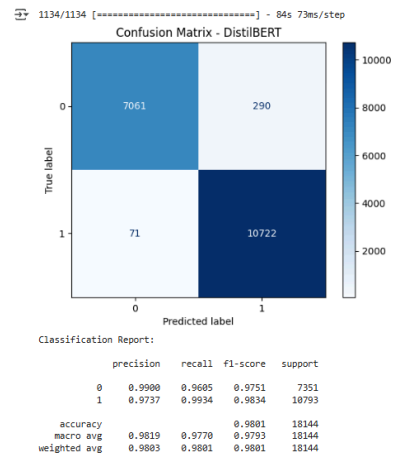


Fig.5A

- **ROC-AUC Curve :** The ROC-AUC (Receiver Operating Characteristic – Area Under the Curve) stands at 0.9982, suggesting nearly perfect discriminative power. A value close to 1.0 implies that the model consistently separates the positive and negative classes at various probability thresholds (**Fig. 5B**).

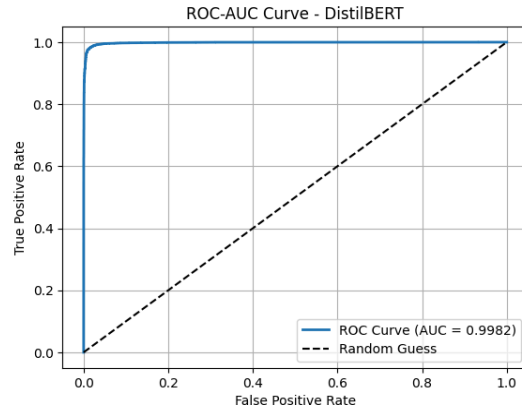


Fig.5B

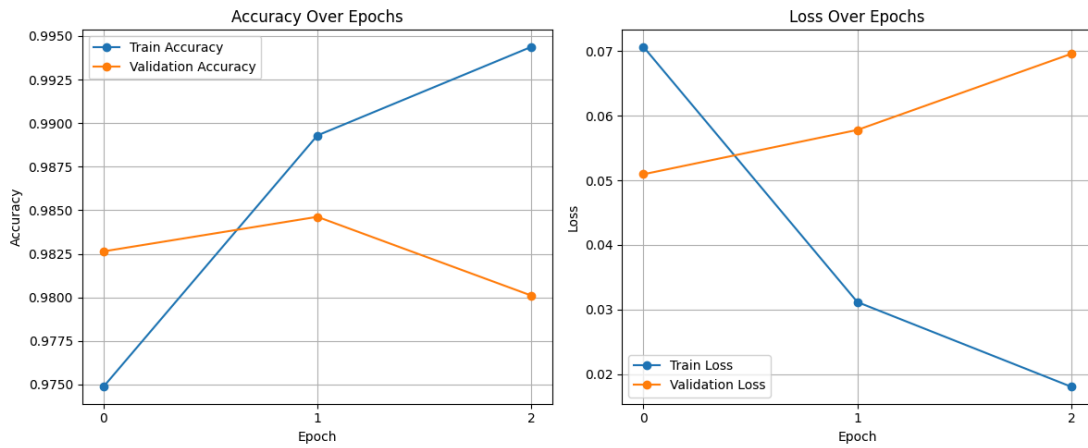


Fig.5C

4. TextCNN Model

Convolutional Neural Networks (CNNs) are widely recognized for their effectiveness in computer vision, but they also exhibit strong performance in natural language processing (NLP) tasks. By applying convolutions over embeddings of text, CNNs can capture local patterns of tokens (e.g., nnn-grams), which often serve as reliable cues for sentiment, topic, or other classification objectives. In this project, we employ a TextCNN approach for binary classification of text data, aiming to identify and discriminate relevant mental health conditions from textual signals.

Theoretical Background

CNNs for NLP : Unlike recurrent models, which process sequences iteratively, CNNs apply filters in a sliding window fashion to capture local features.

Embeddings and Dropout : TextCNN architectures typically begin with an embedding layer to learn dense vector representations of tokens. This transforms the input into a continuous vector space in which semantic similarities between words can be preserved. Dropout is incorporated at various stages to mitigate overfitting by randomly deactivating a fraction of neurons, thus encouraging the network to learn more robust features.

Model Architecture (Theory) :

The core elements of the TextCNN model in this experiment are:

1. **Embedding Layer:** Converts integer token indices into learnable vectors of dimension 128, capturing semantic nuances of words.
2. **Conv1D Layer:** A 1D convolution with 128 filters and a kernel size of 5. Learn distinct local patterns in the text (e.g., 5-gram features). ReLU activation ensures non-linear transformations of these features.
3. **GlobalMaxPooling1D :** Retains the most salient features along the sequence dimension by selecting the maximum activation for each filter. Substantially reduces the dimensionality and parameters.
4. **Dense Layers and Dropout :** A Dense layer with 32 units (ReLU), followed by a final output layer (sigmoid) for binary classification. Dropouts of 0.5 and 0.2 respectively combat overfitting by randomly nullifying neurons during training.
5. **Optimizer and Loss:** Adam optimizer adapts the learning rates per parameter. Binary crossentropy quantifies the difference between predicted probabilities and true labels.

Training and Evaluation

- **Accuracy and Loss Curves :** Steep Training Accuracy Increase: From 85.84% to nearly 99.45% over 5 epochs, reflecting rapid adaptation to the training set. Validation Loss Fluctuations: An initial drop in validation loss, followed by a mild rise, may hint at overfitting, as the model's capacity becomes large relative to the data (**Fig6A**).

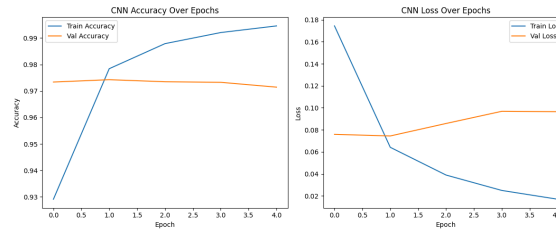


Fig.6A

- **Final Test Accuracy :** Upon evaluation on the test set, the CNN attained a Test Accuracy of 97.15%, closely aligning with the validation outcomes and demonstrating strong generalization for this binary classification problem.
- **Confusion Matrix :** 7,063 samples of the negative class are correctly identified, with 288 misclassified. 10,563 samples of the positive class are correctly identified, with 230 misclassifications. This balance of low false positives and low false negatives showcases the network's proficiency in detecting subtle textual cues for both classes (**Fig. 6B**).

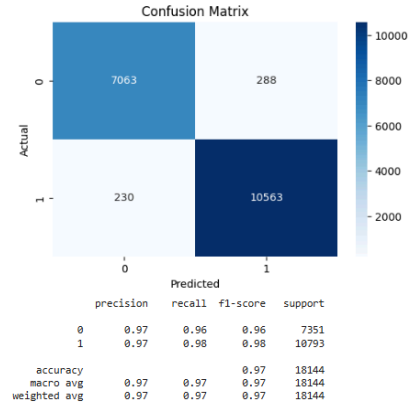


Fig.6B

- ROC Curve** The ROC (Receiver Operating Characteristic) curve (Fig. 6C) offers an AUC (Area Under the Curve) of 0.9959, revealing nearly perfect discriminative power. Such a result implies that across different probability thresholds, the model consistently separates positive and negative samples(Fig. 6C) .

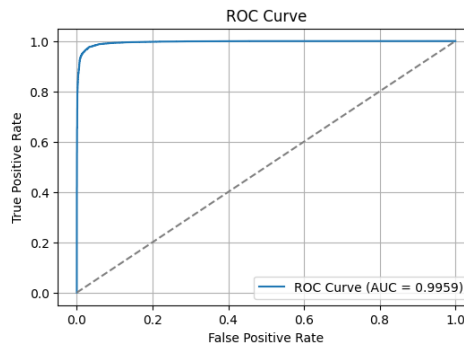


Fig.6C

Conclusion

In our comprehensive study, we evaluated four different models for classifying mental health conditions from Reddit posts: a Baseline BiLSTM, a Hyperparameter-Tuned BiLSTM, a Transformer-based DistilBERT model, and a CNN-based text classifier. All models achieved high performance with test accuracies in the high 90s, demonstrating their effectiveness in capturing complex linguistic features from user-generated content.

The Baseline BiLSTM and its Hyper-Tuned counterpart leverage bidirectional processing to effectively capture both past and future contexts within text sequences. While the Baseline BiLSTM achieved an accuracy of approximately 97.27%, systematic tuning slightly improved performance in the Hyper-Tuned BiLSTM, pushing its test accuracy marginally higher. These models benefit from relatively straightforward architectures and moderate computational requirements, making them suitable for scenarios where model interpretability and resource constraints are important considerations.

The Transformer-based model, built on DistilBERT, delivered the highest performance among the four approaches, achieving a test accuracy of 98.01% with an almost perfect ROC-AUC of 0.998. By harnessing the power of self-attention and pretrained language representations, this model excels at understanding nuanced, context-rich language patterns commonly found in Reddit posts. Although its fine-tuning process is computationally more demanding, DistilBERT's superior accuracy makes it highly attractive for applications where predictive performance is paramount.

On the other hand, the CNN-based text model—designed to capture local n-gram features through convolutional filters—demonstrated competitive performance with a test accuracy of approximately 97.15%. Its fast convergence and lower computational overhead offer advantages in real-time or resource-constrained environments, despite it achieving marginally lower accuracy compared to the BiLSTM and Transformer models.

Overall, while all four models are viable options for detecting mental health conditions from Reddit posts, the Transformer-based DistilBERT model emerges as the best candidate if maximum accuracy and robust language understanding are the primary objectives. Its ability to generalize well to unseen data, as indicated by its high ROC-AUC, makes it ideal for deployment in systems where precision is critical. However, for applications where computational efficiency, speed, or ease of deployment are more critical, the CNN-based or BiLSTM approaches present strong alternatives.

This comparative evaluation underscores the importance of aligning model choice with the specific operational needs and resource constraints of the deployment environment. Future work could explore ensemble techniques or further domain-specific pretraining to enhance performance even further, ensuring that the chosen approach continues to evolve in response to the unique challenges posed by mental health detection from social media text.

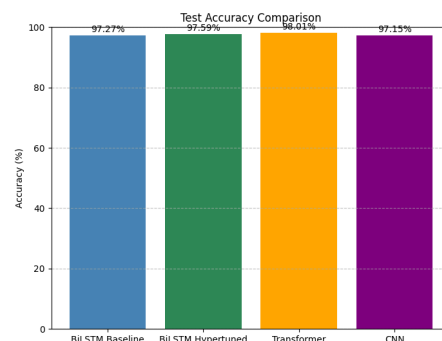


Fig 6: Accuracy evaluation of all four models

References

1. Low, D. M., Rumker, L., Torous, J., Cecchi, G., Ghosh, S. S., & Talkar, T. (2020). Natural Language Processing Reveals Vulnerable Mental Health Support Groups and Heightened Health Anxiety on Reddit During COVID-19: Observational Study. *Journal of medical Internet research*, 22(10), e22635.
2. Yazdavar, A. H., Mahdavinnejad, M. S., Bajaj, G., et al. "Multimodal mental health analysis in social media." *PLOS ONE* 15 (4): e0226248, 2020.