

# Configuration Manual

MSc Research Project  
Master of Science in Artificial Intelligence

**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.....

**Lecturer:** .....Rejwanul Haque.....

**Submission**

**Due Date:** ..... 16 September 2024.....

**Project Title:** ...Comparative Analysis of Transformer Models for Multi-Class Text Classification.....

**Word Count:** .....996..... **Page Count:** .....8.....

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)	<input type="checkbox"/>
<b>Attach a Moodle submission receipt of the online project submission,</b> to each project (including multiple copies).	<input type="checkbox"/>
<b>You must ensure that you retain a HARD COPY of the project,</b> both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer.	<input type="checkbox"/>

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

<b>Office Use Only</b>	
Signature:	
Date:	
Penalty Applied (if applicable):	

# Configuration Manual

Sharanya Neelakanti  
Student ID: X23138220

## 1 INTRODUCTION

This is an inference of sentiment analysis implemented using a battery of fine-tuned pre-trained language models: BERT, RoBERTa, XLNet, DistilBERT, and GPT-3.5 on the GoEmotions dataset. The system shall be made up of the following three main text classifications: positive, negative, and neutral sentiments. Key features of the project include:

- Loading pre-trained models for inference
- Evaluating model performance on a test dataset
- Predicting sentiment for single text inputs
- Handling both traditional transformer models and the OpenAI GPT-3.5 API
- Configs of Models and All

### 1.1 Test Dataset:

- Source: GoEmotions dataset (preprocessed)
- Classes: Mapped to 3 main classes (positive, neutral, negative)

## 2 Connecting Google colab

### 2.1 Open Google Colab

- Go to <https://colab.research.google.com/>
- Create a new notebook or open an existing one

### 2.2 Connect to GPU Session

- Click on "Runtime" in the top menu
- Select "Change runtime type"
- In the pop-up window, set "Hardware accelerator" to "T4 GPU"
- Click "Save"

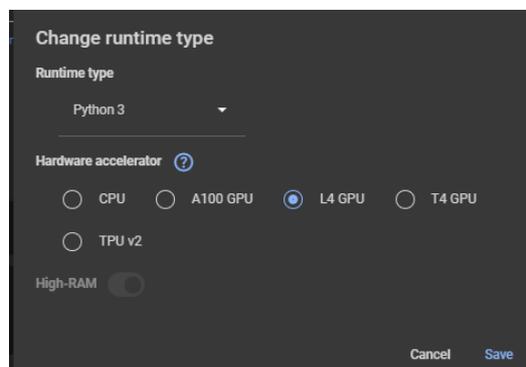
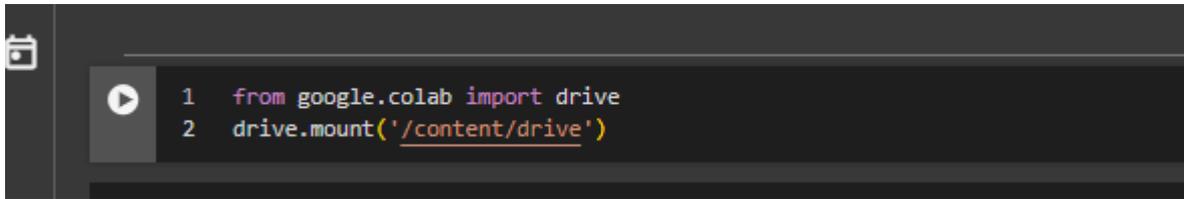


Figure 1 GPU Selection

## 2.3 Mount Google Drive

- Run the following code to mount your Google Drive

A screenshot of a code editor with a dark background. On the left, there is a play button icon. The code is as follows:

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

Figure 2 Google drive Mount

## 3 Setup and Configuration

### 3.1 Directory Structure

Create the following directory structure, replacing "PROJECT\_NAME" with a name of your choice

```
/content/drive/MyDrive/PROJECT_NAME/
├── Dataset/
├── BERT_MODEL/
├── ROBERTA_MODEL/
├── XLNET_MODEL/
└── DISTILBERT_MODEL/
```

Figure 3 Directory Structure

### 3.2 Dataset Upload

Upload the preprocessed GoEmotions dataset:

- File: 'go\_emotions\_merged\_augmented.csv' (or your specific filename)
- Upload to: '/content/drive/MyDrive/PROJECT\_NAME/Dataset/'

### 3.3 Model Upload

Upload the pre-trained models to their respective directories:

- BERT: Upload to '/content/drive/MyDrive/PROJECT\_NAME/BERT\_MODEL/'
- RoBERTa: Upload to  
'/content/drive/MyDrive/PROJECT\_NAME/ROBERTA\_MODEL/'
- XLNet: Upload to '/content/drive/MyDrive/PROJECT\_NAME/XLNET\_MODEL/'
- DistilBERT: Upload to  
'/content/drive/MyDrive/PROJECT\_NAME/DISTILBERT\_MODEL/'

## 3.4 Model Paths

3.4.1 Update the MODEL\_PATHS dictionary in your code to reflect your chosen directory structure:

```
34
35
36 MODEL_PATHS = {
37     'BERT': '/content/drive/MyDrive/PROJECT_NAME/BERT_MODEL/best_model_BERT',
38     'RoBERTa': '/content/drive/MyDrive/PROJECT_NAME/ROBERTA_MODEL/best_model_RoBERTa',
39     'XLNet': '/content/drive/MyDrive/PROJECT_NAME/XLNET_MODEL/best_model_XLNet',
40     'DistilBERT': '/content/drive/MyDrive/PROJECT_NAME/DISTILBERT_MODEL/best_model_DistilBERT',
41     'GPT-3.5': "ft:gpt-3.5-turbo-0125:personal:nci-thesis:9tKr25Ma"
42 }
43
```

Figure 4 Model Paths

3.4.2 Update the dataset path as well like below

```
print("\n")

def main():
    # Load your dataset
    data = pd.read_csv('/content/drive/MyDrive/PROJECT_NAME/Dataset/go_emotions_merged_augmented.csv')
    # Print the first few rows and column names for debugging
    print("Dataset columns:", data.columns)
    print("\nFirst few rows of the dataset:")
    print(data.head())

    while True:
        print("\nChoose an option:")
        print("1. Evaluate all models on the dataset")
        print("2. Evaluate all models on a limited number of rows")
        print("3. Get predictions for a single input")
        print("4. Exit")

        choice = input("Enter your choice (1-4): ")

        if choice == '1':
            evaluate_models(data)
```

Figure 5 Dataset Path change

## 3.4.3 OpenAI API Configuration

```
20
21 # Securely set your OpenAI API key
22 os.environ['OPENAI_API_KEY'] = 'sk-None-QfeHnc5xeLH1HfLT2pC1T3B1bkFJQBZJwd9uH7h5TEOD19ji'
23 client = openai.OpenAI()
24
```

Figure 6 OpenAI API Key

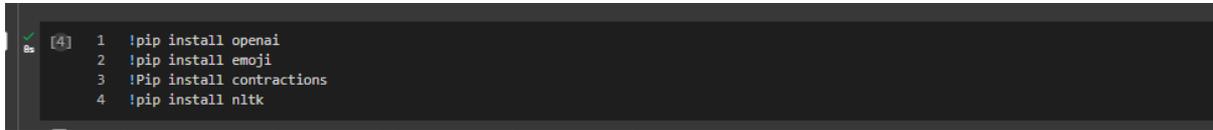
Note: The OpenAI API key provided in the code is specific to the fine-tuned model for this task. Do not change this key.

## 3.4.4 Model Loading

- BERT, RoBERTa, XLNet, and DistilBERT: Loaded using their respective tokenizers and model classes from the transformer's library.
- GPT-3.5: Accessed via OpenAI API using the provided fine-tuned model ID.

### 3.5 Library Installation

- Before running the inference, ensure all required libraries are installed. You can install them using pip

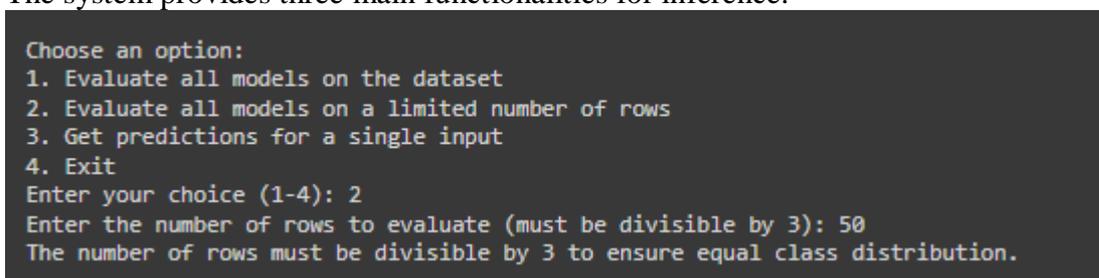


```
[4] 1 !pip install openai
    2 !pip install emoji
    3 !Pip install contractions
    4 !pip install nltk
```

Figure 7 Library installation

## 4 Inference

The system provides three main functionalities for inference:



```
Choose an option:
1. Evaluate all models on the dataset
2. Evaluate all models on a limited number of rows
3. Get predictions for a single input
4. Exit
Enter your choice (1-4): 2
Enter the number of rows to evaluate (must be divisible by 3): 50
The number of rows must be divisible by 3 to ensure equal class distribution.
```

Figure 8 Options

#### 4.1 Evaluate all models on the entire test dataset:

- Use option 1 in the main menu
- Processes the entire dataset and provides classification reports and confusion matrices for each model

#### 4.2 Evaluate all models on a limited number of samples:

- Use option 2 in the main menu
- Enter the number of samples to evaluate (must be divisible by 3 for balanced class distribution)
- Processes a subset of the data and provides classification reports and confusion matrices for each model

#### 4.3 Get predictions for a single input:

- Use option 3 in the main menu
- Enter the text to classify
- Preprocesses the text and provides predictions from all models

```
Choose an option:
1. Evaluate all models on the dataset
2. Evaluate all models on a limited number of rows
3. Get predictions for a single input
4. Exit
Enter your choice (1-4): 3
Enter the text to classify: NCI is a Good Girl;
Input text: NCI is a Good Girl;
Preprocessed text: nci is a good girl
BERT prediction: positive
RoBERTa prediction: positive
XLNet prediction: positive
DistilBERT prediction: positive
GPT-3.5 prediction: positive
```

#### 4.4 To run the inference:

1. Ensure all required libraries are installed
2. Create the directory structure as described in section 3.1
3. Upload the dataset and pre-trained models as described in sections 3.2 and 3.3
4. Set up the OpenAI API key as described in section 3.5
5. Run the main() function
6. Choose the desired option from the menu and follow the prompts

Note: This code is designed for inference only. It assumes that the models have been previously fine-tuned on the GoEmotions dataset and saved. Ensure that the model paths are correct, and the pre-trained models are available before running the inference.

## References

Sanh, V., Debut, L., Chaumond, J., & Wolf, T. (2019). DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. arXiv (Cornell University).

<https://doi.org/10.48550/arxiv.1910.01108>

Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. arXiv (Cornell University).

<https://doi.org/10.48550/arxiv.1810.04805>

Yang, Z., Dai, Z., Yang, Y., Carbonell, J. G., Salakhutdinov, R., & Le, Q. V. (2019). XLNet: Generalized Autoregressive Pretraining for Language Understanding. arXiv (Cornell University).

<https://doi.org/10.48550/arxiv.1906.08237>

Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., & Stoyanov, V. (2019). RoBERTa: A Robustly Optimized BERT Pretraining Approach. arXiv (Cornell University).

<https://doi.org/10.48550/arxiv.1907.11692>

Latif, E., & Zhai, X. (2023). Fine-tuning ChatGPT for Automatic Scoring. arXiv (Cornell University).

<https://doi.org/10.48550/arxiv.2310.10072>

<https://huggingface.co/datasets/mrm8488/goemotions>