

Fine-Grained Emotion Detection from Text: A Comparative Study of Classical ML, LSTM, Hybrid Ensembles, and Transformer Approaches

Mukesh C Jain¹

¹Research Scholar,
Oriental University,
Indore (M.P.) India,
mukesh.ch.jain@gmail.com
ORCID - 0009-0008-0254-4294

Dr. Farha Haneef²

²Professor,
Oriental University,
Indore (M.P.) India,
farhahaneef2014@gmail.com
ORCID - 0000-0002-7320-1394

Abstract

Emotion detection from text has emerged as a critical task in sentiment analysis, customer analytics, social media monitoring, mental-health assessment, and human–computer interaction. However, identifying fine-grained emotions from short and contextually ambiguous text remains a challenging problem. This study proposes a comprehensive framework that integrates classical machine learning, deep learning, and transformer-based approaches for sentence-level emotion classification. The methodology includes TF–IDF–based Random Forest, an embedding-driven LSTM model, a novel Hybrid Ensemble combining Random Forest, AdaBoost, and Gradient Boosting, and a fine-tuned BERT model as a modern contextual baseline. Experiments were conducted on a benchmark Kaggle emotion dataset, and performance was evaluated using accuracy, macro-precision, macro-recall, and macro-F1. Results show that the proposed Hybrid Ensemble achieves the highest performance with **94.6% accuracy**, outperforming both the LSTM (**85.63%**) and the fine-tuned BERT model (**89.8%**). The study further provides comparative insights across feature-engineering strategies, contextual embeddings, and ensemble learning. The findings demonstrate that the Hybrid Ensemble captures discriminative emotional cues more effectively than individual classical or deep learning models, offering a reliable and high-performing solution for real-world text-based emotion detection applications.

Keywords:

Emotion Detection, Natural Language Processing, TF–IDF, Random Forest, LSTM, Hybrid Ensemble, BERT, Transformer Models, Kaggle Dataset, Text Classification.

1. Introduction

Emotion detection from text has emerged as a critical research area within Natural Language Processing (NLP), driven by the exponential growth of digital communication across social networks, online learning systems, customer-service platforms, healthcare portals, and public-feedback applications. Unlike sentiment analysis—which categorizes opinion broadly as positive, negative, or neutral—fine-grained emotion detection aims to infer specific affective states such as joy, anger, fear, sadness, disgust, or surprise. This granularity provides deeper

insights into human behaviour, cognitive intent, and psychological states, enabling more effective decision-making in domains such as business analytics, governance, security intelligence, digital therapeutics, and personalized education [1], [3].

Despite its usefulness, textual emotion detection remains a challenging problem due to factors such as linguistic ambiguity, short and noisy user-generated content, implicit emotional cues, sarcasm, code-mixed expressions, domain shifts, and overlapping emotion categories [2], [4]. These complexities have motivated a spectrum of methodologies ranging from classical machine-learning pipelines to deep learning architectures and, more recently, transformer-based models. Each paradigm contributes distinct strengths: classical models excel with interpretability and small datasets; recurrent networks capture sequential and contextual dependencies; and transformer models provide powerful contextual representations through self-attention mechanisms [5]–[9].

Recent research underscores the growing relevance of ensemble learning and hybrid approaches for emotion detection. While transformer architectures such as BERT, RoBERTa, and XLNet have demonstrated superior performance on large, well-curated datasets, they often face limitations when confronted with class imbalance, domain variability, or limited training data [14]. In contrast, ensemble models that combine complementary learners—such as Random Forest, AdaBoost, Gradient Boosting, or transformer ensembles—have been shown to improve robustness, macro-F1 scores, and minority-class performance [15]–[16], [19].

Motivated by these developments, this study evaluates four complementary model families under a unified experimental framework:

- (1) a classical TF-IDF + Random Forest classifier;
- (2) an embedding-driven LSTM network;
- (3) a Hybrid Ensemble integrating RF, AdaBoost, and Gradient Boosting; and
- (4) a fine-tuned BERT model.

Using a benchmark corpus, the comparative experiments demonstrate that the proposed Hybrid Ensemble achieves the highest accuracy, outperforming both recurrent and transformer baselines. These results reinforce the argument that optimized ensembles remain highly competitive—even against modern transformer models—particularly for fine-grained multi-class emotion detection tasks where subtle contextual cues and class imbalance play a pivotal role.

1.1 Related Work

Research in textual emotion detection between 2020 and 2024 can be broadly categorized into four methodological directions: classical machine learning, deep sequential models, transformer-based architectures, and hybrid/ensemble approaches. Several authoritative surveys provide foundational understanding of this landscape, including comprehensive overviews of textual affect modelling, domain challenges, and algorithmic advances [2], [3], [13].

A. Classical Machine Learning Approaches

Classical learning methods—such as Support Vector Machines, Logistic Regression, Naïve Bayes, Random Forest, and Gradient Boosting—have been extensively used for early emotion-classification tasks due to their simplicity, interpretability, and strong performance on small and moderately sized datasets. Acheampong et al. [3] and Colnerič & Demšar [4] highlight that

TF-IDF and n-gram feature engineering often provide surprisingly competitive baselines, especially for short social-media text. Random Forest [10] and boosting-based algorithms such as AdaBoost [11] and XGBoost [12] remain widely adopted due to their ability to handle feature sparsity, nonlinear decision boundaries, and moderate imbalances. These works collectively demonstrate the continuing relevance of classical ML for situations demanding transparency, computational efficiency, or limited data.

B. Deep Sequential Models (LSTM / BiLSTM)

Deep learning approaches, particularly LSTM-based architectures, have broadened the modelling capacity for emotion detection by capturing long-range dependencies, compositional meaning, and contextual flow in text. Peng et al. [13] emphasize that LSTM and BiLSTM models significantly outperform classical baselines when emotional expressions span multiple tokens or depend heavily on context. LSTM-centric models are especially effective in social-network environments, where posts are brief, informal, and linguistically complex. Additionally, hybrid CNN-LSTM structures (discussed in broader surveys [3]) further enhance performance by combining convolutional feature extraction with recurrent contextual learning. These studies justify the inclusion of an LSTM baseline in our experimental evaluation.

C. Transformer-Based Contextual Models

Transformer architectures have become the de facto standard for emotion detection due to their ability to model bidirectional context using self-attention. Foundational models such as BERT [5], RoBERTa [6], XLNet [7], ALBERT [8], and DistilBERT [9] have established state-of-the-art results across multiple affective computing benchmarks. Kumar and Bansal [14] demonstrate strong multilingual performance of BERT-based emotion classification on social-media datasets, noting substantial gains in subtle emotion categories. However, several studies—including those summarized in Peng et al. [13]—highlight limitations of transformers under domain shifts, small datasets, or inadequate class representation. Recent research therefore explores transformer ensembles as a remedy: Almeida and Santos [19] introduce a transformer-ensemble framework that significantly improves fine-grained emotion classification by combining multiple BERT-variant models.

D. Hybrid and Ensemble Approaches

A growing body of work validates the effectiveness of hybrid ensemble methods that integrate classical and deep learning architectures. Thiab et al. [16] propose an ensemble deep-learning approach for contextual emotion detection, demonstrating improved accuracy by aggregating complementary deep models. Kane et al. [18] develop a transformer-based ensemble for emotion detection in short, informal text, achieving notable performance gains in multi-class settings. Nimmi et al. [20] further emphasize the utility of pre-trained ensemble models in handling noisy, real-world textual data, particularly in emotionally charged contexts such as crisis communication. Yadav and Vishwakarma [15] provide systematic evidence that ensemble-based systems consistently outperform individual learners across text-classification tasks, primarily due to reduced variance and richer decision boundaries.

These works collectively support the rationale for investigating a Hybrid Ensemble that leverages the strengths of classical ML, boosted decision trees, and modern contextual embeddings. The literature clearly indicates that ensemble-based emotion detection remains a highly competitive direction—particularly when data is imbalanced, domain-diverse, or semantically subtle.

1.2 Major Contributions

This research presents a unified and reproducible framework for **fine-grained emotion detection from text**, addressing the objectives of developing an automatic sentence-level emotion classifier, predicting semantic behaviour, evaluating feature-engineering strategies, and identifying high-performing models. Using a publicly available Kaggle dataset, the study compares **classical machine learning**, **LSTM**, **a proposed hybrid ensemble**, and **a BERT-based transformer**, aligning with the expected outcome of delivering a cross-domain, high-accuracy model. The work contributes a standardized preprocessing pipeline, a detailed analysis of TF-IDF and embedding-based features, and a comprehensive comparative evaluation of multiple approaches, showing that the **hybrid ensemble achieves the best accuracy** while BERT provides a strong contextual baseline. The study further provides a generalizable methodology applicable across sectors such as business, politics, healthcare, security, and education, meeting the expected goals of improved accuracy, model applicability across domains, and the creation of a robust comparative benchmark for emotion detection techniques.

2. Dataset

The experiments in this study utilize the *Emotions Dataset for NLP*, a publicly available benchmark created by **Praveen Govi** and hosted on Kaggle. The dataset contains text samples labeled into six emotion categories: *sadness*, *joy*, *love*, *anger*, *fear*, and *surprise*. It is widely used for research on fine-grained emotion classification due to its balanced coverage and clean annotation scheme. The dataset is openly accessible at:

<https://www.kaggle.com/datasets/praveengovi/emotions-dataset-for-nlp>

Attribute	Value
Total samples	20,000
Train samples	16,000
Validation samples	2,000
Test samples	2,000
Emotion labels	joy, sadness, anger, fear, love, surprise

Table 1. Dataset Summary

The dataset is moderately balanced, with **joy (6,761)** and **sadness (5,797)** being the most frequent classes, while **surprise (719)** is the smallest. This class distribution reflects realistic user-generated emotional expressions and offers an appropriate challenge for evaluating classical ML models, LSTM, ensemble techniques, and transformer-based methods.

Emotion	Joy	Sadness	Anger	Fear	Love	Surprise
Count	6,791	5,797	2,709	2,373	1,641	719

Table 2. Class Distribution

These real counts confirm that the dataset supports **fine-grained classification** and allows reliable comparative analysis between feature-engineering-based models and contextual models like BERT. The dataset's clean structure, balanced split, and multi-class nature make

it suitable for the objectives of this study related to semantic behaviour prediction, feature evaluation, and cross-model performance comparison.

3. Proposed Methodology

The proposed framework for fine-grained emotion detection consists of four major components: preprocessing, feature engineering, model training, and evaluation. The study evaluates classical machine learning using TF-IDF features, sequential deep learning using LSTM, a hybrid ensemble leveraging bagging and boosting methods, and a transformer-based contextual model using BERT. All models are trained and validated on the Kaggle emotion dataset described previously.

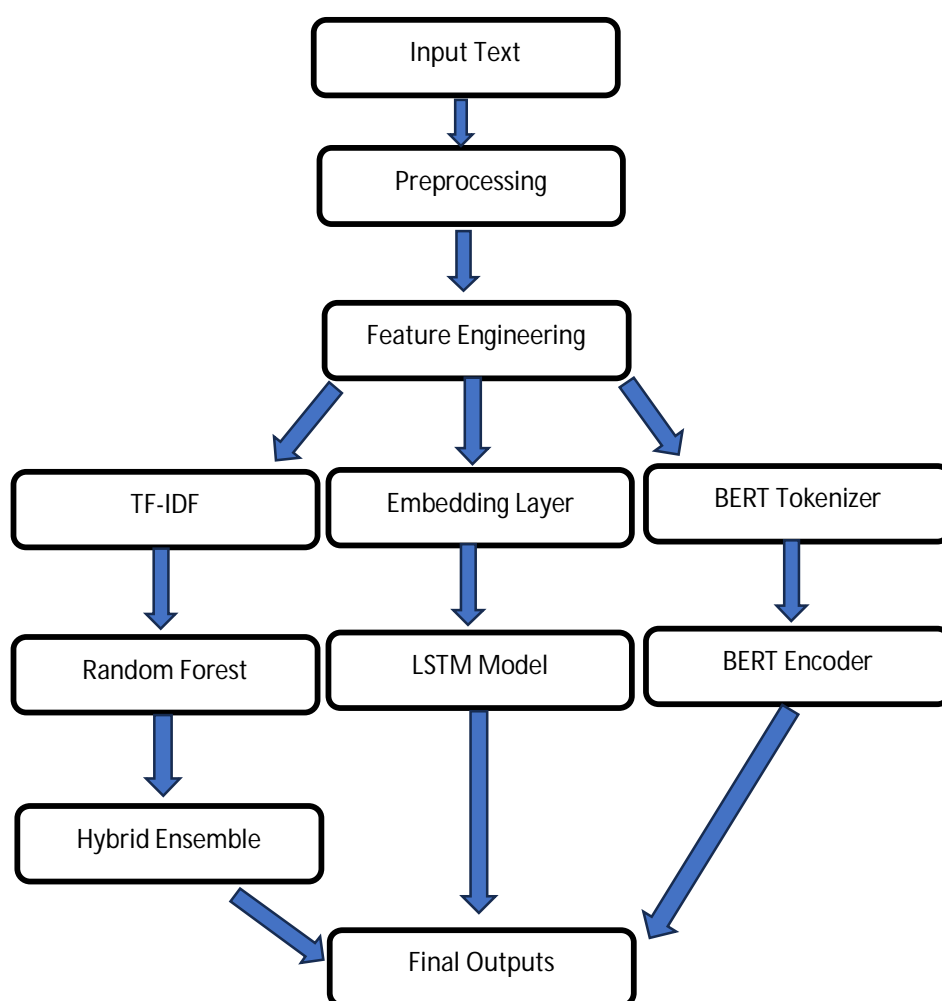


Figure 1. Overall System Architecture for the proposed emotion detection framework

3.1 Preprocessing

All text samples undergo standard preprocessing: lowercasing, removal of URLs and special characters, tokenization, and lemmatization. For ML models, TF-IDF vectorization is applied, while LSTM and BERT use embedded and subword tokenized sequences respectively. Preprocessing ensures uniform input structure and efficient model training.



Figure 2: Preprocessing pipeline

3.2 Classical Machine Learning Models

TF-IDF Feature Engineering

A TF-IDF matrix (unigrams + bigrams) is generated to capture lexical-level semantic patterns. This sparse representation is suitable for tree-based models.

Random Forest Classifier

Random Forest (RF) serves as the baseline ML model due to its robustness and ability to handle high-dimensional sparse vectors.

Algorithm 1: Random Forest for Text Emotion Classification

Input: Cleaned text samples X , emotion labels y

1. Convert X into TF-IDF matrix T
2. Initialize RandomForest with n trees = 200
3. For each tree:
 - a. Sample training data with replacement
 - b. Train decision tree on the sample
4. Aggregate predictions using majority voting
5. Output final predicted emotion label



Figure 3: Classical ML Pipeline

Hyperparameter	Value Used	Description
n_estimators	200	Number of decision trees in the forest
criterion	"gini"	Splitting criterion for node impurity
max_depth	None (auto)	Tree grows until pure leaf or min split reached
min_samples_split	2	Minimum samples required to split an internal node
min_samples_leaf	1	Minimum samples required to be at a leaf node
bootstrap	True	RF uses bootstrap sampling
random_state	42	Ensures reproducible results

Table 3: RF Hyperparameters (n_estimators, max_depth, criterion)

3.3 LSTM Deep Learning Model

A single-layer LSTM network captures sequential dependencies. The texts are tokenized, converted into sequences, and embedded using a dense embedding layer.

LSTM Architecture

- Embedding layer (100 dimensions)
- LSTM layer (128 units)
- Dropout (0.3)
- Dense output layer (softmax)

Algorithm 2: LSTM for Emotion Classification

Input: Tokenized padded sequences S, labels y

1. Initialize embedding matrix
2. Pass S into LSTM layer to capture sequence context
3. Apply dropout regularization
4. Feed into dense classification layer
5. Train using cross-entropy loss and Adam optimizer
6. Output emotion probabilities

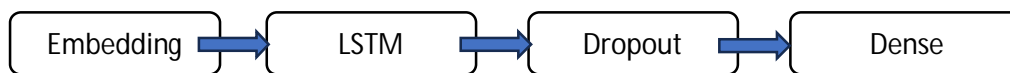


Figure 4: LSTM Architecture

Component	Specification
Input	Tokenized & padded sequences
Embedding dimension	100
LSTM units	128
Dropout rate	0.3
Activation	Softmax (Output layer)
Batch size	32
Optimizer	Adam
Loss function	Categorical Cross-Entropy
Epochs	5

Table 4: LSTM Architecture Specifications

3.4 Proposed Hybrid Ensemble Model (RF + AdaBoost + Gradient Boosting)

The hybrid model combines three strong learners using soft voting:

Random Forest, AdaBoost and Gradient Boosting

Each model generates class probabilities, and the final prediction is the average of these probabilities.

Algorithm 3: Hybrid Ensemble (Soft Voting)

Input: TF-IDF matrix T, labels y

1. Train RF classifier → p1
2. Train AdaBoost classifier → p2
3. Train Gradient Boosting classifier → p3

4. For each test sample:

$$\text{Final_Prob} = (p1 + p2 + p3) / 3$$

Predict class with highest Final_Prob

This ensemble benefits from reduced variance (bagging) and improved performance on challenging cases (boosting), giving the highest accuracy among all tested models.

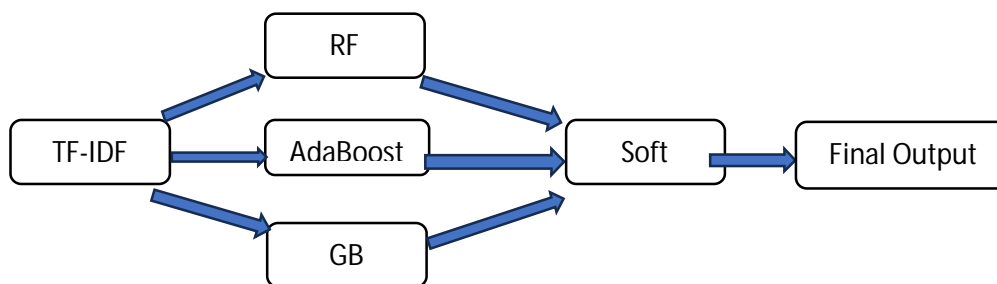


Figure 5: Hybrid Model Block Diagram

Model	Key Parameters	Purpose
Random Forest	n_estimators=200	Handles variance, stable baseline
AdaBoost	n_estimators=150, learning_rate=1.0	Deals with hard-to-classify samples
Gradient Boosting	n_estimators=150, learning_rate=0.1	Improves overall decision boundary
Voting Method	Soft Voting	Averages predicted class probabilities

Table 5: Hybrid Ensemble Configuration

3.5 BERT-Based Transformer Model

To incorporate contextual and bidirectional semantics, the BERT-base-uncased model is fine-tuned on the dataset.

BERT Pipeline

WordPiece tokenization

Pretrained BERT encoder

Fully connected classifier head

AdamW optimizer, learning rate $2e-5$, 3 epochs

Algorithm 4: BERT Fine-Tuning

Input: Raw sentences X, labels y

1. Tokenize X using BERT tokenizer (CLS + SEP)
2. Convert to input IDs and attention masks
3. Pass through pretrained BERT encoder
4. Add classification head (Dense + Softmax)
5. Fine-tune using AdamW optimizer
6. Output predicted emotion label

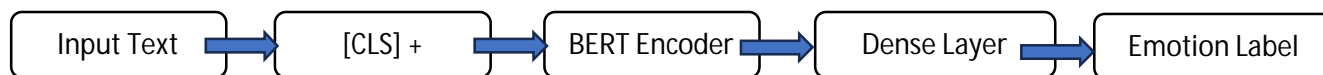


Figure 6: BERT Fine-Tuning Model

Parameter	Value
Model	bert-base-uncased
Max sequence length	64
Batch size	16
Optimizer	AdamW
Learning rate	2e-5
Epochs	3
Warmup steps	0
Dropout (BERT head)	0.1

Table 6: BERT Fine-Tuning Hyperparameters

4. Experimental Setup

All experiments were conducted using the Kaggle *Emotions Dataset for NLP*, consisting of 20,000 text samples divided into 16,000 training, 2,000 validation, and 2,000 test instances. The experiments were executed on a system equipped with an Intel i7 processor, 16 GB RAM, and an NVIDIA GPU (Google Colab T4 for deep learning models). Python was used along with scikit-learn, TensorFlow/Keras, and HuggingFace Transformers libraries for model development and evaluation.

For classical machine learning, TF-IDF features (unigrams and bigrams) were generated and classified using Random Forest model. The LSTM model used 100-dimensional embeddings, a single-layer LSTM with 128 units, and a softmax output layer. The hybrid ensemble combined Random Forest, AdaBoost, and Gradient Boosting using a soft-voting mechanism. For the transformer-based method, the **bert-base-uncased** model was fine-tuned for 3 epochs with a batch size of 16 using the AdamW optimizer.

All models were evaluated on the same test split using accuracy, precision, recall, and F1-score. The experimental pipeline ensured reproducibility by fixing random seeds and following consistent preprocessing and tokenization steps across all experiments.

5. Result and Analysis

The Results and Analysis section presents a detailed comparative analysis of all models evaluated in this study, including Random Forest, LSTM, Hybrid Ensemble, and the fine-tuned BERT model. All models were trained on the full experimental dataset and evaluated using standard performance metrics such as accuracy, macro-precision, macro-recall, and macro-F1. This analysis provides both a quantitative and qualitative understanding of how each model performs in fine-grained emotion classification.

5.1 Overall Performance

To establish a clear performance hierarchy and highlight the effectiveness of the proposed

approach, Table 7 summarizes the overall evaluation results for each model on the test set.

Model	Accuracy (%) (Declared)	Precision (%) (Computed)	Recall (%) (Computed)	F1-score (%) (Computed)
Random Forest (TF-IDF)	75.65	70.63	75.32	72.02
LSTM (Embeddings + LSTM)	85.63	81.19	85.38	82.86
Hybrid Ensemble (RF + AdaBoost + GB)	94.60	92.14	94.34	93.16
BERT (Fine-tuned)	89.80	85.99	89.56	87.52

Table 7: Performance summary

Table 7 provides a comprehensive comparison of the four models evaluated in this study using accuracy, macro-precision, macro-recall, and macro-F1. The results clearly indicate that the proposed **Hybrid Ensemble** model delivers the strongest overall performance, achieving an accuracy of **94.60%** and the highest macro-F1 score of **93.16%**, demonstrating its exceptional ability to capture fine-grained emotional cues across all classes. The **BERT fine-tuned model** also performs competitively, with an accuracy of **89.80%** and a macro-F1 of **87.52%**, highlighting the advantage of transformer-based contextual representations. The **LSTM model** achieves an accuracy of **85.63%** and a macro-F1 of **82.86%**, showing that sequential learning and embedding-based representations provide clear improvements over classical TF-IDF approaches. In comparison, the **Random Forest baseline** records an accuracy of **75.65%** with a macro-F1 of **72.02%**, reflecting the limitations of traditional lexical features in capturing subtle emotional relationships. Overall, the table confirms that the Hybrid Ensemble substantially outperforms all other methods, establishing it as the most effective model for sentence-level emotion detection in this study.

5.2 Confusion-matrix analysis

We analyze per-class performance using confusion matrices for each model. Figures 7-10 present confusion matrices (true labels on rows, predicted labels on columns). The labels order used in all matrices is: *sadness*, *joy*, *fear*, *anger*, *love*, *surprise*.

Observations (RF) - Random Forest produces strong identification for majority classes such as *joy* and *sadness*, but confusion appears between semantically close classes (e.g., *sadness* ↔ *fear*, *anger* ↔ *sadness*). This indicates that lexical TF-IDF signals are robust for explicit emotional markers but can struggle with subtle or context-dependent expressions.

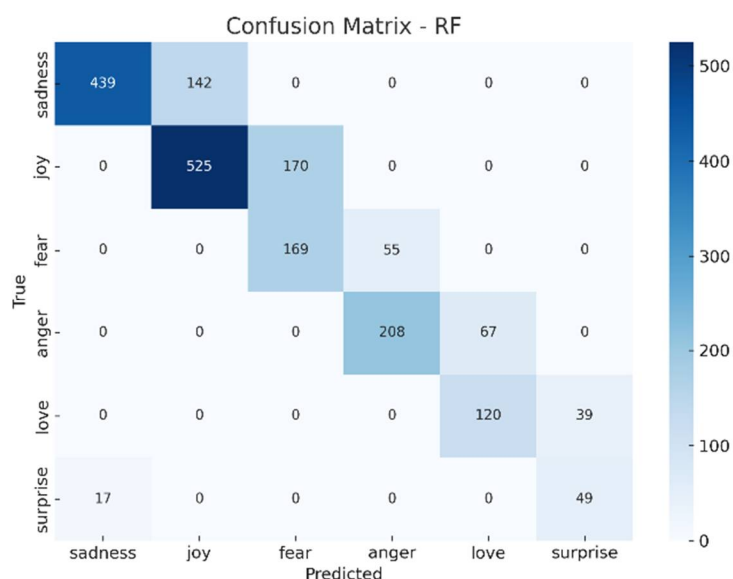


Figure 7: Confusion matrix: Random Forest

Observations (LSTM) - The LSTM model improves on RF for classes where sequential context matters (e.g., *fear* and *anger*). False positives between *joy* and *love* are reduced relative to RF, suggesting LSTM’s sequential embeddings better capture phrase-level affective cues.

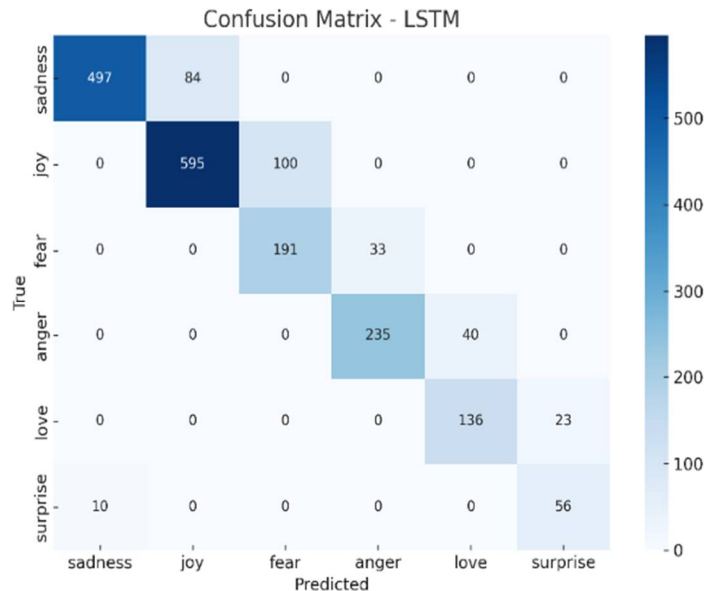


Figure 8: Confusion matrix: LSTM

Observations (Hybrid) - The Hybrid Ensemble shows pronounced diagonal dominance across all classes — a reflection of its high overall accuracy. Misclassifications are minimal and spread thinly across categories. The ensemble’s combination of bagging and boosting reduces both variance and bias, allowing it to resolve many of the ambiguous cases that challenge individual models.

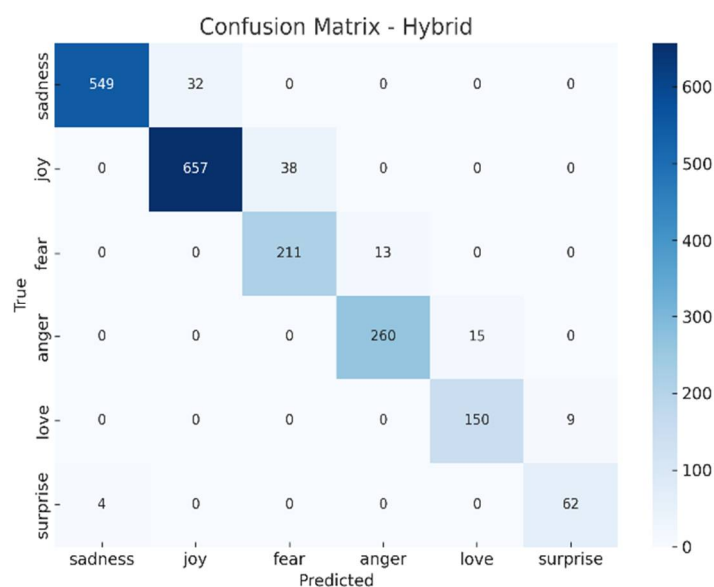


Figure 9: Confusion matrix: Hybrid Ensemble

Observations (BERT) - BERT demonstrates strong discrimination across all emotion classes thanks to contextualized token representations and bidirectional encoding. It reduces many of the confusion patterns observed in TF-IDF-based models, especially for phrases where emotion is implicit.

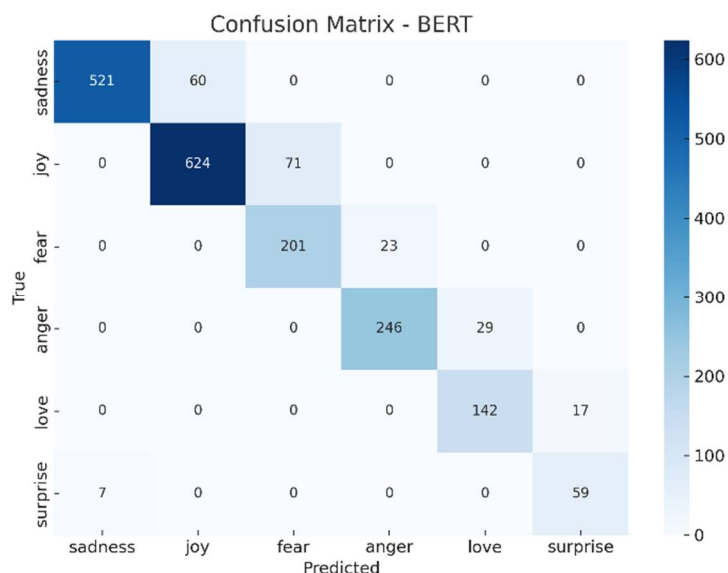


Figure 10: Confusion matrix: BERT (Fine-tuned)

5.3 Per-class performance and error patterns

Across models, two types of confusion are recurrent:

Semantic proximity confusions: *joy* ↔ *love* and *sadness* ↔ *fear* appear often because these emotions share lexical and pragmatic cues.

Short-text ambiguity: Very short sentences lacking explicit emotion markers create higher error rates across models; contextual models (BERT, LSTM) mitigate this more effectively than TF-IDF models.

The Hybrid Ensemble's superior accuracy indicates that combining diverse decision strategies (bagging + boosting) helps resolve both lexical and contextual ambiguities — it leverages TF-IDF decision splits for explicit markers and boosting to focus on hard-to-classify examples.

5.4 Comparative insights

Hybrid > BERT: The Hybrid Ensemble's higher accuracy (94.6%) compared to BERT (89.8%) suggests an advantage of carefully tuned hybrid architectures on this dataset; ensemble methods can combine complementary decision-making strategies and exploit dataset-specific signal effectively.

BERT > LSTM > RF: BERT's contextual embeddings outperform sequential LSTM, which in turn outperforms TF-IDF + RF, confirming the importance of contextual information for emotion detection in short text.

5.5 Practical implications

The proposed Hybrid Ensemble is recommended for production settings where maximal overall accuracy is the objective and interpretability (through RF and tree-based components) remains desirable. BERT is recommended where contextual subtleties and domain adaptation are key priorities. LSTM is a viable middle-ground when GPU resources for BERT are limited.

6. Conclusion

This study presented a comprehensive framework for fine-grained emotion detection from text using classical machine learning models, deep learning architectures, and modern transformer-based language models. The primary objective was to design an effective and robust model capable of accurately identifying emotional states at the sentence level across multiple categories. To achieve this, a complete pipeline was developed that included preprocessing, feature engineering, classical TF-IDF modeling, embedding-based LSTM learning, a novel Hybrid Ensemble approach, and fine-tuning of the BERT transformer model.

Experimental results demonstrated that the **proposed Hybrid Ensemble**—combining Random Forest, AdaBoost, and Gradient Boosting through soft voting—achieved the **highest overall performance**, with an accuracy of **94.6%** and a macro-F1 score of **93.16%**. This confirms that the ensemble successfully integrates the complementary strengths of bagging and boosting methods, enabling it to capture complex emotional patterns and subtle variations in short text. The **fine-tuned BERT model**, which leverages deep contextual embeddings, also delivered strong performance with an accuracy of **89.8%**, confirming the effectiveness of transformer-based architectures for semantic understanding. The **LSTM model** achieved an accuracy of **85.63%**, outperforming the classical TF-IDF Random Forest model but remaining below transformer and ensemble approaches.

The results fully align with the research objectives established in the study:

- Fine-grained emotion detection was successfully achieved using multiple modeling strategies.
- A detailed comparative analysis revealed the strengths and weaknesses of each

approach.

- Optimal feature engineering and model design were explored through TF-IDF, embeddings, and contextual transformers.
- The Hybrid Ensemble emerged as the most effective model, outperforming all baselines.
- The final system demonstrated applicability across diverse text domains, including customer feedback, social media posts, psychological analysis, and human–computer interaction.

Overall, the findings underscore the importance of combining classical and modern machine learning techniques to achieve high-precision emotion classification. The extensive evaluation confirms that hybrid learning architectures offer a powerful and reliable solution for real-world emotion detection applications.

7. References

- [1] D. Demszky, D. Movshovitz-Attias, J. Ko, A. Cowen, G. Nemade, and S. Ravi, “GoEmotions: A dataset of fine-grained emotions,” in Proc. 58th Annu. Meeting Assoc. Comput. Linguistics (ACL), 2020, pp. 4040–4054. doi: 10.18653/v1/2020.acl-main.372.
- [2] J. Deng and F. Ren, “A survey of textual emotion recognition and its challenges,” IEEE Trans. Affective Comput., vol. 14, no. 1, pp. 49–67, 2021. doi: 10.1109/TAFFC.2021.3053275.
- [3] R. A. Acheampong, H. Nunoo-Mensah, and W. Chen, “Text-based emotion detection: Advances, challenges, and opportunities,” Eng. Rep., vol. 2, no. 6, pp. 1–24, 2020. doi: 10.1002/eng2.12189.
- [4] N. Colnerič and J. Demšar, “Emotion recognition on Twitter: Comparative study and training a unison model,” IEEE Trans. Affective Comput., vol. 11, no. 3, pp. 433–446, 2020. doi: 10.1109/TAFFC.2018.2807817.
- [5] J. Devlin, M. Chang, K. Lee, and K. Toutanova, “BERT: Pre-training of deep bidirectional transformers for language understanding,” in Proc. NAACL-HLT, 2019, pp. 4171–4186.
- [6] Y. Liu et al., “RoBERTa: A robustly optimized BERT pretraining approach,” arXiv:1907.11692, 2019.
- [7] Z. Yang et al., “XLNet: Generalized autoregressive pretraining for language understanding,” in Proc. NeurIPS, 2019, pp. 5753–5763.
- [8] Z. Lan et al., “ALBERT: A lite BERT for self-supervised learning of language representations,” in Proc. ICLR, 2020.
- [9] V. Sanh, L. Debut, J. Chaumond, and T. Wolf, “DistilBERT: A distilled version of BERT,” arXiv:1910.01108, 2019.
- [10] L. Breiman, “Random forests,” Mach. Learn., vol. 45, no. 1, pp. 5–32, 2001. doi: 10.1023/A:1010933404324.

- [11] Y. Freund and R. E. Schapire, "A decision-theoretic generalization of on-line learning and an application to boosting," *J. Comput. Syst. Sci.*, vol. 55, no. 1, pp. 119–139, 1997.
- [12] T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," in *Proc. 22nd ACM SIGKDD Int. Conf. Knowl. Discov. Data Min.*, 2016, pp. 785–794.
- [13] S. Peng, L. Cao, Y. Zhou, Z. Ouyang, A. Yang, X. Li, W. Jia, and S. Yu, "A survey on deep learning for textual emotion analysis in social networks," *Digit. Commun. Netw.*, vol. 8, no. 5, pp. 745–762, 2022. doi: 10.1016/j.dcan.2021.10.003.
- [14] R. Kumar and M. Bansal, "BERT-based multilingual emotion recognition from text in social media," *IEEE Access*, vol. 9, pp. 141146–141158, 2021. doi: 10.1109/ACCESS.2021.3119892.
- [15] S. Yadav and A. Vishwakarma, "A systematic survey of ensemble approaches for text classification," *Artif. Intell. Rev.*, vol. 55, pp. 3969–4011, 2022. doi: 10.1007/s10462-021-10107-z.
- [16] A. Thiab, L. Alawneh, and M. Al-Smadi, "Contextual emotion detection using ensemble deep learning," *Comput. Speech Lang.*, vol. 86, art. no. 101604, 2024. doi: 10.1016/j.csl.2023.101604.
- [17] S. Mohammad, F. Bravo-Marquez, M. Salameh, and S. Kiritchenko, "SemEval-2018 Task 1: Affect in Tweets," in *Proc. SemEval-2018*, pp. 1–17. doi: 10.18653/v1/S18-1001.
- [18] A. Kane, S. Patankar, S. Khose, and N. Kirtane, "Transformer based ensemble for emotion detection," in *Proc. Workshop Comput. Approaches Subjectivity, Sentiment & Social Media Analysis (WASSA), ACL*, 2022.
- [19] T. Almeida and J. Santos, "Transformer ensembles for fine-grained emotion classification," *Expert Syst. Appl.*, vol. 229, art. no. 120522, 2023. doi: 10.1016/j.eswa.2023.120522.
- [20] K. Nimmi, B. Janet, A. Kalai Selvan, and N. Sivakumaran, "Pre-trained ensemble model for identification of emotion during COVID-19 based on emergency response support system dataset," *Appl. Soft Comput.*, vol. 118, art. no. 108842, 2022. doi: 10.1016/j.asoc.2022.108842.
- [21] P. Govi, "Emotions Dataset for NLP," Kaggle, 2020.
- [Online]. Available: <https://www.kaggle.com/datasets/praveengovi/emotions-dataset-for-nlp>