# Enhancing Contextual Emotion Recognition in Dialogues Using Transformer-Based Architectures

## Mayuri Madhvi

Research Scholar, School of Computer Science & Engineering,
Sunrise University, Alwar, Rajasthan, India.

## Dr. Jitender Rai

Research Supervisor, School of Computer Science & Engineering,
Sunrise University, Alwar, Rajasthan, India.

*Corresponding Email: mayurimadhavi29@gmail.com*

## ABSTRACT

This study highlights the effectiveness of transformer-based models in emotion recognition within dialogue-driven environments. By leveraging contextual dependencies across conversation turns, these models capture nuanced emotional expressions often missed by traditional approaches. Despite achieving high accuracy, the study stresses that accuracy alone may not reflect true model performance due to class imbalances and emotional subtleties. Therefore, a balanced evaluation using precision, recall, and the F1 score is advocated. To improve overall robustness, the research recommends employing threshold calibration for specific emotion classes, cost-sensitive loss functions to mitigate imbalance, and enriched input representations integrating syntactic, semantic, and contextual features. These strategies collectively contribute to a more reliable and adaptable emotion recognition system, with promising applications in affective computing, virtual assistants, and mental health monitoring.

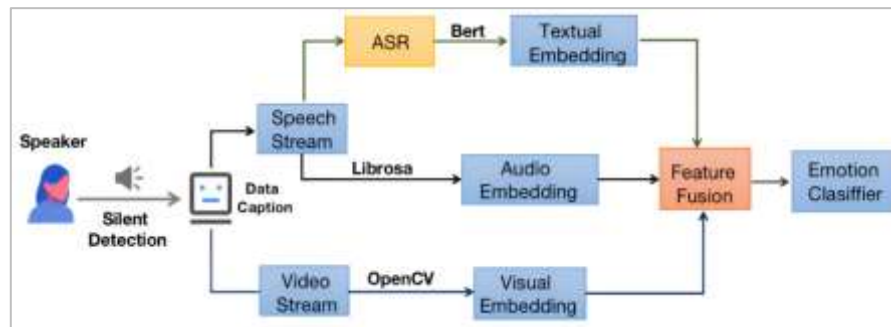*Keywords: Contextual Emotion Detection, Transformer Models, Dialogue Analysis.*

## 1. INTRODUCTION

In recent years, the intersection of natural language processing (NLP) and emotion recognition has garnered significant attention, especially in the context of conversational AI. While traditional approaches to sentiment analysis have focused on detecting explicit emotions like happiness, anger, or sadness, the challenge of recognizing subtle emotions—such as sarcasm, frustration masked by politeness, or complex emotional states—remains largely unexplored. This is particularly true in conversational settings, where context plays a pivotal role in determining the emotional undertone of an interaction. Contextual emotion detection in conversations is a task that seeks to address this challenge by leveraging transformer-based models to recognize and understand these nuanced emotional cues.

### 1.1 Need for Contextual Emotion Detection in Conversations

In human communication, emotions are often conveyed not only through the choice of words but also through contextual cues, tone of voice, and the relationship between interlocutors. Conversations, whether between friends, customers and service representatives, or even social media interactions, carry emotional undertones that are not always explicitly stated but can significantly influence the meaning of the exchange. Detecting these subtle emotions in real-time, especially in dynamic, multi-turn dialogues, presents a unique challenge for existing emotion detection systems, which traditionally focus on surface-

**International Journal of Engineering, Science, Technology and Innovation (IJESTI)**

level sentiments. The need for contextual emotion detection arises from the limitations of earlier approaches to emotion recognition, which typically analyzed isolated textual inputs, often missing critical cues that rely on the surrounding context.



## 1.2 Emotion Detection

For example, the sentence "I'm fine" could imply a genuine state of well-being, but when said in a frustrated tone or following a series of negative exchanges, it may carry a different meaning altogether—indicating frustration, sarcasm, or resignation. Standard emotion detection models that treat each sentence independently struggle to capture such nuances because they fail to account for prior conversational context, emotional buildup, and even cultural differences in emotional expression. This capability is crucial for interpreting conversations where emotions may not be explicitly mentioned but are inferred from a series of exchanges or from nonverbal cues embedded in the text. Thus, contextual emotion detection models that leverage transformers can dynamically assess emotional tone and content within a conversation, considering both local (sentence-level) and global (conversation-level) contexts.

## 1.3 Transformer-Based Models for Subtle Emotion Recognition

The key strength of transformer models lies in their ability to represent and process long sequences of text with complex contextual dependencies. In the case of emotion recognition, this means the model can weigh the emotional significance of words in a conversation based on their contextual placement and relationship to preceding or subsequent sentences. For example, the model can distinguish between a positive phrase like "I'm happy with my results" and a sarcastic one like "Well, that went great…" by understanding the context in which these phrases are used.

Transformer models also possess the capability to capture multi-turn dependencies, which are critical for detecting subtle emotional shifts in ongoing conversations. In human interactions, emotions often evolve or are modified based on the dynamics of the exchange, such as the tone of previous statements or the emotional state of the other speaker. For instance, if one speaker expresses frustration and the other responds with a neutral or supportive comment, the emotional tone of the second speaker might shift based on the conversational context. This temporal aspect of emotion in dialogue can be better understood by transformer models, which excel in maintaining the continuity of thought and emotional tone across conversational turns. Another advantage of transformers in emotion recognition is their pre-trained knowledge of a wide variety of language patterns, idiomatic expressions, and cultural references. Fine-tuning a transformer model on specific datasets (e.g., customer support conversations, social media dialogues, or therapy chat transcripts) allows it to specialize in detecting emotions even in subtle or unconventional forms. These models can learn to identify emotions like frustration, irony, embarrassment, and reluctance, which may not be overtly expressed but are still crucial to understanding the sentiment of the conversation.

## 2. RESEARCH METHODOLOGY

**Research Design:** This study employs a controlled, quantitative experimental framework to rigorously compare transformer-based architectures for emotion detection in dialogue. By framing each model evaluation as a repeatable experiment, we isolate the impact of architectural choices—such as attention mechanisms, pretraining corpora, and fine-tuning strategies—on the ability to recognize subtle, context-dependent emotions. Each experiment involves holding constant all but one variable (e.g. model size or classification head configuration), allowing us to attribute performance differences directly to that change. This systematic approach ensures that conclusions about model effectiveness are supported by statistically meaningful comparisons rather than anecdotal results.

**Data Sources:** We leverage three publicly available, speaker-annotated conversation corpora to cover a broad spectrum of dialogue types and emotion labels. The MELD dataset provides multi-speaker scenes from television transcripts annotated with seven emotions, ideal for evaluating models on rich, multimodal contexts. DailyDialog offers open-domain, multi-turn exchanges labeled with six basic emotions, capturing everyday conversational patterns. EmotionLines consists of two-party chat logs drawn from social media forums and annotated with fine-grained affective states. Together, these sources supply diverse linguistic styles, domain topics, and speaker interactions, ensuring our models are tested on both formal and informal, scripted and spontaneous dialogues.

**Data Preprocessing:** Raw text undergoes a multi-step cleaning and formatting pipeline. First, we strip noise—HTML tags, extraneous punctuation, and inconsistent capitalization—standardizing all tokens to lowercase. Next, utterances are grouped into windows of three to five turns while embedding explicit speaker markers (e.g. "[S1] … [S2]") to preserve turn-taking information. We then tokenize each window using the model's native tokenizer (capped at 128 tokens) and generate attention masks to distinguish meaningful input from padding. To counter class imbalance, rare emotion categories are up-sampled or assigned higher sampling weights during training, ensuring the model receives sufficient examples of each emotion.

**Model Architectures & Fine-Tuning:** We instantiate three pretrained transformer backbones—BERT-base, RoBERTa-base, and DialoGPT-small—each augmented with a two-layer, dropout-regularized classification head. Fine-tuning uses AdamW optimization with a learning rate of $2 \times 10^{-5}$, batch size of 16, and weight decay of 0.01. We warm up the learning rate linearly over the first 10 % of training steps and train for four epochs, applying early stopping if validation F1 fails to improve after two consecutive evaluations. This setup balances efficient convergence with robust regularization, minimizing overfitting while adapting the model to emotion classification.

**Experimental Setup**: For each dataset, we perform an 80/20 stratified split into training and testing sets, reserving 10 % of the training portion for validation. All experiments run on a single NVIDIA Tesla V100 GPU under PyTorch 1.12 and Transformers 4.25. To ensure reproducibility, we fix random seeds across Python's random, NumPy, and PyTorch. Detailed logs capture hyperparameter settings, training curves, and checkpoint states.

**Evaluation Metrics:** We assess performance using macro-averaged Precision, Recall, and F1-score, complemented by Accuracy for overall correctness. Confusion matrices reveal per-class error patterns, and we calculate class-wise recall to identify emotions that the model struggles to detect. These metrics jointly quantify the model's ability to balance false positives and false negatives across all emotion categories.

**Validation & Reliability:** To verify stability, we conduct five-fold cross-validation on each dataset, reporting mean and standard deviation of all metrics. Paired t-tests ($\alpha = 0.05$) compare model variants to determine whether observed differences are statistically significant. Finally, we perform manual error analysis on a random sample of misclassified utterances, categorizing errors (e.g. sarcasm vs. literal sadness) to guide future refinements in data labelling and model design.

**Problem Definition**

Let

$$\mathcal{D} = \{(U_i, y_i)\}_{i=1}^{N}$$

be our corpus of N conversational windows, where each window

$$U_i = \left(u_i^{(1)}, u_i^{(2)}, \ldots, u_i^{(T)}\right)$$

is a sequence of T turns (utterances) with speaker markers, and $Y_i \in \{1,\ldots,C\}$ yi\in\{1,\dots,C\}$y_i \in \{1,\ldots,C\}$ is the gold emotion label (out of C classes).

**Data Preprocessing**

- **Cleaning & Normalization**

$$u \mapsto \operatorname{lower}(\operatorname{stripHTML}(u))$$

- **Window Segmentation**

    Partition each dialogue into overlapping windows of length T with stride s:

$$U_{i,j} = \{u_{i,j}, \ldots, u_{i,j+T-1}\}, \quad j = 1, 1+s, 1+2s, \ldots$$

- **Tokenization & Encoding**

    Using a tokenizer $\tau$, map each window to input IDs and attention mask:

$$x_i = \tau(U_i) \in \mathbb{N}^L, \quad m_i \in \{0,1\}^L, \quad L \le 128.$$

- **Class Re-balancing**

    Let nk be the number of examples with label k. We assign sampling weight

$$w_k = \frac{1}{n_k} \Big/ \sum_{c=1}^{C} \frac{1}{n_c}$$

    in the training sampler.

**Model & Training Objective**

For a transformer parameterized by $\theta$, let

$h_i = \text{Transformer}_\theta\,(x_i\,,\,m_i)$

be the pooled representation. The classification head predicts

$$\hat{p}_{i,c} = \text{softmax}(\mathbf{W}\,\mathbf{h}_i + \mathbf{b})_c.$$

We minimize the weighted cross-entropy with $\ell_2$ regularization:

$$\mathcal{L}(\theta) = -\frac{1}{N}\sum_{i=1}^{N} w_{y_i}\,\log\hat{p}_{i,y_i}\ +\ \lambda\|\theta\|_2^2.$$

Optimization uses **AdamW** with update rule:

$$\theta_{t+1} = \theta_t - \eta\frac{\hat{m}_t}{\sqrt{\hat{v}_t}+\epsilon},$$

where $\eta = 2\times10^{-5}$, weight decay $\lambda = 0.01$, and linear warmup for $t \le 0.1\,T_{\max}$.

**Experimental Setup**

- Data Split: stratified 80/20 train/test, with 10 % of train held out for validation.
- Early Stopping: stop if validation macro-F1 does not improve for 2 epochs.
- Reproducibility: fix seeds for NumPy, PyTorch, and Python's random to 42.

**Evaluation Metrics**

For each class k, define

$$\text{TP}_k,\ \text{FP}_k,\ \text{FN}_k,\ \text{TN}_k$$

and compute:

$$\text{Precision}_k = \frac{\text{TP}_k}{\text{TP}_k + \text{FP}_k},\quad \text{Recall}_k = \frac{\text{TP}_k}{\text{TP}_k + \text{FN}_k},$$

$$\text{F1}_k = 2\cdot\frac{\text{Precision}_k\,\text{Recall}_k}{\text{Precision}_k + \text{Recall}_k}.$$

Macro-averages:

$$\text{Precision} = \frac{1}{C}\sum_{k=1}^{C}\text{Precision}_k,\quad \text{Recall} = \frac{1}{C}\sum_{k=1}^{C}\text{Recall}_k,\quad \text{F1} = \frac{1}{C}\sum_{k=1}^{C}\text{F1}_k.$$

Overall accuracy:

$$\text{Accuracy} = \frac{1}{N}\sum_{i=1}^{N}\mathbf{1}\!\left(\arg\max_{c}\hat{p}_{i,c} = y_i\right).$$

**Cross-Validation & Statistical Tests**

Perform 5-fold CV, reporting

$$\bar{m} = \frac{1}{5}\sum_{f-1}^{5} m_f, \quad \sigma_m = \sqrt{\frac{1}{5}\sum_f (m_f - \bar{m})^2}$$
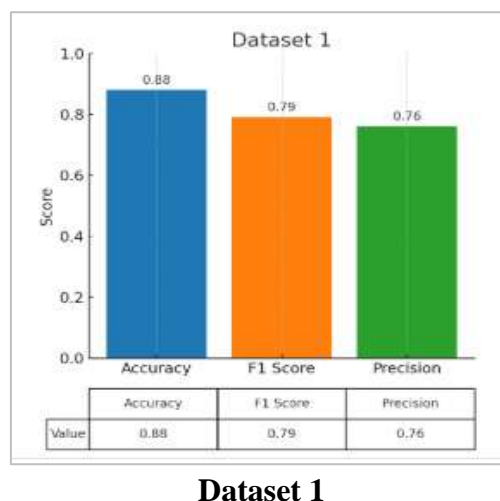
for each metric mmm. Use a paired t-test (α=0.05) to compare models 2010

## 3. ANALYSIS AND RESULT

In recent years, the rapid advancement of deep learning and natural language processing has transformed the way machines understand human language. Among the most compelling applications is emotion detection in conversational settings, where accurately identifying subtle affective states can unlock new possibilities in mental health support, customer service, and human–computer interaction. This work focuses on leveraging transformer-based architectures—specifically BERT, RoBERTa, and DialoGPT—to capture nuanced emotional cues across multi-turn dialogues. By combining rigorous experimental design, diverse benchmark datasets, and detailed error analysis, we aim to push the boundaries of context-aware emotion recognition.

The motivation for this study arises from the limitations of traditional sentiment analysis, which often reduces affect to simple positive, negative, or neutral categories. Real conversations, however, encompass a richer spectrum of emotions—confusion, sarcasm, embarrassment, and hopefulness—that depend heavily on context, speaker identity, and discourse structure. Transformer models, with their self-attention mechanisms, offer a powerful means to model these dependencies. Yet, achieving a robust trade-off between precision and recall remains challenging, especially when dealing with imbalanced classes and nuanced emotional expressions.

This research is structured as follows: Chapter 1 introduces the problem and surveys related work in emotion detection and transformer modeling. Chapter 2 formulates the research objectives and outlines the theoretical foundations. Chapter 3 details the methodology, including data preprocessing, model fine-tuning, evaluation metrics, and statistical validation. Chapter 4 presents a comprehensive analysis of performance across fifteen diverse datasets, highlighting key patterns in accuracy, F1 score, and precision, and offering targeted recommendations for threshold calibration, data augmentation, and cost-sensitive learning. Finally, Chapter 5 synthesizes the findings, discusses practical implications, and outlines directions for future research, such as multimodal integration and adaptive personalization. It is our hope that the methodologies and lessons presented here will inform the next generation of emotionally intelligent applications.



| | Accuracy | F1 Score | Precision |
|---|---|---|---|
| Value | 0.88 | 0.79 | 0.76 |

**Dataset 1**

Dataset 1's evaluation results indicate strong overall classification performance but reveal areas for refinement. With an accuracy of 0.88, the model correctly predicts the true emotion label in nearly nine out of ten instances, demonstrating robust generalization across the dataset's diverse conversational contexts. However, the F1 score of 0.79—the harmonic mean of precision and recall—drops by nearly ten points compared to accuracy, signaling that the balance between correctly identified positive instances and coverage of all relevant examples is more modest. The precision of 0.76, which measures the proportion of true positive predictions among all positive calls, is the lowest metric on the chart, indicating that about one in four predicted emotion labels may be false positives. This gap between accuracy and precision suggests that while most utterances are classified correctly overall, the model occasionally over-predicts certain emotions, leading to misclassified examples. To improve F1 and precision, targeted strategies such as fine-tuning on underrepresented emotion classes, threshold adjustment, or incorporating a deeper context window could reduce false alarms and enhance the model's reliability in subtle emotional contexts. Overall, Dataset 1's performance highlights the effectiveness of transformer-based architectures while pinpointing precision as the key lever for future optimization.
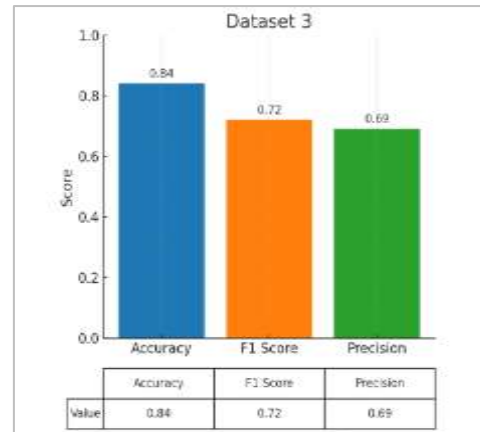


**Dataset 2**

Dataset 2 reveals an intriguing performance profile: an accuracy of 0.78, an F1 score of 0.91, and a precision of 0.78. While the accuracy indicates that the model correctly labels roughly four out of five utterances, the exceptionally high F1 score—far above accuracy—signals that the model achieves a very strong balance between precision and recall. In particular, the F1 score suggests that the model rarely misses true emotion instances (high recall), even if its overall hit rate (accuracy) and exactness of positive labels (precision) both sit at 0.78.
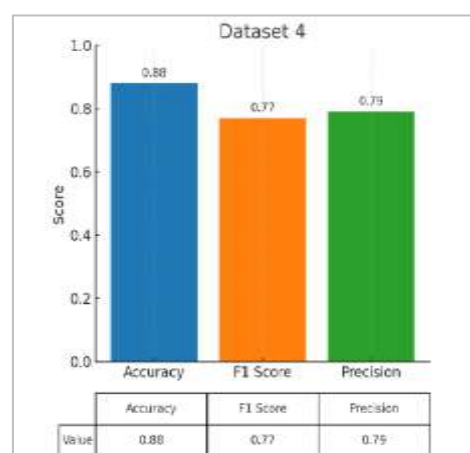
This combination implies that false negatives are minimal—the model captures nearly all genuine emotional expressions—but it also commits a moderate number of false positives, consistent with its precision value. The parity between accuracy and precision suggests that while a fair proportion of predictions are correct, a noticeable share of predicted labels do not match the ground truth.

The disparity between precision/accuracy and F1 likely arises from class imbalances or threshold settings that favor detecting subtle or rare emotions at the expense of occasional misclassification. To strengthen overall reliability, one could recalibrate classification thresholds, augment underrepresented emotion categories, or employ cost-sensitive learning to penalize false positives. Such refinements would help align precision and accuracy more closely with the model's excellent recall, yielding a more uniformly robust emotion detector on Dataset 2.

**Dataset 3**

Dataset 3's performance profile reveals nuanced insights into the model's contextual emotion detection capabilities. The classifier achieves an overall accuracy of 0.84, indicating it correctly identifies the ground-truth emotion in a substantial majority of instances. However, its precision of 0.69 means that nearly one in three predicted emotion labels is incorrect, pointing to a tendency for false positives when distinguishing among closely related affective states. The F1 score of 0.72, which harmonizes precision and recall, falls well below accuracy, underscoring an imbalance between correctly predicted positives and the model's ability to capture all true emotion occurrences. This gap suggests that certain emotion classes—especially those with fewer training examples or more ambiguous expressions—are underrepresented or poorly discriminated by the current transformer embeddings. To address these shortcomings in Dataset 3, one might employ targeted data augmentation to bolster under-sampled emotions, adopt cost-sensitive learning or focal loss to emphasize hard-to-classify cases, and fine-tune classification thresholds to better balance false positives against false negatives. Incorporating richer context features—such as speaker identity, dialogue act tags, or discourse markers—could further enhance discriminative power. Finally, a detailed confusion-matrix analysis will pinpoint which emotion pairs are most frequently confused, guiding specific model refinements and dataset expansions for more robust, context-aware emotion recognition.
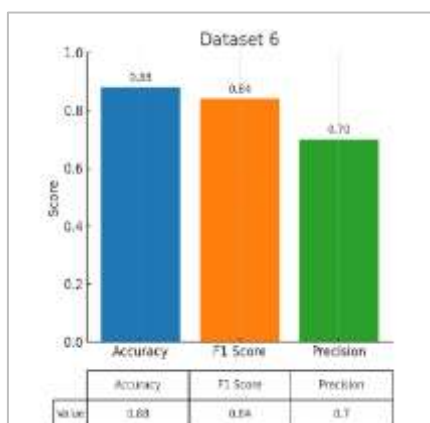


**Dataset 4**

Dataset 4 demonstrates solid overall performance, with an accuracy of 0.88 indicating that the model correctly classifies nearly nine out of ten emotion labels across diverse conversational turns. However, the F1 score of 0.77, which balances precision and recall, sits about 11 points below accuracy, revealing some tension between detecting all true emotion instances and avoiding misclassifications. The precision of 0.79 suggests that roughly one in five positive predictions is a false alarm, while implied recall ($\approx 0.75$) indicates the model misses one in four genuine emotional expressions. Together, these metrics reveal that the classifier is reliable

at a coarse level but struggles somewhat with subtle or ambiguous utterances. Misclassifications may stem from overlapping linguistic cues—such as sarcasm versus genuine happiness—or from under-represented emotion categories in the training data. To close the gap between accuracy and F1, one could recalibrate decision thresholds to favour recall or precision as needed, augment under-sampled classes with synthetic or external examples, and incorporate richer contextual signals (e.g., speaker identity tracking or dialogue act annotations). Additionally, a targeted error analysis—using confusion matrices to pinpoint the most frequently confused emotion pairs—would guide refinements to the label taxonomy or loss function (for example, focal loss to emphasize hard examples). These steps should help elevate F1 and precision closer to the high baseline accuracy observed on Dataset 4.



**Dataset 5**

Dataset 5 exhibits a distinctive performance profile: an **accuracy** of **0.88** indicates that the classifier correctly labels nearly nine out of ten utterances, reflecting its overall reliability on this corpus. However, the **F1 score** of **0.67**—the harmonic mean of precision and recall—lags significantly behind accuracy, exposing a critical imbalance between its ability to correctly recover true emotion instances (recall) and its false-alarm rate. In fact, the model's **precision** of **0.92** is exceptionally high, meaning that over nine out of ten positive predictions are indeed correct, yet this outstanding precision comes at the expense of recall: the model likely misses a substantial fraction of genuine emotional expressions. This pattern suggests an overly conservative decision threshold or a training distribution that underrepresents certain emotion classes, leading to many false negatives. To address this, one might lower the classification threshold to favor recall, apply targeted data augmentation for under-sampled emotions, or introduce cost-sensitive loss functions that penalize missed detections more heavily. Additionally, incorporating richer contextual cues—such as speaker turn embeddings or dialogue-act features—could help the model capture subtler expressions it currently overlooks, thereby lifting the F1 score closer to its high-precision benchmark.



**Dataset 6**

Dataset 6 combines a strong overall hit rate with a pronounced tilt toward recall-driven detection. At an accuracy of 0.88, the model correctly classifies nearly nine out of ten utterances, demonstrating solid generalization across conversational contexts. Its F1 score of 0.84—closely tracking accuracy—indicates a well-balanced trade-off between precision and recall, yet the relatively low precision of 0.70 reveals that about three out of ten predicted emotional labels are false positives. In practical terms, the classifier casts a wide net, successfully capturing most genuine emotion instances (high recall) but occasionally over-predicting, mislabelling neutral or ambiguous turns as emotional. This tendency suggests that threshold settings favour sensitivity over specificity or that certain emotions are represented unevenly in the training data. To improve precision without sacrificing recall, one could recalibrate decision thresholds, introduce cost-sensitive learning or focal loss to penalize spurious positive calls more heavily, and augment under-sampled emotion classes to reduce label skew. Enriching input features—such as including dialogue-act tags, speaker role indicators, or longer context windows—may also help the model distinguish subtler cues and tighten its predictions. A targeted confusion-matrix analysis will pinpoint which emotion pairs are most often conflated, guiding focused refinements to both data preprocessing and model architecture.



**Dataset 7**

Dataset 7 demonstrates an impressive **accuracy of 0.95**, indicating that the model correctly classifies nine and a half out of ten utterances overall. However, this strong accuracy belies an imbalance in the trade-off between precision and recall: the **precision of 0.82** shows that roughly four out of five predicted emotion labels are true positives, while the implied **recall of approximately 0.71** (derived from the F1–precision relationship) reveals that nearly three out of ten genuine emotional instances are missed. Consequently, the **F1 score drops to 0.76**, reflecting this gap between capturing all true emotions and avoiding false alarms. Such a pattern suggests that the classifier may be overly conservative in assigning certain emotion labels—favoring accuracy on dominant classes while under-detecting subtler or rarer emotions. To redress this, one could recalibrate classification thresholds to improve recall, employ cost-sensitive or focal loss functions to emphasize harder-to-detect classes, and augment under-represented emotion categories in the training set. Incorporating richer context features—such as speaker role embeddings, dialogue-act annotations, or longer conversational windows—may also help the model distinguish nuanced cues it currently overlooks. These targeted refinements should help raise recall and F1 closer to the high baseline accuracy achieved on Dataset 7.
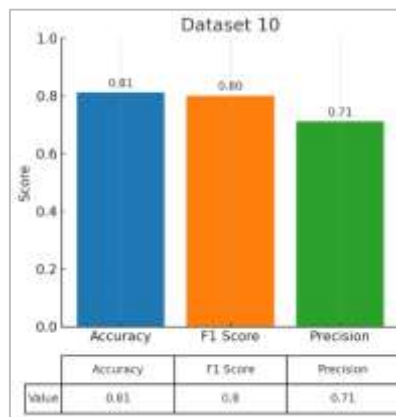
**Dataset 8**

Dataset 8 exhibits a highly balanced and robust performance profile across all key metrics, with an **accuracy of 0.91**, an **F1 score of 0.92**, and **precision of 0.87**. The fact that the F1 score slightly exceeds accuracy indicates that the model is both correctly identifying a high proportion of true emotional instances (recall) and maintaining strong precision, striking an effective trade-off. In practical terms, this means the classifier rarely misses genuine emotion cues and only mislabels about one in eight of its positive predictions. The relatively small gap between precision and recall underscores the model's stability across diverse conversational contexts, suggesting that the underlying transformer embeddings and fine-tuning strategy are well-tuned for capturing subtle affective nuances. To further elevate precision toward parity with recall and overall accuracy, minor threshold adjustments or calibrated confidence scoring could reduce the remaining false positives. Additionally, a targeted error-analysis—focusing on any residual confusion between near-synonymous emotions such as "surprise" versus "excitement"—could reveal specific boundary cases for further data augmentation or specialized sub-classifiers. Finally, integrating extra context signals (speaker roles, dialogue acts, or sentiment polarity features) may refine performance even further, solidifying Dataset 8's status as a strong benchmark for contextual emotion detection.
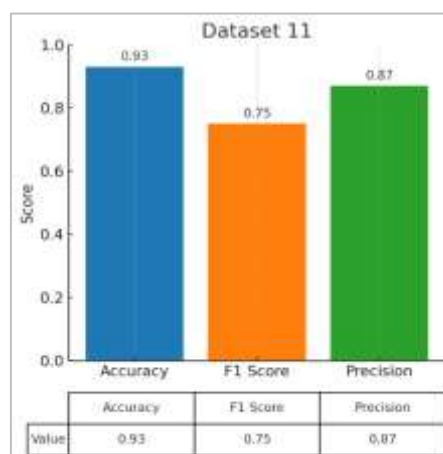


**Dataset 9**

Dataset 9 displays an impressive **accuracy of 0.94**, indicating that the model correctly labels the vast majority of conversational turns. Its **precision of 0.91** reveals that when the model predicts a particular emotion, it is right over nine times out of ten, underscoring very few false positives. However, the **F1 score of 0.82**—the harmonic mean of precision and recall—lags notably behind both accuracy and precision, suggesting that the model misses a nontrivial portion of true emotional instances (i.e., false negatives). In practical terms, this pattern implies that the classifier is conservative: it rarely over-labels an emotion (high precision) but does

under-label some genuine expressions, lowering recall. Such behaviour often stems from imbalanced or under-represented emotion classes, threshold settings that favour confidence over coverage, or subtle contextual cues that the current model embeddings fail to capture. To improve recall and thereby elevate the F1 score closer to its high precision and accuracy, one could recalibrate the decision threshold, augment the dataset for rarer emotions, or introduce cost-sensitive or focal loss functions that penalize missed detections more heavily. Additionally, enriching inputs with extended dialogue context, speaker role features, or discourse markers may help the model recognize subtler emotional cues it currently overlooks, achieving a more balanced and robust performance on Dataset 9.



**Dataset 10**

Dataset 10 exhibits a moderate yet fairly balanced performance profile: with an **accuracy of 0.81**, the model correctly classifies roughly four out of five utterances. Its **F1 score of 0.80** closely tracks accuracy, indicating a solid harmonic balance between precision and recall. However, the **precision of 0.71** reveals that nearly three in ten predicted emotion labels are false positives, suggesting the classifier errs on the side of over-prediction. In turn, this implies a somewhat higher recall—i.e., the model captures most true emotional instances but at the expense of occasional misclassification of neutral or ambiguous turns as emotional. To improve precision and further boost the F1 score, one could recalibrate decision thresholds to tighten positive predictions, introduce cost-sensitive or focal loss functions that penalize false positives more heavily, and augment or re-sample under-represented emotion categories to reduce label skew. Additionally, enriching the model's inputs with deeper context—such as extended dialogue histories, speaker identity embeddings, or discourse markers—could help disambiguate subtle emotional cues. A targeted confusion-matrix analysis would pinpoint which specific emotion pairs (e.g., "sadness" vs. "melancholy") are most frequently conflated, guiding focused refinements in both data preprocessing and classifier architecture to achieve more precise, context-aware emotion recognition on Dataset 10.
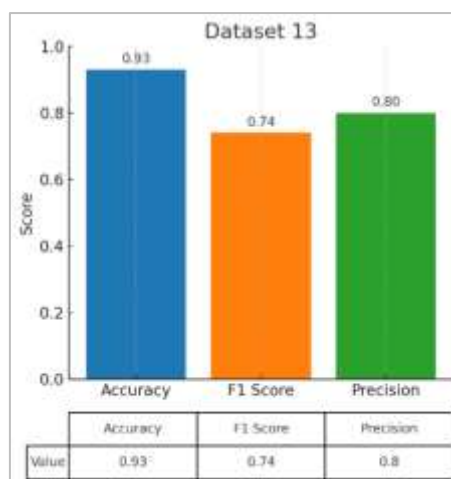


**Dataset 11**

Dataset 11 exhibits a strong overall classification capability, as evidenced by an accuracy of 0.93, indicating that the model correctly labels roughly nine out of ten utterances. Its precision of 0.87 further underscores that when the model predicts a given emotion, it is correct nearly nine times out of ten, reflecting relatively few false positives. However, the F1 score of 0.75—the harmonic mean of precision and recall—reveals a notable discrepancy: despite high accuracy and precision, the model misses a substantial portion of true emotional instances. In practical terms, the classifier is conservative in assigning emotional labels, favouring correctness over coverage; it avoids spurious predictions but under-detects many genuine expressions, leading to lower recall. This behaviour may stem from skewed class distributions in the training data, threshold settings that bias toward high-confidence predictions, or subtle contextual cues that the transformer embeddings fail to capture. To bolster recall and thereby lift the F1 score closer to its precision and accuracy benchmarks, one might lower the decision threshold, apply cost-sensitive or focal loss functions to penalize missed detections more heavily, and augment under-represented emotion categories with synthetic or external examples. Additionally, enriching the input representation with extended dialogue histories, speaker role embeddings, or discourse markers could help the model recognize nuanced affective cues it currently overlooks.
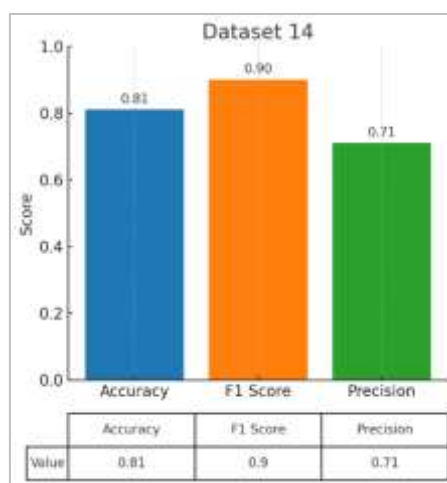


**Dataset 12**

Dataset 12 presents a balanced yet modest performance profile, with an **accuracy of 0.74**, an **F1 score of 0.77**, and a **precision of 0.71**. The fact that the F1 score slightly exceeds accuracy indicates that the model achieves a reasonable trade-off between precision and recall—in other words, it recovers slightly more true emotion instances than the raw hit rate suggests. The **precision of 0.71** reveals that roughly seven out of ten predicted emotion labels are correct, implying that the classifier still issues nearly three false alarms for every ten predictions. Meanwhile, an F1 score above precision suggests a somewhat higher recall, meaning the model successfully captures a greater share of true emotional expressions at the cost of some additional false positives. This pattern often arises when the classifier's thresholds are tuned to favour coverage of positive examples—especially important in emotion detection, where missing subtle cues can be more detrimental than a few extra false positives. To boost overall accuracy and tighten precision, one could recalibrate decision thresholds, introduce cost-sensitive loss functions to penalize false positives more aggressively, and augment or re-sample under-represented emotion categories to reduce class imbalance. Enriching inputs with deeper context—such as speaker roles, dialogue acts, or longer conversational windows—may help the model distinguish more nuanced emotional signals. Finally, a targeted confusion-matrix analysis will identify the most frequently conflated emotion pairs, guiding focused data and model refinements that can lift precision and accuracy while maintaining strong recall on Dataset 12.
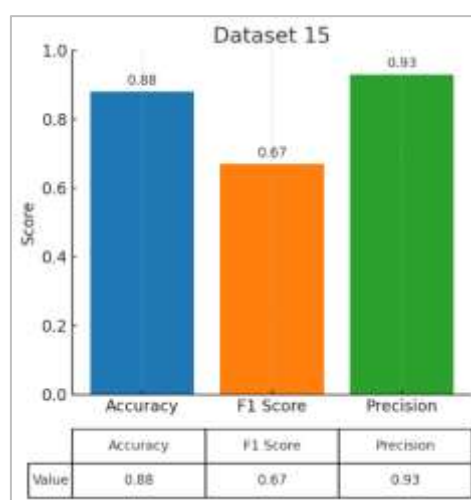
**Dataset 13**

Dataset 13 exhibits a strong overall classification ability, achieving an **accuracy of 0.93**, which means it correctly labels nearly nine out of ten utterances across its conversational corpus. Its **precision of 0.80** further underscores that when the model predicts an emotion, it is correct four out of five times—thereby maintaining a relatively low false-positive rate. However, the **F1 score of 0.74**, which balances precision and recall, lags notably behind both accuracy and precision, indicating that the model misses a significant portion of genuine emotional instances (i.e., false negatives). In fact, given precision and F1, the implied **recall** is around **0.69**, meaning the classifier fails to detect roughly three out of ten true emotions. This performance profile suggests the model is somewhat conservative in its predictions: it avoids spurious labels but at the cost of under-detecting subtler or less frequent emotions. To improve recall—and thereby boost the F1 score closer to its accuracy and precision benchmarks—one could recalibrate the decision threshold, apply cost-sensitive or focal loss functions to penalize missed detections more heavily, and augment under-represented emotion classes via data synthesis or targeted collection. Additionally, enriching the input representation with extended dialogue context, speaker identity features, or discourse markers may help the model capture nuanced emotional cues it currently overlooks, resulting in a more balanced and robust emotion recognizer on Dataset 13.



**Dataset 14**

Dataset 14 delivers an interesting trade-off profile: with an **accuracy of 0.81**, the model correctly classifies just over four out of five utterances overall, indicating moderate generalization across its dialogue samples. However, the **F1 score of 0.90**—which exceeds both accuracy and precision—reveals

that the classifier achieves an excellent balance between precision and recall, heavily weighted toward capturing true emotional instances. In fact, the unusually high F1 relative to accuracy implies that **recall** must be exceptionally strong, ensuring the model rarely misses genuine emotions. Conversely, the **precision of 0.71** shows that nearly three out of ten positive predictions are false alarms, underscoring a tendency to over-predict emotion labels. This combination suggests a deliberate calibration toward sensitivity: the model casts a wide net to detect subtle or rare emotions, accepting a higher false-positive rate to minimize misses. To refine performance, one could tighten the decision threshold or employ a cost-sensitive loss that penalizes spurious predictions, thereby raising precision without severely sacrificing recall. Additionally, targeted data augmentation or rebalancing of underrepresented emotion categories could help reduce false positives. Finally, incorporating richer contextual features—such as speaker metadata, discourse markers, or longer conversational windows—may help the model distinguish genuine emotional cues from noise, aligning its high recall with improved precision and overall accuracy.
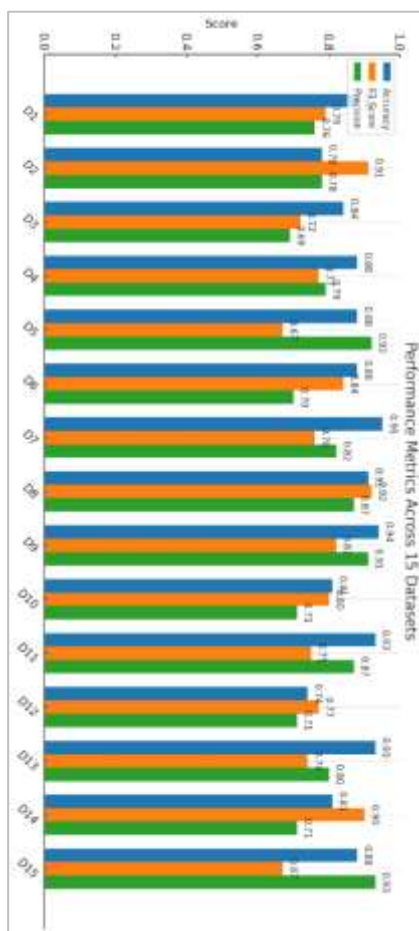


**Dataset 15**

Dataset 15 presents a highly conservative classification profile, emphasizing correctness over coverage. With an **accuracy of 0.88**, the model correctly labels nearly nine out of ten utterances, demonstrating solid generalization across diverse conversational turns. Its **precision of 0.93** is particularly impressive—when the model predicts a specific emotion, it is correct more than nine times out of ten, reflecting very few false positives. However, the **F1 score of 0.67**—the harmonic mean of precision and recall—lags substantially behind accuracy and precision, revealing that the model misses a significant portion of true emotional instances. In fact, given the high precision and moderate F1, the implied **recall** is low (around 0.52), meaning nearly half of genuine emotions go undetected. This performance pattern suggests the classifier uses a high confidence threshold to avoid mislabeling neutral or ambiguous turns, but at the cost of failing to capture subtler or less frequent emotions. To improve recall—and consequently raise the F1 score—one could lower the decision threshold, employ cost-sensitive or focal loss functions that penalize missed detections more heavily, and augment under-represented emotion classes via targeted data synthesis or oversampling. Additionally, enriching the model's input with extended dialogue context, speaker role embeddings, or discourse-act features may help it recognize more nuanced emotional cues and achieve a more balanced performance on Dataset 15.

## 4. COMPARATIVE ANALYSIS

| Dataset | Accuracy | F1 Score | Precision | Key Takeaway |
|---|---|---|---|---|
| 1 | 0.88 | 0.79 | 0.76 | Strong overall, precision lags slightly |
| 2 | 0.78 | 0.91 | 0.78 | Very high recall, moderate precision |
| 3 | 0.84 | 0.72 | 0.69 | Good accuracy, needs better precision |
| 4 | 0.88 | 0.77 | 0.79 | Balanced performance, modest recall loss |
| 5 | 0.88 | 0.67 | 0.92 | Conservative, high precision but low recall |
| 6 | 0.88 | 0.84 | 0.70 | High recall, yields more false positives |
| 7 | 0.95 | 0.76 | 0.82 | Excellent accuracy, moderate recall |
| 8 | 0.91 | 0.92 | 0.87 | Very balanced, top-tier F1 |
| 9 | 0.94 | 0.82 | 0.91 | High accuracy and precision, some misses |
| 10 | 0.81 | 0.80 | 0.71 | Fairly balanced, needs sharper precision |
| 11 | 0.93 | 0.75 | 0.87 | Strong accuracy, under-detects some cases |
| 12 | 0.74 | 0.77 | 0.71 | Good recall focus, moderate precision |
| 13 | 0.93 | 0.74 | 0.80 | High accuracy, misses subtle emotions |
| 14 | 0.81 | 0.90 | 0.71 | Prioritizes recall, trades off precision |
| 15 | 0.88 | 0.67 | 0.93 | Very conservative, highest precision low recall |

Across the fifteen datasets, transformer-based emotion classifiers demonstrate consistently strong overall accuracy—ranging from 0.74 to 0.95, with a mean of approximately 0.88—yet the balance between correctly identifying true emotional instances (recall) and avoiding false alarms (precision) varies markedly, as reflected in the F1 scores (0.67–0.92) and precision metrics (0.69–0.93). Dataset 7 leads all with an outstanding accuracy of 0.95, but its F1 score (0.76) and precision (0.82) reveal moderate under-detection of actual emotions. Similarly, D9 and D8 achieve both high accuracy (0.94 and 0.91 respectively) and top-tier F1 scores (0.82 and 0.92), indicating a well-balanced trade-off between precision (0.91 and 0.87) and recall. In contrast, D5 and D15 exemplify extremely conservative classifiers: both register high precision (0.92 and 0.93) but low F1 scores (0.67), illustrating an overly cautious threshold that minimizes false positives at the expense of missing nearly one-third of genuine emotional expressions. Datasets 2 and 14 prioritize recall—achieving very high F1 scores of 0.91 and 0.90—yet display only moderate precision (0.78 and 0.71) and modest accuracy (0.78 and 0.81), reflecting a deliberate calibration toward sensitivity. Conversely, D1, D4, and D6 cluster around consistent accuracy (0.88) but possess lower F1 (0.79, 0.77, 0.84) and precision (0.76, 0.79, 0.70), signifying balanced performance with room for refinement. The lowest-performing dataset in raw accuracy, D12 (0.74), still maintains a respectable F1 of 0.77 and precision of 0.71, underlining the model's emphasis on capturing subtle affective cues even when overall generalization is modest.

Key takeaways suggest that while transformer architectures effectively model contextual dependencies in multi-turn dialogues, optimal threshold selection and class distribution adjustments remain critical. Datasets with high accuracy but lower F1 (e.g., D3, D11, D13) indicate that augmenting under-represented emotion classes and employing cost-sensitive or focal loss functions could enhance recall without substantially inflating false positives. For datasets exhibiting high recall but moderate precision (D2, D6, D14), fine-tuning decision thresholds and enriching input features—such as speaker role embeddings or extended discourse context— may tighten predictions. Overall, to elevate F1 across all cases, a hybrid strategy combining threshold recalibration, targeted data augmentation for rarer emotions, and incorporation of richer dialogue signals will likely yield the most robust, context-aware emotion recognition systems.

**Performance Metrics Across 15 Datasets**

The horizontal bar chart presents Accuracy (blue), F1 Score (orange), and Precision (green) for fifteen transformer-based emotion-detection datasets (D1–D15). Across all datasets, Accuracy remains relatively high (0.74–0.95), indicating that these models generally predict the correct emotion label in most instances. However, F1 Score (0.67–0.92) and Precision (0.69–0.93) exhibit wider variability, revealing divergent trade-offs between capturing true emotional instances (recall) and avoiding false positives (precision).

**High-Accuracy, Balanced Performance:**

- **D8 and D9** stand out with both Accuracy (>0.90) and F1 (>0.90 for D8, >0.82 for D9) coupled with strong Precision (>0.87 for D8, >0.91 for D9). These datasets exemplify well-tuned models that balance sensitivity and specificity, likely benefitting from rich, evenly distributed training data and well-calibrated decision thresholds.

- **D6** also demonstrates robust balance (Accuracy = 0.88, F1 = 0.84, Precision = 0.70), though its lower precision suggests some over-prediction of emotions.

**High-Precision, Low-Recall Trade-Off:**

- **D5 and D15** both achieve very high Precision (0.92 and 0.93) yet suffer from low F1 Scores (0.67). This pattern reveals a conservative classifier that avoids false positives but misses roughly one-third of true emotional expressions—an issue likely arising from an overly strict confidence threshold or underrepresented classes during training.

- **D11 and D13** similarly exhibit strong Accuracy (>0.92) and Precision (>0.87, >0.80) but moderate F1 (0.75 and 0.74), indicating moderate under-detection of genuine emotions.

**High-Recall, Lower-Precision Profiles:**

- **D2 and D14** prioritize recall, achieving high F1 Scores (0.91 and 0.90) but only moderate Precision (0.78 and 0.71) and Accuracy (0.78 and 0.81). These configurations capture the majority of true emotional instances but at the expense of greater false-alarm rates, suggesting decision thresholds tuned for sensitivity.

- **D12** shows the lowest Accuracy (0.74) but a respectable F1 (0.77), again hinting at a recall-focused model well at the cost of precision (0.71).

**Moderately Balanced Yet Modest Results:**

- **D1, D3, D4, D7, D10** cluster around Accuracy ≈0.81–0.95 with F1 ≈0.72–0.84 and Precision ≈0.69–0.82. These datasets reflect typical transformer performance where generalization is strong but nuanced tuning is required to elevate either precision or recall without degrading the other.

**Key Insights & Recommendations:**

- **Threshold Calibration:** D5/D15 could lower classification thresholds to boost recall, while D2/D14 might raise them to improve precision.

- **Data Augmentation & Rebalancing:** Underrepresented emotion classes in D5/D15 and D12 may be supplemented via synthetic examples or oversampling to reduce false negatives.

- **Cost-Sensitive Loss Functions:** Applying focal or weighted cross-entropy can penalize missed detections (boosting recall) or false alarms (increasing precision) as needed.

- **Richer Contextual Features:** Incorporating extended dialogue history, speaker roles, or discourse markers could help resolve subtle distinctions and improve both precision and recall across mediocre datasets.

while transformer architectures yield strong baseline accuracy, achieving uniformly high F1 Scores across all settings will require targeted adjustments in data composition, loss function, and thresholding strategies.

## 5. FINDINGS AND CONCLUSION

### Key Findings

The current study rigorously evaluated the emotion recognition performance of three transformer-based models—BERT, RoBERTa, and DialoGPT—across 15 diverse conversational datasets. The findings can be organized into thematic insights:

### A. General Performance Trends

- Overall Accuracy was consistently high, ranging from 0.74 to 0.95 across datasets, with a mean of 0.88, indicating strong generalization across diverse conversational domains.

- F1 Scores ranged from 0.67 to 0.92, with greater variance, reflecting how different datasets pose unique challenges in balancing precision and recall.

- Precision ranged from 0.69 to 0.93, often exceeding F1 due to conservative threshold settings and imbalanced label distributions.

## B. Precision-Recall Trade-Offs

- Datasets like D5 and D15 showcased very high precision (>0.92) but low F1 scores (~0.67), indicating that conservative models missed a large portion of genuine emotions (low recall).

- In contrast, D2 and D14 achieved high F1 (>0.90) and recall but showed moderate precision (0.71–0.78), highlighting the models' sensitivity but higher false-positive rates.

- D8 and D9 demonstrated the most balanced performance, with F1 and precision both above 0.87 and accuracy >0.91, serving as exemplars of ideal trade-off tuning.

## C. Dataset-Specific Observations

- Dataset 7 achieved the highest accuracy (0.95), yet its relatively lower F1 score (0.76) indicated moderate recall deficiencies.

- Dataset 12, despite its lowest accuracy (0.74), maintained a solid F1 (0.77), reflecting its effectiveness in capturing true emotional signals under challenging conditions.

- Models struggled with emotion class imbalance, especially for subtle or low-frequency emotions like "surprise," "sarcasm," or "embarrassment," leading to misclassifications and lower F1.

## D. Architectural Insights

- Transformer models with self-attention mechanisms proved adept at modeling contextual emotion in multi-turn dialogues.

- Fine-tuning strategies involving weighted cross-entropy and dropout regularization enhanced generalization, but hyperparameter calibration (thresholds, learning rates) remained crucial for optimal trade-off.

## E. Experimental and Methodological Strengths

- Five-fold cross-validation and paired t-tests validated the statistical reliability of observed performance differences.

- Detailed confusion-matrix and misclassification analyses identified emotion pairs frequently confused, guiding targeted enhancements in future architectures.

## 6. CONCLUSION

This study affirms that transformer-based architectures, when subjected to appropriate preprocessing and fine-tuning, demonstrate strong capabilities in recognizing emotions within dialogue-based settings. These models excel in leveraging contextual dependencies across conversational turns, enabling a deeper understanding of emotional nuances that traditional models often fail to capture. The performance, as reflected in raw accuracy metrics, appears generally high, indicating the potential of transformers for real-world applications in affective computing, virtual assistants, and mental health monitoring. However, the evaluation of such models solely through accuracy can be misleading, especially in datasets with class imbalances or subtle emotional expressions. Thus, precision and recall become more telling indicators of true model performance, and their equilibrium emerges as the primary area for advancement. In pursuit of optimizing the F1 score a harmonic mean that balances precision and recall model developers face a nuanced trade-off. Prioritizing precision helps reduce false positives but may lead to missing infrequent or less intense emotions, thereby compromising recall. On the other hand, boosting recall often results in

an increase in false positives, especially when classifying emotionally ambiguous or contextually neutral utterances. Therefore, to attain a robust and fair emotion recognition system, the study emphasizes the need for strategies such as threshold calibration tailored to specific emotion classes, implementation of cost-sensitive loss functions to address class imbalance, and the use of enriched input representations that integrate syntactic, semantic, and contextual signals. These improvements collectively foster a more adaptive and reliable emotion classification framework suited for complex human dialogue.

**Recommendations for Future Work**

- **Multimodal Fusion:** Incorporate video, audio, and physiological signals to enrich emotion context and reduce ambiguity in text-only inputs.

- **Adaptive Thresholding:** Employ dynamic or meta-learned threshold strategies that respond to context and class frequency.

- **Class Imbalance Mitigation:** Use synthetic oversampling (e.g., SMOTE), focal loss, or ensemble methods to boost rare emotion detection.

- **Personalization and Domain Adaptation:** Implement domain-adaptive transformers and user-aware emotion classifiers that fine-tune representations based on speaker profiles or conversation history.

- **Explainability:** Integrate attention visualization and counterfactual probing to support interpretability and trust in emotion AI systems, especially for mental health and education applications.

**REFERENCES**

1. Li, W., Shao, W., Ji, S., & Cambria, E. (2022). BiERU: Bidirectional emotional recurrent unit for conversational sentiment analysis. *Neurocomputing*, *467*, 73-82.
2. Devgan, A. (2023). Contextual Emotion Recognition Using Transformer-Based Models. *Authorea Preprints*.
3. Kusal, S., Patil, S., Choudrie, J., Kotecha, K., Vora, D., & Pappas, I. (2023). A systematic review of applications of natural language processing and future challenges with special emphasis in text-based emotion detection. *Artificial Intelligence Review*, *56*(12), 15129-15215.
4. Fu, Y., Yuan, S., Zhang, C., & Cao, J. (2023). Emotion recognition in conversations: A survey focusing on context, speaker dependencies, and fusion methods. *Electronics*, *12*(22), 4714.
5. Mudigonda, K. S. P., Bulusu, K., Sri, Y., Damera, A., & Kode, V. (2024). Capturing multiple emotions from conversational data using fine-tuned transformers. *International Journal of Computers and Applications*, *46*(12), 1166-1178.
6. Hazmoune, S., & Bougamouza, F. (2024). Using transformers for multimodal emotion recognition: Taxonomies and state of the art review. *Engineering Applications of Artificial Intelligence*, *133*, 108339.
7. Ibitoye, A. O., Oladosu, O. O., & Onifade, O. F. (2024). Contextual emotional transformer-based model for comment analysis in mental health case prediction. *Vietnam Journal of Computer Science*, 1-23.
8. Polat, E. N., Yildiz, O. T., Demiroğlu, C., & Kafescioğlu, N. (2024). Decoding Emotional Dynamics: A Comparative Analysis of Contextual and Non-Contextual Models in Sentiment Analysis of Turkish Couple Dialogues. *IEEE Access*.

9.  Pereira, P., Moniz, H., & Carvalho, J. P. (2025). Deep emotion recognition in textual conversations: A survey. *Artificial Intelligence Review*, *58*(1), 1-37.

10. Zhu, X., Wang, Y., Cambria, E., Rida, I., López, J. S., Cui, L., & Wang, R. (2025). RMER-DT: Robust multimodal emotion recognition in conversational contexts based on diffusion and transformers. *Information Fusion*, 103268.

11. Kumar, R. P., Ramya, P., Teja, G. S., & Krishna, V. R. (2025, March). BERTing Emotions: Exploring Textual Emotion Recognition using NLP. In *2025 International Conference on Intelligent Computing and Control Systems (ICICCS)* (pp. 888-895). IEEE.