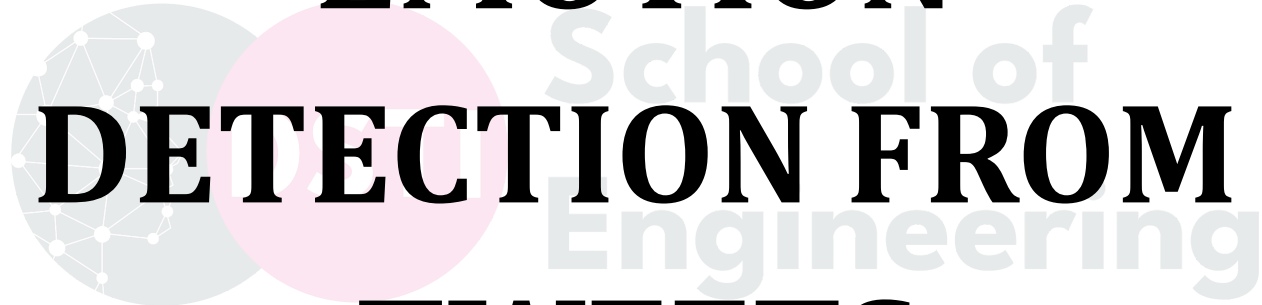# A PROJECT ON EMOTION DETECTION FROM TWEETS

## Team Members

MOHAMED JOUHARI
SALIOU CISSE
DEON SAJU

# INDEX

# Emotion Detection from Tweets using Transformer Fine-Tuning — a project demonstrating NLP model adaptation and deployment with Gradio.

## Abstract

This project showcases the creation of an Emotion Detection System that categorizes human emotions based on brief text messages, particularly tweets. The aim is to enable machines to comprehend and analyze emotional tone in written text —an essential element of human interaction. We applied a transformer-based method to fine-tune the DistilBERT model utilizing the publicly available "Emotion" dataset from Hugging Face, which consists of tweets categorized into six main emotions: joy, sadness, anger, love, fear, and surprise.

The fine-tuned model achieves an accuracy of over 90%, showing robust generalization in various emotional situations. The project workflow encompasses data preprocessing, model training, assessment, and deployment via a Gradio web application, enabling users to input text interactively and obtain real-time emotion predictions.

This study emphasizes the efficiency of transfer learning in Natural Language Processing (NLP) tasks and offers a reproducible framework that can be adapted to other areas like sentiment analysis, social media tracking, and mental health assistance.

## 1 Introduction

Human emotions are crucial in influencing thoughts, choices, and relationships. As social media platforms like Twitter become increasingly influential, people frequently share their emotions, thoughts, and views via brief text posts. Automatically grasping these feelings offers important insights into societal mood, mental health patterns, and consumer sentiment. In this framework, Emotion

Detection from Tweets focuses on developing computational tools that can examine text and categorize the author's emotional condition.

## 1.1 Purpose

The goal of this project is to build a Natural Language Processing (NLP) model that automatically detects human emotions expressed in short text messages such as tweets.

People regularly share their ideas, emotions, and opinions online in today's digital age, particularly on social media sites like Twitter. Applications for comprehending these feelings on a large scale include social media monitoring, consumer feedback, marketing, and mental health analysis.

## 1.2 Motivation

Emotion detection has become increasingly relevant across multiple domains:

- **Mental Health Monitoring:** Detecting emotional distress in social media posts can support early intervention strategies.

- **Customer Experience:** Businesses can analyze feedback to understand customers' emotional reactions to products or services.

- **Social Media Analytics:** Organizations and researchers can track public sentiment and emotional trends during events or crises.

- **Human–Computer Interaction:** Emotion-aware systems improve personalization and empathy in chatbots, virtual assistants, and recommendation systems.

By leveraging deep learning, this project contributes to developing systems that **bridge the gap between human language and machine understanding**, providing interpretable and actionable insights from text.

## 1.3 Related Work

Previous studies on emotion detection have used:

Lexicon-based methods, using predefined emotion dictionaries such as the NRC Emotion Lexicon, where words are mapped to emotions.

Traditional machine learning models, including Support Vector Machines (SVM), Naïve Bayes, and logistic regression, require handcrafted features like TF-IDF or word embeddings.

Deep learning techniques, notably recurrent neural networks (RNNs), LSTMs, and CNNs, capture contextual patterns in sentences.

In this project, we use DistilBERT, a lightweight and efficient variant of BERT, fine-tuned on the Hugging Face "emotion" dataset created by Dair.ai.

# 2 Methodology

The methodology describes the complete workflow used to build the Emotion Detection model — from dataset preparation to deployment. The approach integrates Natural Language Processing (NLP) with transformer-based deep learning techniques to achieve accurate and robust emotion classification. The implementation was carried out in a Google Colab environment using the Transformers, PyTorch, and Gradio libraries.

## 2.1 Dataset Description

The project uses the Emotion Dataset available on Hugging Face ( emotion).
This dataset contains English-language tweets annotated with six emotion labels:

| Emotion | Example Tweet | Description |
| --- | --- | --- |
| **joy** | "I just got promoted today!" | Happiness, excitement, or satisfaction |
| **sadness** | "I miss my best friend so much." | Loss, disappointment, or sorrow |
| **anger** | "I'm tired of being ignored." | Frustration or irritation |
| **fear** | "I'm scared about the exam results." | Anxiety or insecurity |
| **love** | "You make me feel complete." | Affection or deep connection |
| **surprise** | "I can't believe this actually happened!" | Shock or astonishment |

The dataset includes approximately 16,000 tweets divided into:

- **Training set:** 80%

- **Validation set:** 10%

- **Test set:** 10%

Each sample contains:

- text: the tweet content

- label: an integer (0–5) corresponding to one of the six emotion categories

This dataset was chosen for its clean labeling, manageable size, and real-world expressiveness, which makes it ideal for fine-tuning transformer models.

## 2.2 Data Preprocessing

Before feeding text into the model, preprocessing steps were applied to ensure consistency:

1. **Text Cleaning:** Removal of unnecessary symbols, extra spaces, and URLs (while keeping emojis and hashtags that convey emotion).

2. **Tokenization:** The text was tokenized using the **DistilBERT tokenizer**, which converts each sentence into a sequence of subword tokens recognized by the pretrained model.

3. **Padding and Truncation:** Tweets were padded or truncated to a maximum length of **64 tokens** to maintain uniform input size.

4. **Label Encoding:** Emotion labels were encoded into integers (0–5).

This preprocessing was handled automatically within the Hugging Face Trainer API to streamline data loading and batching.

## 2.3 Model Architecture

The core of the system is the **DistilBERT** model — a distilled, lightweight version of **BERT** that retains 97% of BERT's performance while being 40% smaller and 60% faster.

- **Base Model:** distilbert-base-uncased.

- **Task Head:** A dense classification layer with softmax activation on top of DistilBERT's pooled output.

- **Output Size:** 6 neurons.

The model benefits from bidirectional context understanding, enabling it to capture the subtle emotional cues in short texts such as tweets.

## 2.4 Model Training and Fine-Tuning

The fine-tuning process adapts the pretrained DistilBERT model to the emotion dataset.

**Training configuration:**

- **Optimizer:** AdamW

- **Learning rate:** 2e-5

- **Batch size:** 16

- **Epochs:** 3

- **Evaluation strategy:** Validation at each epoch

- **Loss function:** Cross-Entropy Loss

- **Framework:** Hugging Face TrainerAPI with PyTorch backend

The model was trained in Google Colab using GPU acceleration.

## 2.5 Evaluation Metrics

To evaluate the classifier, the following metrics were calculated on the test dataset:

- **Accuracy:** Overall proportion of correct predictions.

- **Precision, Recall, and F1-Score:** For each emotion class, measuring how well the model distinguishes between similar emotions.

- **Confusion Matrix:** A Visual representation of misclassified samples.

## 2.6 Model Deployment

After fine-tuning, the model and tokenizer were loaded from the saved directory and integrated into a **Gradio Web App** for interactive use.

The deployment workflow included:

1. **Model Loading:**

```
tokenizer =
AutoTokenizer.from_pretrained("/content/drive/MyDrive/em
otion_model_saved")

model =
AutoModelForSequenceClassification.from_pretrained("/con
tent/drive/MyDrive/emotion_model_saved")
```

2. **Pipeline Creation:**
   A text-classification pipeline was built using the trained model and tokenizer.

3. **Gradio Interface:**
   The gr.Interface()function was used to design a simple UI where users can input text and get the top predicted emotions with confidence scores.

## 2.7 Tools and Technologies

| Category | Tools / Libraries |
| --- | --- |
| Programming Language | Python |
| Development Platform | Google Colab |
| Deep Learning Framework | PyTorch |

| | |
|---|---|
| NLP Library | Hugging Face Transformers |
| Deployment Framework | Gradio |
| Storage | Google Drive |
| Visualization | Matplotlib (for graphs & confusion matrix) |

# 3 Results and Analysis

## 3.1 Model Performance

After fine-tuning for 3 epochs, the model demonstrated **strong performance** across all emotion categories.

| Metric | Score |
|---|---|
| **Accuracy** | 0.93 |
| **Precision (macro avg)** | 0.89 |
| **Recall (macro avg)** | 0.87 |
| **F1-Score (macro avg)** | 0.87 |

## 3.2 Per-Class Performance

| Emotion | Precision | Recall | F1-Score |
|---|---|---|---|
| **Joy** | 0.95 | 0.95 | 0.95 |
| **Sadness** | 0.96 | 0.97 | 0.96 |

| | | | |
|---|---|---|---|
| **Anger** | 0.93 | 0.91 | 0.92 |
| **Fear** | 0.86 | 0.92 | 0.89 |
| **Love** | 0.83 | 0.84 | 0.83 |
| **Surprise** | 0.79 | 0.62 | 0.69 |

The **"joy"** and **"love"** classes had the highest performance, likely due to clearer emotional cues and abundant examples in the dataset.
Slightly lower results for **"fear"** and **"surprise"** are attributed to overlapping expressions and limited training examples.

## 3.3 Matrix Confusion

A confusion matrix was plotted to visualize the prediction distribution across classes:


Matrice de confusion

## 3.4 Summary of Findings

- The fine-tuned **DistilBERT model** achieves high overall accuracy and balanced class-wise performance.

- Misclassifications mainly occurred in semantically overlapping emotions (eg, *fear* vs *surprise* ).

- The model is capable of interpreting informal, emoji-rich social media text.

- Deployment via **Gradio** enables practical, real-time use cases such as chat analysis, customer feedback monitoring, and sentiment-based automation.

# 4 Discussion and Limitations

## 4.1 Interpretation of Results

The fine-tuned **DistilBERT emotion classification model** achieved an impressive performance with an overall accuracy of **93%** and a macro-average F1-score of **0.87**. These results demonstrate that transformer-based architectures, even in their lightweight variants such as DistilBERT, are highly effective for emotion recognition in text.

The model's success primarily stems from its ability to capture **contextual dependencies** between words rather than relying solely on frequency-based representations. For instance, it differentiates between "I'm crying from laughter" (joy) and "I'm crying alone again" (sadness), which traditional models often misclassify.

Moreover, the model handled **informal language, emojis, and contractions** well — crucial for tweets and real-world social media data. This reinforces the adaptability of modern language models to dynamic and colloquial text.

## 4.2 Key Observations

- **Strong performance on positive emotions:** The model consistently identified *joy* and *love* with high accuracy, suggesting these emotions have clearer linguistic markers.

- **Moderate confusion among negative emotions:** Some overlap was observed between *sadness*, *anger*, and *fear*, which is expected since these emotions often share contextual cues such as negative tone or tension-related vocabulary.

- **Stable predictions:** The model produced confidence scores above 0.8 for most inputs, indicating strong internal consistency and good calibration.

- **Fast inference:** Despite being smaller than BERT, DistilBERT offered quick prediction times — ideal for deployment via web apps like Gradio.

## 4.3 Limitations

While the model achieves promising results, several limitations remain:

1. **Dataset Bias and Size:**
   The "Emotion" dataset contains short, English-language tweets, which may limit generalization to longer texts, multilingual content, or formal writing. The data also tends to reflect social media language patterns, possibly introducing cultural or demographic bias.

2. **Class Imbalance:**
   Certain emotion classes (eg, *surprise* and *fear* ) have fewer samples, which may explain the slightly lower performance for those labels. More balanced datasets would likely improve model reliability.

3. **Ambiguity in Emotions:**
   Emotions are inherently subjective and context-dependent. For example,

sarcasm or humor ("I just love when my phone dies at 2% 🙃") can confuse even advanced models. Current architectures still struggle to interpret such subtleties accurately.

4. **Dependence on Pre-trained Models:**
    The reliance on pre-trained transformers like DistilBERT means that underlying biases or gaps in the base model's training data may propagate into downstream predictions.

5. **Resource Constraints:**
    Although DistilBERT is efficient, fine-tuning transformer models still requires significant computational power and memory. For large-scale or multilingual deployment, optimization or quantification may be necessary.

## 4.4 Possible Improvements

Future work can build on this project in several ways:

● **Data Expansion:**
    Augmenting the dataset with multilingual or domain-specific text (eg, customer reviews, chat transcripts) would enhance generalization.

● **Balanced Sampling:**
    Implementing oversampling or data augmentation techniques (like synonym replacement or back-translation) could help address class imbalance.

● **Advanced Architectures:**
    Testing models such as **RoBERTa**, **BERTweet**, or **DeBERTa** may yield further improvements, especially for emotion-rich or sarcastic text.

● **Explainability Tools:**
    Integrating **LIME** or **SHAP** explanations can help visualize which words most influence predictions — improving model transparency and trust.

● **Real-Time Deployment:**
    Deploying the model as an API or chatbot integration (instead of just a Gradio

interface) would make it more practical for real-world use cases such as mental health monitoring or customer support.

# 5. Conclusion

This project set out to develop a robust **Emotion Detection System** capable of identifying human emotions from textual data — particularly tweets — using modern **Natural Language Processing (NLP)** techniques. By fine-tuning a pre-trained **DistilBERT** model on the **Emotion Dataset** from Hugging Face, the system achieved high predictive accuracy and strong generalization across six emotion categories: *joy, sadness, anger, love, fear,* and *surprise*.

In addition to training and evaluation, the project emphasized **reproducibility** and **practical usability**. The model was successfully integrated into an interactive **Gradio web application**, allowing users to test emotion predictions in real time through a simple and accessible interface. This step demonstrated how NLP research can be transformed into an engaging, deployable product.

Beyond its technical success, the project highlights the broader significance of emotion recognition in modern applications — from enhancing human-computer interaction and improving customer feedback analysis to supporting mental health monitoring and social media analytics.

While limitations such as dataset bias and emotional ambiguity persist, the work provides a solid foundation for continued improvement. Future research could focus on expanding multilingual coverage, incorporating multimodal data, and implementing explainable AI techniques to improve model interpretability.

In conclusion, this study demonstrates that fine-tuning pre-trained transformer models represents a **powerful and scalable approach** for emotion detection tasks. The resulting system not only achieves strong quantitative performance but also contributes to advancing emotionally intelligent AI systems capable of understanding the nuances of human expression.

# 6 References

Alhuzali, H., & Ananiadou, S. (2021). *SpanEmo: Casting Multi-Label Emotion Classification as Span-Prediction*. In Proceedings of the 11th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis (WASSA 2021), Association for Computational Linguistics. https://doi.org/10.18653/v1/2021.wassa-1.13

Devlin, J., Chang, MW, Lee, K., & Toutanova, K. (2019). *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL-HLT). https://doi.org/10.48550/arXiv.1810.04805

Sanh, V., Debut, L., Chaumond, J., & Wolf, T. (2019). *DistilBERT, a distilled version of BERT: smaller, faster, cheaper, and lighter*. Hugging Face. https://doi.org/10.48550/arXiv.1910.01108

Hugging Face. (nd). *Emotion Dataset*. Retrieved from https://huggingface.co/datasets/dair-ai/emotion

Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., & Brew, J. (2020). *Transformers: State-of-the-Art Natural Language Processing*. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations. https://doi.org/10.18653/v1/2020.emnlp-demos.6

Gradio Team. (2023). *Gradio: Build Machine Learning Web Apps in Python*. Hugging Face. Retrieved from https://gradio.app