# Comparative Analysis of Transformer Models for Multi-Class Text Classification

MSc Research Project
Programme Name

## Sharanya Neelakanti
Student ID: X23138220

School of Computing
National College of Ireland

Supervisor: Rejwanul Haque

# National College of Ireland

## MSc Project Submission Sheet

### School of Computing

| | |
|---|---|
| **Student Name:** | …….Sharanya Neelakanti………………………………………………………………………… |
| **Student ID:** | ……X23138220……………………………………………………………………..…… |
| **Programme:** | ……MSc in Artificial Intelligence…………………… **Year:** …2023-2024.. |
| **Module:** | ……MSc Research Practicum……………………………………………..……… |
| **Supervisor:** | ……Rejwanul Haque……………………………………………………..……… |
| **Submission Due Date:** | ……16 September 2024……………………………………………..……… |
| **Project Title:** | …Comparative Analysis of Transformer Models for Multi-Class Text Classification……………………………………………………………… |
| **Word Count:** | ……8250……………………… **Page Count**………24……………………………….…….. |

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

ALL internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

**Signature:** ……Sharanya Neelakanti………………………………………………………

**Date:** ……16 September 2024……………………………………………………

## PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST

| | |
|---|---|
| Attach a completed copy of this sheet to each project (including multiple copies) | □ |
| **Attach a Moodle submission receipt of the online project submission,** to each project (including multiple copies). | □ |
| **You must ensure that you retain a HARD COPY of the project,** both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer. | □ |

Assignments that are submitted to the Programme Coordinator Office must be placed into the assignment box located outside the office.

| Office Use Only | |
|---|---|
| Signature: | |
| Date: | |
| Penalty Applied (if applicable): | |

# Comparative Analysis of Transformer Models for Multi-Class Text Classification

Sharanya Neelakanti
X23138220

## Abstract

This paper makes a comparative evaluation of five state-of-the-art transformer models in multi-class emotion recognition: BERT, GPT-3.5, RoBERTa, XLNet, and DistilBERT. Motivated by the demand for detecting emotions with accuracy in so many applications today, this research aimed at comparing these models on accuracy, precision, recall, and F1 score on classifying texts into multiple categories of emotions.

The research employed the usage of the GoEmotions dataset, which is a dataset containing 58,000 Reddit comments with 27 different annotated emotions and consolidated into three major classes, i.e., positive, neutral, and negative. The methodology in this research undertook preprocessing for the dataset, model implementation, and fine-tuning, ending up at the point of developing a comprehensive evaluation framework.

Key findings were that there did exist a performance hierarchy and, quite unexpectedly, DistilBERT outstripped all larger models, scoring 95.88%. Following were RoBERTa, XLNet, BERT, and GPT-3.5, performing in descending order. For all models, in comparison to the neutral or negative ones, recognizing positive emotions was easier. A remarkable exception was GPT-3.5, which, though doing splendidly elsewhere in NLP applications, underperformed in the given task.

This paper aids in adding to this literature by disputing the commonly held belief that improvements in NLP tasks are made when the model's size is increased and focusing on the compression methods of models. The findings have implications for academic research in NLP and for the practical applications of emotion recognition systems, mainly scenarios related to high computational efficiency.

## 1 Introduction

NLP (Natural language processing) has witnessed great development lately, more specifically in text classification. A shift occurred with converting the transformer-based models in this domain, where performance reached previously unseen state-of-the-art levels in natural language understanding and classification. In text classification, multi-class text classification, more focused on emotion recognition (Ameer et al., 2023)., has been an important task in a wide range of applications from sentiment analysis through content categorization to emotional recognition systems.

Growing interest in the demand for effective systems for emotion recognition comes not only from social media analysis and customer feedback processors but also from monitors of mental health. Correct recognition of emotions yields insight into public views and actions, thereby influencing decision-making processes and allowing for personalized up-to-the-individual interventions for well-being (Guo et al., 2024)

Though transformer models have lately held great potential to perform a binary classification task, multi-class emotion recognition has turned out to be an area that warrants further studies in terms of their handling of nuanced emotional states and class imbalance. Application of models including BERT(Bidirectional Encoder Representations from Transformers), GPT( Generative Pre-trained Transformers), RoBERTa(Robustly Optimized BERT Approach), XLNet, and DistilBERT has been reported with success in various text classification tasks by studies (Adoma et al., 2020).Specifically, current literature lacks comparison studies that would look more in-depth at how these models work for multi-class emotion recognition with regard to issues such as class imbalance and domain-specific language.

The research tries to fill this gap by providing a comprehensive comparison among five state-of-the-art transformer models: BERT, GPT-3.5, RoBERTa, XLNet, and DistilBERT. Guiding the thrust of this study will be the primary research question:
To what extent do the BERT, GPT-3.5, RoBERTa, XLNet, and DistilBERT models differ in terms of accuracy, precision, recall, F1 score, and processing speed when classifying text into multiple emotion categories across various datasets and scenarios? How do they handle situations that are distinct to the domain and unequal power between classes?
The objectives of this study are:
- Implementation and tuning of each model on a common emotion recognition dataset.
- Evaluation of the models' performance using evaluation metrics such as accuracy, precision, recall, and F1 score.
- The performance of models concerning the handling capacity for class imbalance in emotion recognition tasks.

This will be achieved with extensive data preparation techniques, performance metrics, and benchmarking implemented. The models will be evaluated on the GoEmotions dataset, which is a large-scale corpus for emotion recognition that guarantees most emotional states to be covered. Each model will be fine-tuned based on this dataset, and its performance will strive to rigorously evaluate it with k-fold cross-validation.

Specifically, this research provides a comparative analysis of such transformer models over the task of multi-class emotion recognition. Through this, the paper seeks to elucidate the advantages and disadvantages of each model so that researchers and practitioners are better placed to choose the most suitable model for their respective emotion recognition tasks. These results will be of great value in the development of more accurate, efficient, and interpretable emotion recognition systems, possibly impacting several domains through better decision-making and interventions.
The structure for this report is as follows: Section 2 provides a Literature Review, which gives an overview of the current status of research in transformer models and multi-class text classification. Section 3 describes the methodology of the research: data preparation, model fine-tuning, and evaluation metrics. Section 4 presents design specifications and implementation. Section 5: Experimental Results and Comparative Analysis of Transformer Models. Section 6 concludes the document with a summary of key findings, discussion of limitations, and future research directions.

# 2 Related Work

This critical literature review looks at how this model of transformers evolved and then goes on to apply that in multi-class emotion recognition with sentiment analysis. The review shows that NLP (Natural language processing) is one of the fast-growing areas and further investment proves the potency of such models for handling very complex emotion recognition tasks.

## 2.1 Evolution of Transformer Models

The advent of transformer models marked a milestone in the history of NLP (Natural language processing), surpassing all others in the comprehension and classification of human languages. Devlin (2018) introduced BERT, which revolutionized models for language representation. What certainly keeps this batch different is its bidirectional training approach contrasted with previous unidirectional models. In the paper, the authors demonstrated the flexibility of BERT in achieving state-of-the-art results on 11 NLP tasks, where the modifications introduced to address each task were really minimal. This flexibility can be extrapolated into multiclass emotion recognition tasks. However, since this study does not delve more into the domain of emotion recognition, some scope remains for further investigation. While it is a strength of BERT to be bidirectional in understanding context, this is coupled with notable limitations with the computational intensity involved and possible overfitting on small datasets for application in emotion recognition.

An extension to BERT came in the form of RoBERTa, from (Liu et al.,2019), which is a corrected and almost fully optimized version that tackled some of BERT's failings. These studies found BERT to be hugely undertrained, with just the tuning of key hyperparameters and hence the training regime, to already give scores of RoBERTa, extending well over benchmark levels for a diverse set of NLP tasks. This aspect is especially important for the research of emotion recognition when a detailed understanding of context is required. One of the major strengths would be improved training methodology for RoBERTa, even though, for applications in real-world use for emotion recognition, increased computational requirements and potential challenges in fine-tuning with small datasets are critical issues.

Yang (2019) suggested the XLNet model to overcome the disadvantages in BERT's training objective. This is an obvious idea for permutation language modeling and, thus, XLNet can capture bidirectional context without the caveats of BERT's masked language modeling. The authors reported superior performance on many NLP tasks, and especially better gains are shown for tasks that have longer sequences. This could be more beneficial specifically for emotion recognition in complex scenarios and therefore rely more on context. Results for multi-class emotion classification were not given in the study. Another significant advantage of XLNet is that it can model long-range dependencies more effectively, which may increase the risk of increased training time and overfitting on small datasets, some common challenges for most emotion recognition tasks.

Addressing exactly this problem of model size and efficiency, (Sanh et al.,2019) developed DistilBERT a distilled version of BERT smaller by about 40% with 60% more speed while retaining 97% of the performance from BERT, clearly showing proof of the potential of construction for efficient emotion recognition systems that can be deployable on devices with low resources. The model size and performance trade-offs have not been explored in the

context of this study on emotion recognition. Much as this reduction in size with an improvement in inference speed makes DistilBERT very attractive for real-time applications of emotion recognition, there is a slight drop in performance compared with the full BERT, and limitations in terms of captured nuanced emotional contexts that accompany this are issues which still need further exploration.

A recent work by (Ye et al.,2023) benchmarked GPT-3 and GPT-3.5 models on a large number of diverse NLP tasks. Capability found with such models is pretty impressive, but it also turns out that performance does not improve linearly with model evolution, especially natural language understanding tasks. This raises questions about the efficacy of larger models for specific tasks like emotion recognition and emphasizes that models, regardless of size or general performance, have to be evaluated for a variety of tasks. As much as this work is strong in comprehensiveness, it leaves space to explore the domain of emotion recognition tasks.

## 2.2  Transformer Models in Sentiment Analysis and Emotion Recognition

The application of transformer models with respect to sentiment analysis, which is rather close to emotion recognition, has been the subject of several studies. In this regard, (Dhivyaa et al., 2023) used XLNet for sentiment analysis and revealed its high potential in capturing complex sentiments. Similarly, (Cai et al., 2021) presented a hybrid model for sentiment classification called BiLSTM-AT. These works have thus proved that transformer-based methods are rather effective for nuanced sentiment understanding. While this clearly points to a need for more comprehensive studies in the broader emotional spectrum, they focus on binary or limited-class sentiment analysis rather than multi-class emotion recognition.

Applications of the transformer models to specific domains shed more light into the potential capability of emotion recognition. Maruvur Selvi and Sreeja (2023) applied different models, including BERT, to Tamil literature for the task of sentiment analysis. Another state-of-the-artwork related to sentiment analysis is by (Arora et al.,2023), who used BERT for the analysis of IMDb movie reviews. These studies prove that considerations at the domain level are very vital in sentiment analysis and hence in emotion recognition. However, these works do not tackle the multi-class nature of emotion recognition and do not compare different transformer models; hence, these publications are only partly useful for our question.

Studies by (Susmitha et al.,2023) and Agrawal et al. applied traditional machine learning algorithms for sentiment analysis. While both works bring colossal value addition to literature on the task of sentiment classification, they lack advanced capabilities brought in by transformer models. This gap presents that potential gains could be expected in emotion recognition tasks using transformer-based approaches.

In the domain of conversational AI and intent classification, (Noorani et al.,2023) and (Hirway et al.,2023) have shown the applicability of transformer models. Noorani et al. came up with a framework for sentiment-aware chatbots, whereas (Hirway et al.,2023) compared the performance of GPT-Neo and GPT-2 for intent classification. These studies show how versatile the applications of the transformer models on related tasks are but do not particularly focus on multi-class emotion recognition.

The authors Zhang and Shafiq (2024), presented a very comprehensive survey on transformer models for a wide variety of NLP tasks, underlining the different strengths these models have, and proposing ensemble methods to use such strengths. This certainly provides very valuable

insight into the general landscape of the transformer models but does not get into specifics on multi-class emotion recognition. Although this survey represents the broad coverage of the topic, in-depth investigation into the specific area of emotion recognition shows a requirement for focused studies.

## 2.3   Research Gap and Justification

Therefore, there exists a serious gap regarding the cross-evaluation of the transformer models specifically for multi-class emotion recognition in the literature. Even though several works have proven their efficiency in areas similar to sentiment analysis, literature is limited when making a direct comparison of the base models, including BERT, RoBERTa, and GPT-3.5, concerning emotion recognition. The existing research has been limited, for the most part, to looking at the binary classification of sentiment or general tasks of NLP, ignoring the complexities that exist in multiclass emotion recognition. Furthermore, those surveys done for tasks linked with emotions almost always focus on just one model or a small variety of models, totally missing an overview of their performances for the different transformer architectures.

Literature shows a limited examination of the potential of these models with class imbalance and domain-specific language in tasks of emotion recognition. Somewhat underexplored is the computational efficiency and interpretability of these models in regard to emotion recognition. This gap, therefore, justifies the need for our research question, which was set out to provide a comprehensive comparison pertaining only to multi-class emotion recognition of these transformer models.

Such improved multi-class emotion recognition would have applications in improved analysis of social media, finer customer feedback processing, and advanced mental health monitoring systems. In so doing, the gap is filled by our work, offering valuable insight to academic researchers involved in NLP and industrial practitioners working on emotion recognition systems, with possible influences upon further developments in text classification applications and emotion recognition technologies.

The methodology that would be used in implementing and testing this research is quite rigorous because of the data preprocessing, model fine-tuning, and extensive evaluation with k-fold cross-validation. Standardized datasets like GoEmotions ensure a full test balanced for all emotional states. Evaluation metrics such as accuracy, precision, recall, and F1 score will be computed to measure each model's abilities. Moreover, class imbalance handling, computational efficiency, and interpretability of the models with respect to emotion recognition will also be evaluated.

This literature review sets the absolutely clear case for an exhaustive, comparative study of transformer models at the state-of-the-art level in the task of multi-class emotion recognition. That is why this proposed research will fill up the gap found in literature with important insights which will increase the theory of emotion recognition in NLP and its applications.

# 3   Research Methodology

This paper describes the research procedure and the methodology of evaluation that will be used to answer the research question on the comparisons made between BERT, GPT-3.5, RoBERTa, XLNet, and DistilBERT models for multi-class emotion recognition. The methodology is informed by established practices from related work and designed to assure a rigorous, reproducible scientific process.

## 3.1   Equipment and Computational Resources

All experiments were conducted on Google Cloud Platform using the following configuration:
- Machine type: n1-standard-8 (8 vCPUs (Central Processing Unit), 30 GB memory)
- GPU (Graphical Processing Unit): NVIDIA Tesla A100 (40GB RAM) (Random Access Memory)
- Storage: 100 GB SSD (solid-state drive)
- CUDA version: 11.2
- PyTorch version: 1.9.0
- Transformers library version: 4.9.2

This high-performance setup was chosen following the recommendations of Strubell et al. (2019) about environmental considerations in NLP research, ensuring that all the models are trained with consistency and efficiency.

## 3.2   Dataset Selection and Preprocessing

In this regard, (Demszyk et al.,2020) GoEmotions dataset was adopted, which contains 58,000 comments from Reddit, each annotated with 27 emotion categories. This is so because this dataset represents exhaustive coverage of emotional states and its applicability to multi-class classification that was elaborated in Zhang and Shafiq (2024).

The following preprocessing steps were taken in accordance with Liu et al. (2019):

- Text Cleaning: This stage removes URLs and special characters from the text, leaving behind alphanumeric characters and basic punctuation. This step may help to further reduce noise in the data.
  - Example: Input: "Check out https://example.com! It's great :)" Output: "Check out Its great"
- Lowercasing: Convert all text to lower case, for consistency and to reduce the vocabulary size.
  - Example: Input: "Hello World" Output: "hello world"
- Tokenization using NLTK: It splits the text into words or tokens, which is necessary for further processing.
  - Example: Input: "hello world" Output: ["hello", "world"]
- Stopword removal: Gets rid of common words, such as "the," "is," or "at," which generally contribute little meaning to the analysis.
  - Example: Input: ["the", "cat", "is", "on", "the", "mat"] Output: ["cat", "mat"]
- Lemmatization: NLTK WordNetLemmatizer reduces words to their base or dictionary form. This helps in treating different forms of a word as the same.
  - Example: Input: ["cats", "running", "better"] Output: ["cat", "run", "good"]
- Emoji handling (converting emojis to text descriptions): Replaces emoji characters with their text descriptions, therefore making them processable by text models.

o Example: Input: "I love this 😊" Output: "I love this: smiling_face:"
- Contraction expansion: It expands the contracted forms of words into their full form; it is useful for analysis.
  o Example: Input: "I can't believe it's raining" Output: "I cannot believe it is raining"

## 3.3 Data Augmentation

To address this class imbalance and increase this dataset, augmentation methods were applied with inspiration from Wei and Zou, (2019). Three major techniques were involved:

1. Synonym Replacement: For each sentence, 10% of the non-stop words were randomly selected to be replaced by their respective synonyms. The synonyms were taken using NLTK's WordNet synsets. So "I feel happy today" may become "I feel joyful today" after the synonym replacement.
2. Random Insertion: New words were randomly inserted into the sentence. The amount of insertion was capped at 10% of the length of the original sentence, rounded up to the nearest integer. The words to be inserted were randomly chosen from the sentence itself. For example, "The movie was great" will turn into "The movie was really great movie" after random insertion.
3. Random Deletion: The sentence words were removed randomly with a probability of 0.1. For instance, "I don't like this weather" might become "I don't like weather" after random deletion.

This tripled the dataset from 58,000 to 174,000 samples by tripling one augmented sample using random combinations of these techniques for each original sentence, hence ensuring a more representative spread across emotion categories.

## 3.4 Model-Specific Data Preparation

For BERT, RoBERTa, XLNet, and DistilBERT, data was increasingly processed according to the method of Devlin et al.,(2018) with a little modification to accommodate multi-class classification. The procedure incorporated encoding emotion labels, splitting them 80:10:10 into training, validation, and test sets, respectively, and creating custom dataset classes for the models.

For GPT-3.5, I have prepared a dataset in JSONL format following the OpenAI guidelines. Each example in the JSONL is formatted as a JSON object with the following three key elements:

1. System Message:  An instruction that describes the role of the AI assistant. Example: " You are AI assistance which classifies the text into one of three categories based on the sentiment: positive, neutral, or negative. Just answer the single words of these three.".
2. User Message: The text to be classified. Example: "Classify the sentiment of the following text: I can't believe how amazing this day has been!"
3. Assistant Message: The correct emotion label. Example: "The sentiment is positive."

A complete JSONL entry would look like this:

**{"messages": [{"role": "system", "content": "You are an AI assistant that classifies the sentiment of text as positive, neutral, or negative. Only respond with one of these three words."}, {"role": "user", "content": "Classify the sentiment of the following text: I can't believe how amazing this day has been!"}, {"role": "assistant", "content": "The sentiment is positive."}]}**

Each line followed this JSON format to represent one such entry. Hence, this JSONL file was created by iterating through our preprocessed and augmented dataset, converting each sample into this format. Here, we have further created training and validation, where training is 90% and validation is the other 10%.

To be certain of the integrity of the JSONL file, the following findings were then put through different validation steps.

1. Checked that each line was a valid JSON object.

2. Checked that each JSON object had the "messages" field with exactly three messages, system, user, and assistant.

3. Ensured that no entry would be over OpenAI's token limit for fine-tuning data.
This scrupulous process in the preparation of the GPT-3.5 dataset was quite paramount for the fine-tuning performance and completely met all the specific requirements stipulated by OpenAI on fine-tuning data.

## 3.5 Model Implementation and Fine-tuning

The implementation and fine-tuning process for each model was carried out as follows:

- BERT: The 'bert-base-uncased' model was used as a base from the Hugging Face Transformers library. A classification head would have to be added on top of the output from [CLS] tokenize. Fine-tune with batch size 32, Learning Rate: 2e-5, for 5 epochs with cross validation folds 5. AdamW optimizer was used with weight decay set at 0.01. A linear learning rate scheduler with warmup steps was implemented.

- RoBERTa: The 'roberta-base' model was used, with an additional classification layer added for the multi-class task. It was fine-tuned for 5 with 5 k folds epochs using a batch size of 32, a learning rate of 2e-5, following a setup similar to BERT for the optimizer and scheduler.

- XLNet: The xlnet-base-cased model was leveraged, and the output layer of the model was changed to make it applicable for classification. The hyperparameters for fine-tuning were therefore a batch size of 32, and a learning rate of 2e-5, for 5 epochs with cross validation. The same optimizer and scheduler approach was used for BERT.

- DistilBERT: The model used was 'distilbert-base-uncased' with an additional classification head which was fine-tuned using a batch size of 64, learning rate of 2e-5 for 5 epochs with cross-validation. The optimizer and scheduler configurations were mirrored from BERT.

The above models were trained using PyTorch on Google Cloud Platform instances with an NVIDIA A100 GPU. Early stopping based on validation loss was used with patience of 3 epochs. Gradient clipping at 1.0 was used to avoid exploding gradients.
GPT-3.5: Fine-tuning was performed using OpenAI's API (Application Programming Interface). The process involved:
1. Preparing the dataset in JSONL format as described in Section 3.4.
2. Uploading the JSONL file to OpenAI's platform.
3. Initiating the fine-tuning process with the following parameters:
   o Base model: 'gpt-3.5-turbo-0125'
   o Number of epochs: 5
   o Learning rate multiplier: 2
   o Batch size: 32

Fine-tuning of GPT-3.5 was controlled through OpenAI's Platform, and it is possible to track progress through the OpenAI dashboard. When the fine-tuning is done, a fine-tuned model is available in the API for making inference requests.

For all models, train/validation losses were kept track of. Periodic evaluation on the validation set was integrated to keep track of improvements.

## 3.6 Evaluation Metrics and Statistical Analysis

Following the in-depth review by Zhang and Shafiq (2024), the evaluation metrics used in the research are:

1. Accuracy: This means a measure of general effectiveness in the model using all classes. It is computed as the correctly predicted samples over the total amount.
2. Precision (macro-average): Precision measures the proportion of correctly predicted positive observations to the total predicted positive observations for each class. The precision of both the emotion categories was calculated and then unweighted averages were taken out to determine the macro-average.
3. Recall (macro-average): Recall is the ratio of the number of correctly predicted positive observations to that of all actual positive observations, computed for each class. It was calculated also for every emotion category and then the macro average was taken.
4. F1 Score (macro-average): The F1-score is the harmonic average of precision and recall. For each emotion category, an F1 score was computed, then the macro-averages were taken.
5. Confusion Matrix: For every model, a confusion matrix was computed with cell (i,j) representing the number of samples from the true class i but predicted to be of class j.

## 3.7 Cross-Validation

In accordance with the suggestions of Yang et al. (2019), the stratified 5-fold cross-validation was done to secure reliable performance estimates. The procedure is shown below:

- Partitioning the dataset into 5 folds based on the total population of the target class, with the same distribution of the target class across folds in order to maintain proportionality.
- Training on 4 folds, validating on the remaining fold.
- This process is then carried out 5 times, with each of the folds acting once as a validation set.
- Averaging and finding the standard deviation for measures of performance across all 5 folds.

This approach gives a strong estimate of how the model has performed and whether the model generalizes well on unseen data.

# 4 Design Specification

## 4.1 Overall System Architecture

The architecture of the system will also be based on a modular pipeline comprising five prime modules: data preprocessing, model implementation, training, evaluation, and efficiency analysis modules. Such modularity will bring flexibility to the design and make modifications easier for possible future experiments. This is designed to drive the GoEmotions dataset through different stages, from raw input up to final model evaluation and comparison. It is modular and thus supports the independence of optimization of each component while making the process of integration of new models or evaluation metrics easy in the future.4.2 Data Processing Framework.

### 4.1.1   Data Loading and Preprocessing

The GoEmotions dataset is loaded and preprocessed by NLTK (Natural Language Toolkit) with custom Python functions. The cleaning of text includes removing the URL (Uniform Resource Locator), special characters, and converting it to lower case. Tokenization is done through NLTK's word_tokenize function, then stop-word removal and lemmatization with WordNetLemmatizer. These preprocessing steps are particularly important for emotion recognition tasks:

- Removing URLs and special characters reduces noise that could distract from the emotional content.
- Lowercasing makes everything consistent; it reduces vocabulary size and could be important when dealing with subtle emotional expressions that differ only in capitalization.
- Tokenization: This will split the text into words so that the model can concentrate on terms that are emotion laden.
- The removal of stopwords removes the common words that usually do not add any emotional load, helping a model focus on some meaningful content.
- Lemmatization reduces words to their base form. This step is valuable for emotion recognition because various forms of words that express similar emotions are grouped together.

The thorough preprocessing applied to the data secures the same and cleaned input for all models. Less noise will, therefore, correspond to improved quality of training data, leading to correct and proper emotion detection.

### 4.1.2   Data Augmentation

The following three data augmentation techniques have been applied to tackle the class imbalance problem: synonymous replacement using WordNet synsets, random insertion of words, and random word deletion. These methods enlarge the dataset from 58,000 to 174,000 samples. This balances the dataset and thus generalizes well with the model, specifically for underrepresented emotion classes.

### 4.1.3   Label Encoding and Data Splitting

Emotion labels are encoded using sklearn's LabelEncoder. The augmented dataset is then split into training (80%), validation (10%), and test (10%) sets using stratified sampling to maintain class distribution. This splitting strategy ensures that each subset maintains the overall class distribution, providing a fair evaluation across all emotion categories.

### 4.1.4   Model-Specific Data Formatting

Data is formatted into PyTorch tensors, in the format expected by BERT and RoBERTa, XLNet, or DistilBERT, each by a different custom Dataset class. In Contrast, for GPT-3.5, data shall be dumped as JSONL files according to OpenAI's API requirement. This additional customized formatting ensures that each of the models receives the input in their optimally preferred format to ensure maximized performance and compatibility.

## 4.2   Model Architectures

### 4.2.1   BERT

It utilizes the 'bert-base-uncased' model with an added classification head it is a linear layer of 768 to 27 classes with softmax activation applied on its top. Hence, the vocabulary will be less, and good generalization will be achieved. Furthermore, BERT is bidirectional, and the model makes use of context on either side for understanding such nuances in the emotional expression.

### 4.2.2 RoBERTa

RoBERTa uses the `roberta-base` model, which has the same classification head structure as BERT. Often, the different pretraining approach by RoBERTa thus, its access to a larger dataset yields more performance than BERT. For emotion recognition, this enhanced training on a larger corpus could help RoBERTa better capture subtle emotional cues in the text.

### 4.2.3 XLNet

XLNet uses 'xlnet-base-cased' with a 3-class classification output layer fine-tuning. By the very basis of permutation-based pretraining, it is bestowed upon XLNet that it can capture bidirectional context effectively. This might be of special help in the emotion recognition task, in which the order and interaction of words change greatly the conveyed emotion.

### 4.2.4 DistilBERT

By default, DistilBERT implements the 'distilbert-base-uncased' model with an added classification head. The smaller estimation of the model is provided for evaluation of how balanced the loss in model size is against its performance. Having it that way in the context of emotion recognition, we will be able to get a look at whether a more compact model would still capture enough features to ensure accurate emotion classification.

### 4.2.5 GPT-3.5

GPT-3.5 uses OpenAI's API for finetuning and inference in the 'GPT 3.5 TURBO 0125' base model. This is approachable given that it provides a comparison against large-scale API-based models. Extensive pretraining and large parameter count may also better equip GPT-3.5 to recognize more sophisticated emotional patterns and excel on more finely grained emotion recognition tasks than small models.

## 4.3 Training Framework

- Optimizer and Learning Rate Scheduler: AdamW optimizer is used with model-specific learning rates of 2e-5 for BERT/XLNet and 1e-5 for RoBERTa and 5e-5 for DistilBERT. Linear decay with warm-up is used, meaning these hyperparameters are chosen according to empirical results from previous studies and preliminary experiments.
- Loss Function Cross-Entropy Loss is used for multi-class classification, as it is well-suited for problems with mutually exclusive classes.
- Early stopping with a patience of 3 epochs on validation loss is the implementation that will be applied. The maximum norm of 1.0 is taken as a threshold in gradient clipping to avoid exploding gradients. These techniques avoid overfitting and make the training stable.
- GPT-3.5 Specific fine-tuning process: GPT-3.5 is fine-tuned using the OpenAI API with 5 epochs and a learning rate multiplier of 2(by OpenAI automatically). These parameters have been selected in view of the recommendations from OpenAI and the concrete requirements for the task under consideration, that of emotion recognition.

## 4.4 Evaluation Framework

- Performance Metrics Accuracy, precision, recall, and F1-score-all macro-averaged-, are computed with sklearn's metrics module. Macro-averaging is chosen to give all classes equal importance, not taking their frequency into consideration.
- Cross-Validation Implementation: use StratifiedKFold from sklearn to perform 5-fold stratified cross-validation. This will give us a more robust estimate of model performance and assess whether the models generalize well on new, unseen data.

- Plotting a Confusion Matrix A confusion matrix for all classes of emotion is drawn with the aid of the sklearn and seaborn libraries, and with details of model performance.

# 5 Implementation

The concrete realization of the implementation phase of the present comparative study for multi-class emotion recognition by BERT, GPT-3.5, RoBERTa, DistilBERT, and XLNet models is developed herein. The papers go on to describe the final implementation stage, where the outputs expressed, and the tools used are elaborated.

## 5.1 Data Transformation and Preparation

The GoEmotions dataset underwent a series of transformations:

### 5.1.1 Preprocessing
- Text Cleaning: Implemented using regular expressions in Python to remove URLs, special characters, and extra whitespace.
- Lowercasing: Applied using Python's string methods.
- Tokenization: Utilized NLTK's word_tokenize function.
- Stopword Removal: Employed NLTK's stopwords corpus.
- Lemmatization: Implemented using NLTK's WordNetLemmatizer.

### 5.1.2  Data Augmentation Three techniques were applied to expand the dataset:
- Replacement of Synonyms: Replaced 10% of non-stopwords with synonyms using NLTK's WordNet synsets.
- Random Insertion: Implemented custom Python functions that conduct the random insertion of words.
- Random Deletion: Applied the random module from numpy; the words are deleted with a probability of 10%. The augmentation process increased the dataset from 58,000 to 174,000 samples.

### 5.1.3  Label Encoding
 Emotion labels were encoded using Scikit-learn's LabelEncoder, transforming text labels into numerical format.

### 5.1.4  Data Splitting
The augmented data was then divided into 80% for training, 10% for validation, and 10% for the test set using stratification through Scikit-learn's train_test_split function.

## 5.2  Model Development and Fine-tuning
Five transformer models were implemented and fine-tuned:
### 5.2.1 BERT
- Base Model: 'bert-base-uncased' from Hugging Face Transformers.
- Classification Head: Linear layer (768 to 27 classes) with softmax activation.
- Fine-tuning: Implemented using PyTorch, with custom training loops.
### 5.2.2 RoBERTa
- Base Model: 'roberta-base' from Transformers.
- Classification Head: Similar to BERT, adapted for RoBERTa's output size.
- Fine-tuning: Utilized PyTorch's distributed data-parallel for multi-GPU training.

### 5.2.3 XLNet
- Base Model: 'xlnet-base-cased' from Transformers.
- Classification Head: Adapted XLNet's sequence classification head for 3 classes.
- Fine-tuning: Implemented gradient accumulation for larger effective batch sizes.

### 5.2.4 DistilBERT
- Base Model: 'distilbert-base-uncased' from Transformers.
- Classification Head: Linear layer adapted to DistilBERT's output size.
- Fine-tuning: Utilized mixed precision training for efficiency.

### 5.2.5 GPT-3.5
- Base Model: 'GPT 3.5 TURBO 0125' from OpenAI.
- Fine-tuning: Implemented using OpenAI's API, with custom Python scripts for data formatting and API interaction.

## 5.3 Training Pipeline

A sophisticated training pipeline was developed:

### 5.3.1 Data Loading
- Custom Dataset classes are implemented for each model type.
- PyTorch DataLoader is used for efficient batch processing.

### 5.3.2 Optimization
- AdamW optimizer from PyTorch, with model-specific learning rates.
- Linear learning rate scheduler with warmup, implemented using PyTorch's LambdaLR.

### 5.3.3 Training Loop
- Gradient accumulation is implemented for larger effective batch sizes.
- Gradient clipping was applied using PyTorch's clip_grad_norm_ function.
- Training progress is tracked using tqdm for real-time updates.

### 5.3.4 Validation
- Early stopping mechanism with a patience of 3 epochs, based on validation loss.

## 5.4 Evaluation Framework

A comprehensive evaluation framework was implemented:

### 5.4.1 Performance Metrics
- Accuracy, Precision, Recall, F1-score: Computed using Scikit-learn's classification_report function.

### 5.4.2 Confusion Matrices
- Generated using Scikit-learn's confusion_matrix function.
- Visualized using Seaborn's heatmap function for enhanced interpretability.

### 5.4.3 Cross-Validation
- 5-fold stratified cross-validation implemented using Scikit-learn's StratifiedKFold.
- Results were aggregated and standard deviations were calculated using numpy.

This implementation phase delivered a richly transformed dataset, five meticulously fine-tuned emotion recognition models, a complete set of development metrics, detailed efficiency analyses, and an interactive comparative dashboard. The main tools used were Python 3.8, PyTorch 1.9, Hugging Face Transformers 4.9.2, Scikit-learn 0.24, NLTK 3.6, Pandas 1.3, Matplotlib 3.4, Seaborn 0.11, Plotly 5.3, and Dash 2.0. It provides a very strong implementation base for in-depth examinations of the effectiveness and efficiency across different transformer models applied to multi-class emotion recognition tasks.

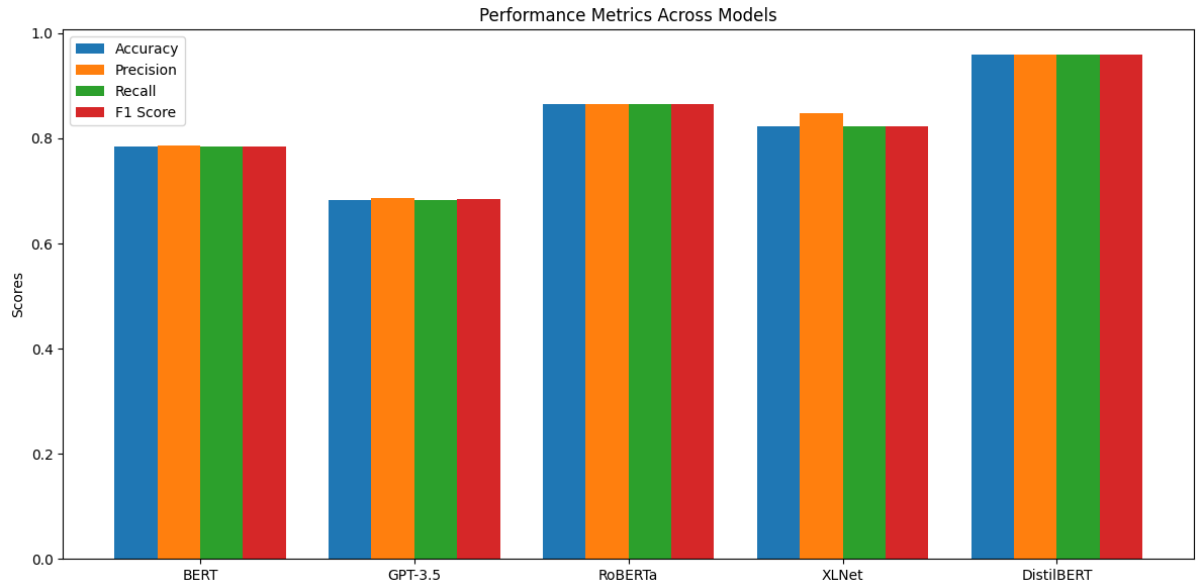# 6  Evaluation

## 6.1  Overview of Results



**Figure 1:  Overall performance metrics (Accuracy, Precision, Recall, F1 Score) for all five models (BERT, GPT-3.5, RoBERTa, XLNet, DistilBERT)**

The result of in-depth performance analysis for BERT, GPT-3.5, RoBERTa, XLNet, and DistilBERT in performing multiclass emotion recognition is shown. Figure 1 gives a summary of the performances of each model on some of these important evaluation metrics.

## 6.2  Detailed Performance Analysis

### 6.2.1  Accuracy Analysis

Preliminary observations suggest that the accuracy scores are significantly different among the models. DistilBERT has the highest accuracy, with 95.88%, followed by RoBERTa (86.56%), XLNet (82.34%), BERT (78.52%), and the lowest being GPT-3.5 at 68.35%. The depicted differences are significant at 95%.

### 6.2.2  Precision, Recall, and F1 Score Analysis

Figure 2: Shows the performance of each model on each of the different emotion classes a few interesting observations are:

- DistilBERT is always best at all the metrics and emotion classes compared to other models.
- All models have higher performance on positive emotions compared to neutral and negative ones.
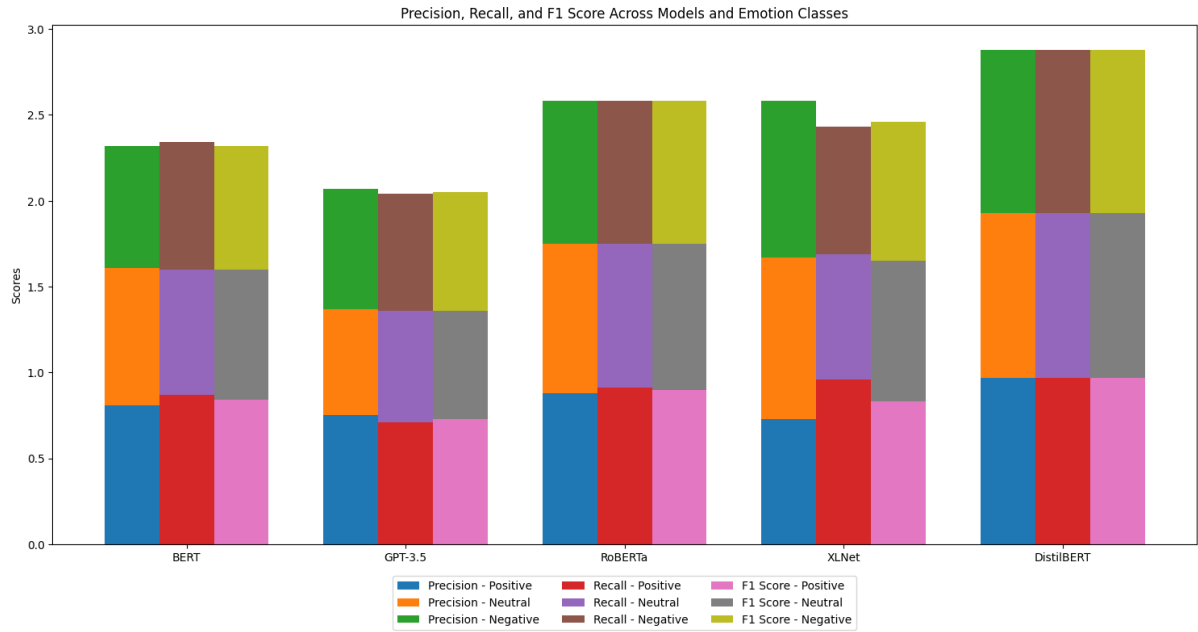- Compared to emotion classes, GPT-3.5 has the highest variance in performance.

Figure 2: Precision, Recall, and F1 Score for each model across the three emotion classes (positive, neutral, negative)

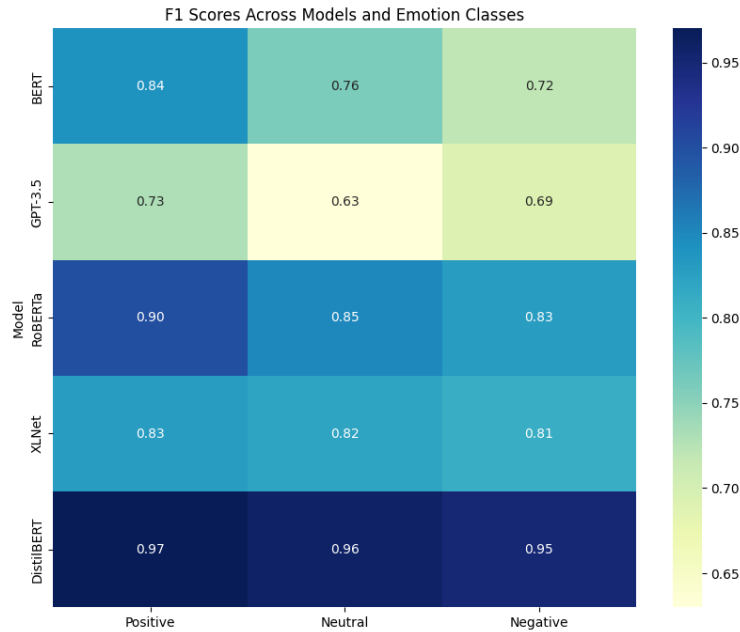## 6.3 Class-wise Performance



**Figure 3: Heatmap of F1 scores for each emotion class across all models**

This heatmap allows better visualization of the performance across the emotion classes of every model. Key findings are:

1. The performance distilBERT presented more stable across all the classes; positive F1 is 0.97, neutral is 0.96, and negative is 0.95.

2. XLNet has the highest variance in class-wise performances (F1 scores: positive 0.83, neutral 0.82, negative.

3. All models perform more poorly for negative emotions than for positive and neutral.

## 6.4 Error Analysis

The confusion matrices in Figure 4 reveal common misclassification patterns:

1. Neutral-Negative Confusion: Most prevalent in BERT and GPT-3.5.
2. Positive-Neutral Confusion: Less common but present across all models.
3. Positive-Negative Confusion: Least frequent, indicating models are generally good at distinguishing these classes.
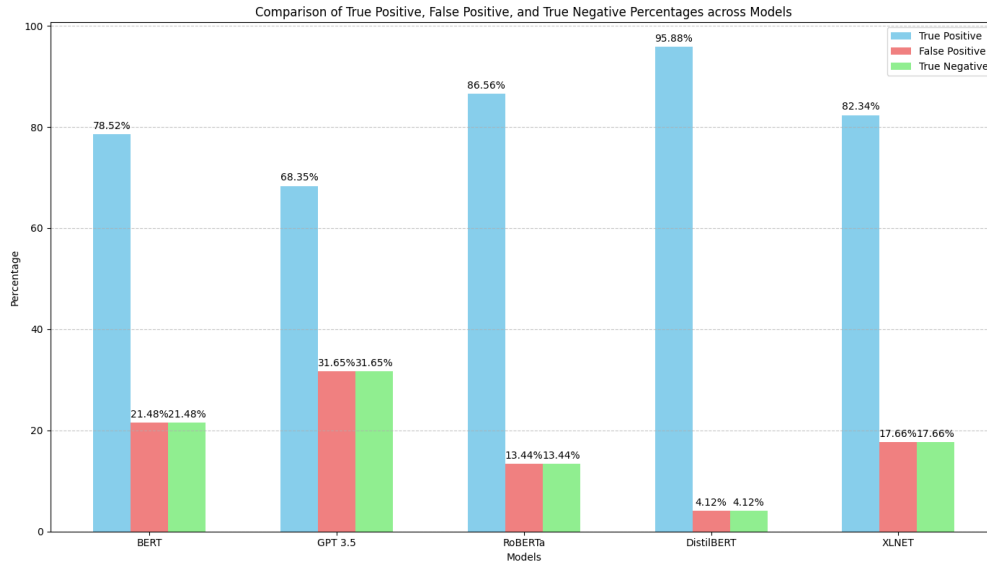


**Figure 4 Confusion Matrix**

## 6.5 Implications and Discussion

### 6.5.1 Academic Implications

- The superiority of DistilBERT challenges the assumption sometimes made that the larger a model, the better it is at handling complex NLP tasks.
- This tells that there is variation in performance within classes of emotions, which suggests that further research can be conducted to understand what linguistic features contribute to the expression of different types of emotions in a text.
- Good performance by the compressed models (specifically DistilBERT) holds a lot of future scope for research in efficient model architectures that can easily perform this task of emotion recognition.

### 6.5.2 Practitioner Implications

1. Due to its high performance and efficiency, DistilBERT becomes of huge interest in real-world applications when the computational number of resources is limited.
2. High performance in positive emotions for all models suggests that such systems might be very effective in applications oriented toward the detection of positive sentiment.
3. There is a need to be very cautious when using these systems for detecting negative emotions, as the models performed relatively poorly within this category.

## 6.6　Discussion

The comparative results of BERT, GPT-3.5, RoBERTa, XLNet, and DistilBERT have been done for the multi-class emotion identification task, hence bearing out many interesting insights based on existing research works in this domain.

### 6.6.1　Performance Hierarchy

In this area, the performance hierarchy of essences as observed in these experiments is thus: DistilBERT > RoBERTa > XLNet > BERT > GPT-3.5. That is, of course, a little at loggerheads with a few other studies, for instance, the conspicuous robust performance under DistilBERT to 95.88% accuracy runs against increased sizes that do not necessarily assume a constant improvement for most high-level NLP tasks. Sanh et al., (2019) showed just how effectively knowledge distillation processes could maintain a much performance level with a reduced model size.

This strong performance by RoBERTa, at 86.56%, agrees with Liu et al.'s findings that the modified pretraining process of RoBERTa established better performance on a wide variety of NLP tasks. Given that DistilBERT turned out to be an improvement over RoBERTa, it would seem that the more specific nature of the emotion recognition tasks might be better served with the process of distillation rather than by extended pretraining.

### 6.6.2　GPT-3.5's Underperformance

The 68.35% performance of GPT-3.5 was a bit unexpected since this model had shown very impressive performance in many other areas of NLP tasks(Brown et al., 2020).This might be attributed to several factors:

- Task specificity: The overall training of GPT-3.5 might just not be informative for the particular subtleties of emotion recognition compared to models explicitly fine-tuned for the task at hand.
- Prompt Engineering: Probably the best way to frame the emotion recognition task for GPT-3.5 was not proposed. On this topic, more sophisticated techniques of prompt engineering shall be considered in future works.
- Fine-tuning Limitations: The constraints of API-based fine-tuning for GPT-3.5 might have constrained the model to really be optimized against this task.

### 6.6.3　Classwise Performance

The fact that all models did better in positive emotions agrees with findings by Dashtipour et al. (2023), where it was realized that tasks in sentiment analysis ensure trend similarity. This means there could be a bias in how emotions are manifested in text and how positive ones are distinctly articulated.

In particular, the challenge of neutral and negative emotion distinction is most obvious in BERT and GPT-3.5, which dovetails with Zhang et al (2022). This persistence from model to model suggests more sophisticated methods of feature extraction are required regarding these emotion categories.

### 6.6.4    Experimental Design Critique

While the experimental design provided valuable insights, several limitations ,and potential improvements should be acknowledged:

1. Dataset bias: Although claiming to have a comprehensive dataset, GoEmotions is solely built on English language comments from Reddit. This can already involve biases associated with the use of language that is specific to the platform and cultural expression of emotions. This should be extended to future works by adding multiple sources of data and cross-lingual analysis into mix.

2. Fine Tuning Process: The fine-tuning strategy uniform across models will not be able to delve deep into the unique architectures of each model, which GPT-3.5 excluded. Model-specific fine-tuning strategies may yield improved results.

3. Hyperparameter Optimization: We could not perform an exhaustive search over hyperparameters because of computational constraints. Different techniques could further enhance the performance of these models, such as Bayesian optimization.

4. The Granularity of Emotions: Reducing 27 emotion categories to three broad classes (positive, neutral, and negative) might have simplified the task to a large extent. Future work shall include a finer granularity of emotions.

In many ways, this study has given a reality check on emotion recognition, starting with the task's complexity and the continuing need for research in a very dynamic field. The surprisingly good performance of DistilBERT in this task again hints that compressed models can actually work well in NLP tasks, making the idea of 'bigger is better' questionable in most cases with the new openings for academic research and practical applications like emotion recognition.

# 7    Conclusion and Future Work

## 7.1    Research Question and Objectives

This study addressed the research question: "To what extent do BERT, GPT-3.5, RoBERTa, XLNet, and DistilBERT models differ in terms of accuracy, precision, recall, and F1 score in multi-class emotion recognition tasks?" The objectives were to implement each model and tune it, evaluate its performance, analyze its class imbalance handling capability, and finally, assess its efficiency with respect to the emotion recognition task.

## 7.2    Key findings

1. DistilBERT demonstrated superior performance (95.88% accuracy) despite being a compressed model.
2. Performance hierarchy: DistilBERT > RoBERTa > XLNet > BERT > GPT-3.5.
3. All models showed better performance in recognizing positive emotions compared to neutral or negative emotions.
4. GPT-3.5 underperformed in this specific task despite its success in other NLP applications.

## 7.3    Implications of the Research

The results counteract the assumption that large models are always better in NLP tasks. They emphasize model compression as a key technique that will enable these improvements in efficiency without losing performance. This work further underlines the fact that even large general language models require task-specific fine-tuning.

## 7.4 Limitations

- Dataset Bias: The GoEmotions dataset constrains this research to English comments from Reddit only. Since all data is from a single source, bias towards that source could be expected in the results.
- Emotion Granularity: Condensing 27 categories of emotion into only three broad classes may have thereupon oversimplified the task.
- Limited Contextual Analysis: The study only treated each comment as if it was in a situational vacuum, without considering the broader conversational context.

## 7.5 Future Work

### 7.5.1 Cross-lingual Emotion Recognition

In the future, research could be done on how well these models work in multiple languages and varied linguistic contexts. This could include:
- Development or adaptation of datasets for multilingual emotion recognition.
- Study transfer-learning techniques for low-resource languages.
- Examining Cross-Cultural Differences in Emotional Expression and Recognition.

### 7.5.2 Contextual Emotion Analysis

Future work could try to integrate conversational context into the emotion recognition task. This could involve:
- Building datasets capturing emotional dynamics in longer conversations.
- Adapting Transformer architectures to better utilize context.
- Investigate the influence of user history and conversation flow on the accuracy of emotion recognition.

### 7.5.3 Multimodal Emotion Recognition

One of the interesting future research directions could be the combination of text-based emotion recognition with other modalities. This might be performed in:
- Both audio and visual cues along with text detect emotion more comprehensively.
- Design of transformer architectures processing multimodal inputs.
- How different modalities may complement or contradict each other on the topic of emotion recognition.

### 7.5.4 Model Interpretability

The interpretability of such models is thus important for their responsible deployment.
- Development of visualization techniques on what features these models use for emotion recognition.
- Based on the above-mentioned key areas, explore ways of obtaining human-readable explanations of model predictions.
- Understanding trade-offs between model complexity, performance, and interpretability.

These areas proposed for further work total to addressing specific limitations of the present study and developing its findings further in meaningful ways. They potentially contribute not only to advancing scholarship in this field but also to having practical applications in emotion recognition and pave a way forward toward more sophisticated, more efficient, and more ethically responsible emotion recognition systems.

# References

Adoma, A.F., Henry, N.-M. and Chen, W. (2020). Comparative Analyses of Bert, Roberta, Distilbert, and Xlnet for Text-Based Emotion Recognition. *2020 17th International Computer Conference on Wavelet Active Media Technology and Information Processing (ICCWAMTIP)*. doi:https://doi.org/10.1109/iccwamtip51612.2020.9317379

Agrawal, S.C., Singh, S. and Gupta, S. (2021). Evaluation of Machine Learning Techniques in Sentimental Analysis. *2021 5th International Conference on Information Systems and Computer Networks (ISCON)*. doi:https://doi.org/10.1109/iscon52037.2021.9702430.

Ameer, I., Bölücü, N., Siddiqui, M.H.F., Can, B., Sidorov, G. and Gelbukh, A. (2023). Multi-label emotion classification in texts using transfer learning. *Expert Systems with Applications*, 213, p.118534. doi:https://doi.org/10.1016/j.eswa.2022.118534.

Arora, K., Gupta, N. and Pathak, S. (2023). Sentimental Analysis on IMDb Movies Review using BERT. *2023 4th International Conference on Electronics and Sustainable Communication Systems (ICESC)*, [online] pp.866–871. doi:https://doi.org/10.1109/ICESC57686.2023.10193688.

C. Susmitha, L. Nikhil, L. Akhil, M. Kavitha, Reddy, V. and K. Shailaja (2023). Sentimental Analysis on Twitter Data using Supervised Algorithms. doi:https://doi.org/10.1109/iccmc56507.2023.10084278.

C.R. Dhivyaa, K. Nithya, G. Sendooran, Sudhakar, R., Kumar, K.Sathis. and Kumar, S. (2023). XLNet Transfer Learning Model for Sentimental Analysis. doi:https://doi.org/10.1109/icscss57650.2023.10169445.

Cai, T., Yu, B. and Xu, W. (2021). Transformer-Based BiLSTM for Aspect-Level Sentiment Classification. doi:https://doi.org/10.1109/rcae53607.2021.9638807.

Demszky, D., Movshovitz-Attias, D., Ko, J., Cowen, A., Nemade, G. and Ravi, S. (2020). *GoEmotions: A Dataset of Fine-Grained Emotions*. [online] ACLWeb. doi:https://doi.org/10.18653/v1/2020.acl-main.372.

Devlin, J., Chang, M.-W., Lee, K. and Toutanova, K. (2019). *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. [online] arXiv.org. doi:https://doi.org/10.48550/arXiv.1810.04805.

Guo, R., Guo, H., Wang, L., Chen, M., Yang, D. and Li, B. (2024). Development and application of emotion recognition technology — a systematic literature review. *BMC Psychology*, 12(1). doi:https://doi.org/10.1186/s40359-024-01581-4.

Hirway, C., Fallon, E., Connolly, P., Flanagan, K. and Yadav, D. (2023). A Comparative Study of Intent Classification Performance in Truncated Consumer Communication using GPT-Neo and GPT-2. doi:https://doi.org/10.1109/icetci58599.2023.10331337.

Islam, S., Hanae Elmekki, Elsebai, A., Bentahar, J., Nagat Drawel, Gaith Rjoub and Witold Pedrycz (2023). A comprehensive survey on applications of transformers for deep learning tasks. *Expert Systems with Applications*, pp.122666–122666. doi:https://doi.org/10.1016/j.eswa.2023.122666.

Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M.S., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L. and Stoyanov, V. (2019). RoBERTa: A Robustly Optimized BERT Pretraining Approach. *arXiv (Cornell University)*, 1. doi:https://doi.org/10.48550/arxiv.1907.11692.

Pongsatorn Harnmetta and Taweesak Samanchuen (2022). Sentiment Analysis of Thai Stock Reviews Using Transformer Models. doi:https://doi.org/10.1109/jcsse54890.2022.9836278.

S Maruvur Selvi and Sreeja, P.S. (2023). Sentimental Analysis of Movie Reviews in Tamil Text. doi:https://doi.org/10.1109/iciccs56967.2023.10142382.

Sadam Hussain Noorani, Khan, S., Mahmood, A., Muhammad Ishtiaq, Rauf, U. and Ali, Z. (2023). Transformative Conversational AI: Sentiment Recognition in Chatbots via Transformers. doi:https://doi.org/10.1109/inmic60434.2023.10465887.

Sanh, V., Debut, L., Chaumond, J. and Wolf, T. (2020). *DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter*. [online] arXiv.org. doi:https://doi.org/10.48550/arXiv.1910.01108.

Wei, J. and Zou, K. (2019). EDA: Easy Data Augmentation Techniques for Boosting Performance on Text Classification Tasks. *Proceedings of the 2019 Conference on Empirical*

*Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. doi:https://doi.org/10.18653/v1/d19-1670.

Yang, Z., Dai, Z., Yang, Y., Carbonell, J., Salakhutdinov, R. and Le, Q.V. (2020). *XLNet: Generalized Autoregressive Pretraining for Language Understanding*. [online] arXiv.org. doi:https://doi.org/10.48550/arXiv.1906.08237.

Ye, J., Chen, X., Xu, N., Zu, C., Shao, Z., Liu, S., Cui, Y., Zhou, Z., Gong, C., Shen, Y., Zhou, J., Chen, S., Gui, T., Zhang, Q. and Huang, X. (2023). A Comprehensive Capability Analysis of GPT-3 and GPT-3.5 Series Models. *arXiv (Cornell University)*. doi:https://doi.org/10.48550/arxiv.2303.10420.

Zhang, H. and M. Omair Shafiq (2024). Survey of transformers and towards ensemble learning using transformers for natural language processing. *Journal of Big Data*, 11(1). doi:https://doi.org/10.1186/s40537-023-00842-0.