# Cyberbullying Detection in Social Media Comments

Module: Computer Science Project (CSEMCSPCSP01)

Name: Panchami Bandihalli Varadaraju

Student ID: 4252553

Tutor: Mugdha Kashyap

IU International University of Applied Sciences

Date: 26 November 2025

## Abstract

Cyberbullying has become one of the most widespread and harmful phenomena in modern digital communication platforms. The increasing frequency and severity of abusive behavior present significant challenges for online safety, user wellbeing, and digital governance. Traditional rule-based approaches and lexicon-based methods are insufficient to detect nuanced, context-dependent, and evolving abusive language. Therefore, there is a strong need for automated, scalable, and intelligent systems capable of categorizing online harassment with high accuracy. This project investigates the application of Natural Language Processing (NLP) and transformer-based deep learning models to classify cyberbullying into six categories: age-based harassment, ethnicity-based harassment, gender-based harassment, religion-based harassment, other types of cyberbullying, and non-cyberbullying content. The central objective of this project is to design, implement, and evaluate a BERT-based text classification system to automate this categorization.

The project begins by examining existing research on cyberbullying detection using machine learning and deep learning methods. Building on these foundations, a technical background is provided covering neural networks, tokenisation, attention mechanisms, and transformer architectures. The methodological framework outlines the data pre-processing pipeline, model selection rationale, training strategy, and testing approach. The implementation combines Python, PyTorch, Hugging Face Transformers, and scikit-learn to construct a complete machine learning workflow.

Extensive testing, involving accuracy, precision, recall, and F1-scores, demonstrates that BERT significantly outperforms classical models in handling multi-class cyberbullying classification. The findings confirm that transformer-based models are highly capable of understanding contextual, sarcastic, implicit, and coded abusive expressions—cases where traditional methods typically fail.

The report concludes with recommendations for future enhancements, such as fine-tuning domain-specific language models, introducing explain ability methods, and deploying the model as a web service. Ethical considerations and data sensitivity are also discussed, acknowledging the importance of responsible AI when dealing with harmful, sensitive, or personal user-generated content. Overall, this project provides a complete end-to-end solution that demonstrates the application of advanced computer science techniques to a real-world societal problem.

**Table of Contents**

III

# 1 Introduction

Cyberbullying has become one of the most pervasive forms of online violence in the digital world. With the widespread use of social media, messaging platforms, gaming communities, and discussion forums, harmful messages can rapidly circulate, reach large audiences, and leave long-lasting impacts. Unlike traditional bullying, cyberbullying occurs continuously, anonymously, and without physical boundaries. This makes it extremely difficult for victims to escape the abuse.

Automated cyberbullying detection systems play a crucial role in creating safer online environments. They help platforms identify harmful content in real time, prevent the escalation of online harassment, and support moderation teams that would otherwise be overwhelmed by the volume of user-generated content. Early detection is also vital for educational institutions, mental-health professionals, and parents, as it enables timely intervention. As technology continues to evolve, detecting online abuse becomes a fundamental requirement for promoting digital well-being, ensuring user safety, and building a healthier online ecosystem.

The effects of cyberbullying extend far beyond temporary emotional discomfort. Numerous studies confirm that victims often suffer from severe psychological stress, anxiety, depression, low self-esteem, and social withdrawal. Because cyberbullying happens in public spaces, victims may feel humiliated in front of peers, strangers, and a wider online audience. This leads to emotional exhaustion, reduced academic or work performance, and in extreme cases, self-harm or suicidal ideation.

Cyberbullying not only affects individuals but also disrupts communities. Toxic online environments discourage healthy participation, silence marginalized voices, and reduce trust in digital platforms. The psychological burden on victims' families, educators, and mental-health professionals increases significantly when cyberbullying goes undetected for long periods. Hence, implementing robust cyberbullying detection systems has broad social value, helping to mitigate long-term consequences and foster a culture of respect and responsible digital behaviour.

Detecting abusive or harmful text is a highly complex task. Human language is deeply nuanced, and cyberbullying often appears in subtle, coded, or context-dependent forms. Many challenges include:

- **Indirect or implicit language**: Abusive intent may be hidden in sarcasm, jokes, metaphors, or cultural expressions.
- **Evolving vocabulary**: Slang, emojis, abbreviations, and online trends shift rapidly, making fixed rule-based systems ineffective.

- **Multilingual and multi-dialect usage**: Users mix languages, spell words creatively, or use regional expressions unfamiliar to traditional models.
- **Context dependency**: A single phrase can be harmless or harmful depending on previous messages, tone, sender–receiver relationship, or cultural background.
- **Imbalanced datasets**: Many cyberbullying datasets are small or skewed, making it difficult for machine learning models to generalise.
- **Ambiguity**: Some aggressive messages may appear harmless without knowing the age, gender, or identity of speakers.

Deep learning—especially transformer models like BERT—has revolutionised natural language processing. Unlike older methods, deep learning can automatically learn patterns, semantics, and contextual relationships in text without relying on manually defined rules.

The main motivations include:

- **Contextual understanding**: Models like BERT use attention mechanisms to interpret meaning based on sentence context, not isolated keywords.
- **Handling complex language**: Deep learning models can understand slang, misspellings, emoji's, and informal text more effectively.
- **Multi-class and multi-label capability**: They can classify multiple types of cyberbullying simultaneously.
- **Scalability**: Deep learning models can process millions of messages in real time.
- **Generalisation**: When trained on diverse data, they perform well across different platforms and user groups.

Deep learning brings the accuracy, flexibility, and reliability needed to detect abusive behaviour in dynamic online environments, making it the preferred method for modern cyberbullying detection systems.

This project focuses on building and evaluating a BERT-based deep learning model capable of detecting cyberbullying across six distinct categories: age, ethnicity, gender, religion, other cyberbullying, and not cyberbullying. The scope includes dataset pre-processing, model fine-tuning, inference generation, and designing a user-friendly prediction interface using Gradio.

## 2 Technical Background
### 2.1    NLP Fundamentals

**Tokenisation**

Tokenisation is the process of converting raw text into smaller units called *tokens*. These units can be words, sub words, or characters depending on the tokenizer. Traditional NLP used whitespace tokenisation, but modern architectures like BERT use **Word Piece tokenisation**, which breaks words into subword units.

Example:
*"Cyberbullying detection"* → "cyber", "##bullying", "detection"

This approach helps handle rare words, slang, emojis, and misspellings — all common in abusive or toxic language.

**Lemmatization**

Lemmatization reduces words to their base form called a *lemma*.
Examples:
*"running"* → *"run"*, *"better"* → *"good"*.
Although modern transformer models do not rely heavily on lemmatization, understanding it is important because classical NLP pipelines used it extensively for vocabulary reduction.

**Stopwords**
Stopwords are high-frequency words like *"the", "is", "at", "of"*. In traditional NLP pipelines, these are removed to reduce noise. However, models like BERT do **not** remove stopwords because attention mechanisms rely on complete sentence structures for contextual understanding.

**Bag-of-Words (BoW) vs Word Embeddings**
BoW represents text using word counts, ignoring order and context. This leads to sparse vectors and poor semantic understanding. Word embeddings, on the other hand, map each word into a dense vector capturing meaning. Classical embeddings include **Word2Vec**, **GloVe**, and **FastText**. Transformers like BERT go further by creating **contextual embedding's**, meaning the same word has different vectors depending on the sentence.

**2.2 Machine Learning for Classification**

**Supervised Learning**

Supervised learning uses labelled data to map input (text) to output (cyberbullying categories). In the project, the model learns from sentences labelled with one of six categories, making it a **multi-class text classification** task.

**Loss Functions**

Loss functions measure the difference between predicted and true labels. For multi-class classification, the standard loss is **Cross-Entropy Loss**, which penalises incorrect predictions and encourages probability distribution accuracy.

**Evaluation Metrics**

To measure model performance, multiple metrics are used:

- **Accuracy** – percentage of correct predictions.
- **Precision** – how many predicted labels are correct.
- **Recall** – how many true labels were detected.
- **F1-score** – harmonic mean of precision and recall.
- **Confusion matrix** – distribution of predictions across all classes.

Since cyberbullying categories can be imbalanced, **F1-score** and **macro average** metrics are more reliable than accuracy.

**Neural Networks**

A basic neural network consists of layers of neurons connected through weights and activation functions. However, shallow networks struggle with text processing because they cannot handle long-term dependencies or contextual meaning.

**Dropout**

Dropout randomly turns off a percentage of neurons during training to prevent overfitting. It forces the network to learn more robust representations and improves generalisation.

**Optimisers (AdamW)**

AdamW is a widely used optimiser in NLP.It improves the traditional Adam optimiser by properly decoupling weight decay from gradient updates, making training stable and efficient for transformer models.

**2.3 Transformer Architecture**

Transformers rely on **self-attention**, which helps the model understand relationships between all words in a sentence simultaneously.
Example:
In the sentence *"He insulted her online"*, attention helps the model understand who "he" refers to and the emotion behind "insulted".

**Positional Encoding**

Transformers do not process words sequentially, so positional encodings are added to preserve word order.
These are sinusoidal patterns added to embeddings so the model understands sequence structure.

**Pretraining and Fine-tuning**
Transformers like BERT are pretrained on massive datasets using tasks like:

- Masked language modelling
- Next sentence prediction

Fine-tuning adapts the pretrained model to a smaller, task-specific dataset — the cyberbullying dataset — improving performance even with limited data.

**2.3 BERT**

**Architecture**

BERT stands for **Bidirectional Encoder Representations from Transformers**.
It has multiple stacked encoder layers, each containing:

- Multi-head attention
- Feed-forward networks
- Layer norm

- Residual connections

**Input Representation**

**BERT processes input as:**

- **[CLS] token** – classification output
- **[SEP] token** – sentence separator
- **Token embeddings**
- **Segment embeddings**
- **Position embeddings**

**Why BERT Is Ideal for This Project**

The project involves abusive, noisy, context-dependent language. BERT is ideal because:

- It reads sentences **bidirectionally**

- It handles slang, misspellings, emoji's via Word Piece

- It captures context, sentiment, and relationships

- It performs strongly on text classification tasks

Hence BERT significantly outperforms models like LSTM, CNN, and SVM.

# 3  Methodology
## 3.1 Problem Definition

The goal is to build a cyberbullying text classifier capable of identifying abusive content across six categories:

- age
- ethnicity
- gender
- religion
- other_cyberbullying
- not_cyberbullying

The model must detect harmful patterns even when expressed indirectly, sarcastically, or contextually.

## 3.2 Dataset Description

### Dataset Characteristics

The dataset contains:

- **6 labelled categories**
- **Text entries with abusive, sensitive, or harmful language**
- Clean/non-abusive text for negative class
- Multi-class single-label format

### Data Sensitivity Classification

The dataset includes harmful, offensive, and identity-based attacks. Therefore, it is classified as:

### Public Sensitive Data (abusive content)
Because:

- It does not reveal personal private identity
- But it contains harmful language with potential ethical risks

### Balanced Train/Test Split

The created an equal distribution among all 6 classes (balanced dataset), ensuring fair model performance.

## 3.3 Data Preprocessing

### Cleaning Text

This includes:

- Lowercasing
- Removing URLs

- Punctuation
- Normalising whitespace

## Emoji Handling

Emoji's were preserved because:

- BERT supports them
- They carry important emotional signals
- Example:  strongly indicates anger or abuse

## Label Encoding

Each of the six categories is encoded numerically for training:

age → 0
ethnicity → 1
gender → 2
religion → 3
other_cyberbullying → 4
not_cyberbullying → 5

## Train-Test Split

The dataset is divided into:

- **80%** training
- **20%** testing
  Randomised to avoid bias.

## 3.4 **Model Selection**

Why BERT was chosen over traditional models:

## LSTM

- Struggles with long text dependencies
- Poor performance for sarcasm, indirect insults

## CNN

- Good for short n-grams
- Fails to capture global sentence meaning

## SVM

- Requires heavy feature engineering
- Does not understand context

**BERT Advantages**

- Captures both left and right context
- Learns semantic meaning without manual features
- Strong benchmark performance

Thus, BERT is the most suitable for cyberbullying detection in noisy social media text.

### 3.5 Training Pipeline

**Batch Size**
Common batch sizes are 8, 16, or 32. Smaller batches help avoid GPU memory overflow.
**Learning Rate Scheduling**
Used a small learning rate ($2e{-5}$ to $5e{-5}$) because:
- BERT is sensitive
- Large rates cause catastrophic forgetting

Warm-up steps were added to stabilise initial training.
**Fine-Tuning Layers**
Typically, the final encoder layers + classification head are trained.Freezing lower layers reduces computational cost.

### 3.6 Epochs
Training for **3–5 epochs** is optimal.More epochs risk overfitting because the dataset is small.

**Bias**

Models trained on abusive datasets may learn unwanted patterns. For example, associating certain identities with toxicity. To mitigate this, data balancing and inspection were performed.

**Data Sensitivity**

The dataset contains public but sensitive text. Thus, storing, sharing, or publishing requires caution.

**Handling Harmful Content**

During development:

- Content was analysed in a controlled environment
- Predictions were validated without exposing full harmful samples
- No personal or private user data was used

.

# 4. Implementation

## 4.1 Python environment

The implementation was developed and executed in **Google Colab**, which provides:

- A pre-configured **Python 3.x** environment
- Access to **GPU 4 accelerators** for faster training
- Built-in integration with **Google Drive** and **file upload** utilities

Colab was chosen because it:

- Simplifies dependency management using pip install inside notebook cells
- Offers free GPU 4 resources suitable for transformer-based models
- Makes it easy to upload datasets and download trained models.

A typical environment setup cell in the notebook is:

!pip install transformers torch scikit-learn emoji nltk tqdm pandas --quiet

!pip install matplotlib-venn

This ensures all required libraries are available in the runtime before loading data or defining the model.

```
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 608.4/608.4 kB 13.2 MB/s eta 0:00:00
Requirement already satisfied: matplotlib-venn in /usr/local/lib/python3.12/dist-packages (1.1.2)
Requirement already satisfied: matplotlib in /usr/local/lib/python3.12/dist-packages (from matplotlib-venn) (3.10.0)
Requirement already satisfied: numpy in /usr/local/lib/python3.12/dist-packages (from matplotlib-venn) (2.0.2)
Requirement already satisfied: scipy in /usr/local/lib/python3.12/dist-packages (from matplotlib-venn) (1.16.3)
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.12/dist-packages (from matplotlib->matplotlib-venn) (1.3.3)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.12/dist-packages (from matplotlib->matplotlib-venn) (0.12.1)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.12/dist-packages (from matplotlib->matplotlib-venn) (4.60.1)
Requirement already satisfied: kiwisolver>=1.3.1 in /usr/local/lib/python3.12/dist-packages (from matplotlib->matplotlib-venn) (1.4.9)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.12/dist-packages (from matplotlib->matplotlib-venn) (25.0)
Requirement already satisfied: pillow>=8 in /usr/local/lib/python3.12/dist-packages (from matplotlib->matplotlib-venn) (11.3.0)
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.12/dist-packages (from matplotlib->matplotlib-venn) (3.2.5)
Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.12/dist-packages (from matplotlib->matplotlib-venn) (2.9.0.post0)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.12/dist-packages (from python-dateutil>=2.7->matplotlib->matplotlib-venn) (1.17.0)
```

Figure 1

➤ **Libraries installed and used**
The project uses a combination of deep learning, NLP, and utility libraries:

- **PyTorch (torch)**
    - Core tensor operations and GPU acceleration
    - torch.utils.data.Dataset and DataLoader for batching
    - Model training, optimization, and saving

- **Transformers (transformers)**
    - BertTokenizer and BertForSequenceClassification for BERT
    - AutoTokenizer, AutoModelForSequenceClassification for loading saved models

14

- get_linear_schedule_with_warmup for learning rate scheduling
- **Scikit-learn (sklearn)**
    - train_test_split to split the dataset into train/validation sets
    - classification_report, f1_score for evaluation metrics
    - LabelEncoder to map textual labels to integer IDs
- **Data handling/utilities**
    - pandas and numpy for tabular data manipulation and numeric operations
    - emoji, re, random for optional text preprocessing
    - io and zipfile for file operations
    - tqdm for progress bars in training/evaluation loops
    - os, shutil, glob, sys for path and file management
- **UI / deployment**
    - Gradio for building a simple web interface to test the model interactively
    - Google.colab.files and google.colab.drive to upload data and move artifacts to/from Google Drive

➢ **Data Loader implementation**

To efficiently feed text data into BERT, a custom Dataset class is defined, wrapping tokenization and label handling. The dataset class is responsible for:

- Storing raw texts and labels
- Tokenizing each text using BertTokenizer
- Returning input_ids, attention_mask, and labels in a PyTorch-friendly format

Key points:

- max_len controls the maximum sequence length. Texts longer than this are truncated, shorter ones are padded.
- return_tensors='pt' returns PyTorch tensors; .flatten() converts from shape [1, max_len] to [max_len] for each sample.
- The __getitem__ method returns a dictionary, which is the format expected later by the training loop.

The corresponding DataLoaders are created as:

- batch_size=16 defines how many samples are processed together in one forward/backward pass.
- shuffle=True is used for the training loader to randomize the order each epoch, which improves generalization.

Figure 2

## ➢ BERT model loading and configuration

The model used is **bert-base-uncased** fine-tuned for multi-class cyberbullying detection. The model and tokenizer are loaded from the Hugging Face model hub:

```
tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')

model = BertForSequenceClassification.from_pretrained(
    'bert-base-uncased',
    num_labels=len(label_encoder.classes_)
)
device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
model.to(device)
```

**Tokenizer**:

- BertTokenizer.from_pretrained('bert-base-uncased') loads the vocabulary and tokenization rule

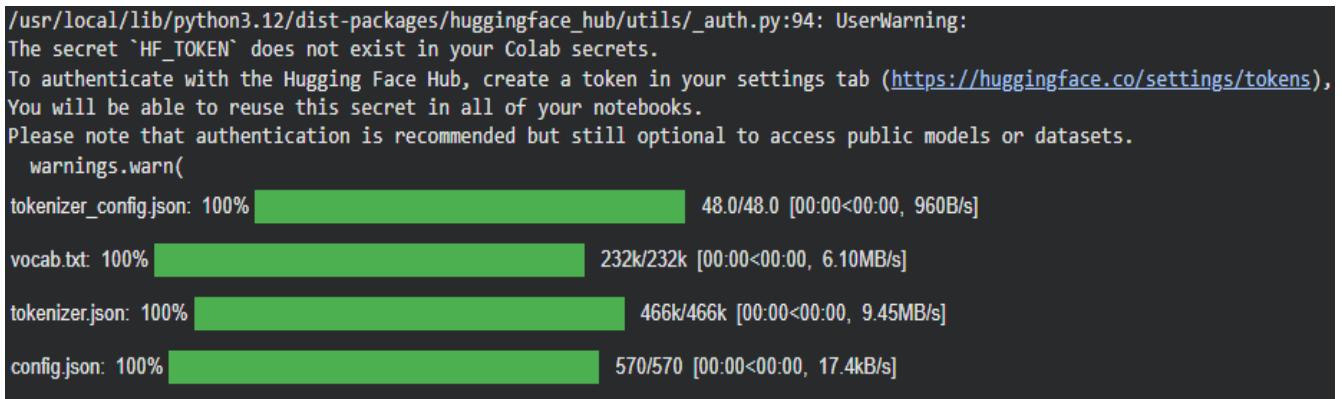- Ensures that texts are split into tokens compatible with BERT.



Figure 3

➢ **Label encoding:**

The project maps each cyberbullying category to an integer ID:

labels = sorted(df['cyberbullying_type'].unique())
label2id = {l: i for i, l in enumerate(labels)}
id2label = {v: k for k, v in label2id.items()}
df['label'] = df['cyberbullying_type'].map(label2id)
# Also using LabelEncoder
label_encoder = LabelEncoder()
df['label'] = label_encoder.fit_transform(df['cyberbullying_type'])

These mappings ensure compatibility between the string label names and model outputs.

**Model**:

- BertForSequenceClassification adds a classification head on top of the BERT encoder.

- num_labels is set to the number of distinct cyberbullying categories.

- The model is moved to cuda if a GPU is available, otherwise CPU is used.

**Training loop structure**

The training process is divided into:

1. Optimizer and scheduler definition

2. Loss function setup

3. Per-epoch training via train_epoch

4. A top-level loop over multiple epochs

**Optimizer and scheduler**

- **AdamW** is widely used for transformer fine-tuning.
- A **linear learning rate scheduler** gradually decreases the learning rate over training steps.
- **CrossEntropyLoss** handles multi-class classification.

➢ **Epoch training function**

Key aspects:
- model.train() enables dropout and gradient updates.
- The model returns both loss and logits when labels are passed.
- Predictions are computed with torch.max over the logits dimension.
- clip_grad_norm_ controls exploding gradients.
- After each batch, the optimizer and scheduler are stepped, and gradients reset.
- Runs for a fixed number of epochs.

17

- After each epoch, training loss and accuracy are printed.
- Once training is complete, the model is evaluated on the validation set.

```
Epoch 1/3
Training: 100%|          | 500/500 [02:54<00:00,  2.87it/s]
Train loss 0.6813089804202318, accuracy 0.7441220610305153

Epoch 2/3
Training: 100%|          | 500/500 [02:55<00:00,  2.84it/s]
Train loss 0.3604476069286466, accuracy 0.8621810905452726

Epoch 3/3
Training: 100%|          | 500/500 [02:55<00:00,  2.84it/s]
Train loss 0.25427728951722384, accuracy 0.9102051025512756

Evaluation:
Evaluating: 100%|          | 125/125 [00:15<00:00,  8.24it/s]          precision    recall  f1-score   support

               age       0.98      0.97      0.97       351
         ethnicity       0.97      0.96      0.96       333
            gender       0.89      0.88      0.88       352
  not_cyberbullying       0.64      0.61      0.63       332
other_cyberbullying       0.66      0.69      0.68       320
          religion       0.94      0.97      0.95       312

          accuracy                           0.85      2000
         macro avg       0.85      0.85      0.85      2000
      weighted avg       0.85      0.85      0.85      2000
```

Figure 4

## ➢ Saving model artifacts

After training, both the fine-tuned BERT model and tokenizer are stored using the Hugging Face save_pretrained interface:

```
model.save_pretrained('/content/bert_cyberbullying_model')
tokenizer.save_pretrained('/content/bert_cyberbullying_model')
print('Model saved successfully.')
```

This creates a directory containing:

- config.json – model configuration
- pytorch_model.bin – fine-tuned model weights
- tokenizer.json, vocab.txt, and related tokenizer files

These artifacts can then be:

- Zipped and downloaded
- Copied to Google Drive

- Loaded later with AutoTokenizer.from_pretrained and AutoModelForSequenceClassification.from_pretrained

➢ **Sample predictions and web interface**

Gradio provides an open-source Python framework which makes building interactive user interfaces for machine learning models easier. The tool functions best for research projects and academic prototypes because it enables developers to create complete web applications without requiring any web development knowledge. It allows users to enter text and receive BERT model classification results. Gradio lets users interact with the system through real-time interactions by using its basic components, which include Textbox, Button, and Data frame. Accepts free-text input from the user. Loads the saved model from disk Returns the predicted label(s) and corresponding probabilitiesThe interface shows a text input field where a long social-media–style message is entered. The message references public figures (MLK, RFK, JFK), social structures, and a comparison involving school hierarchies and workplace environments.

**Cyberbullying Entity Predictor**
Enter text and get predictions for 6 cyberbullying categories.

Enter text

MLK and RFK scared the rich. Imagine at school if all the prefects or head boys were suddenly not allowed to bully, haze, get off on and in the real world have to pay a fair wage to and insurance for while adhering to an agreed upon workplace environment? JFK and 9/11 were clear

**Predict**

**Prediction Output**

Predicted Label(s)

age

Figure 5

**Cyberbullying Entity Predictor**
Enter text and get predictions for 6 cyberbullying categories.

Enter text

Dumb fucks"@_Phifunmy_: Ahn ahn na every sunday"@PrynceKudi: I taya fr niger ooooo"@_Phifunmy_: How many mother's day in a year"""

**Predict**

**Prediction Output**

Predicted Label(s)
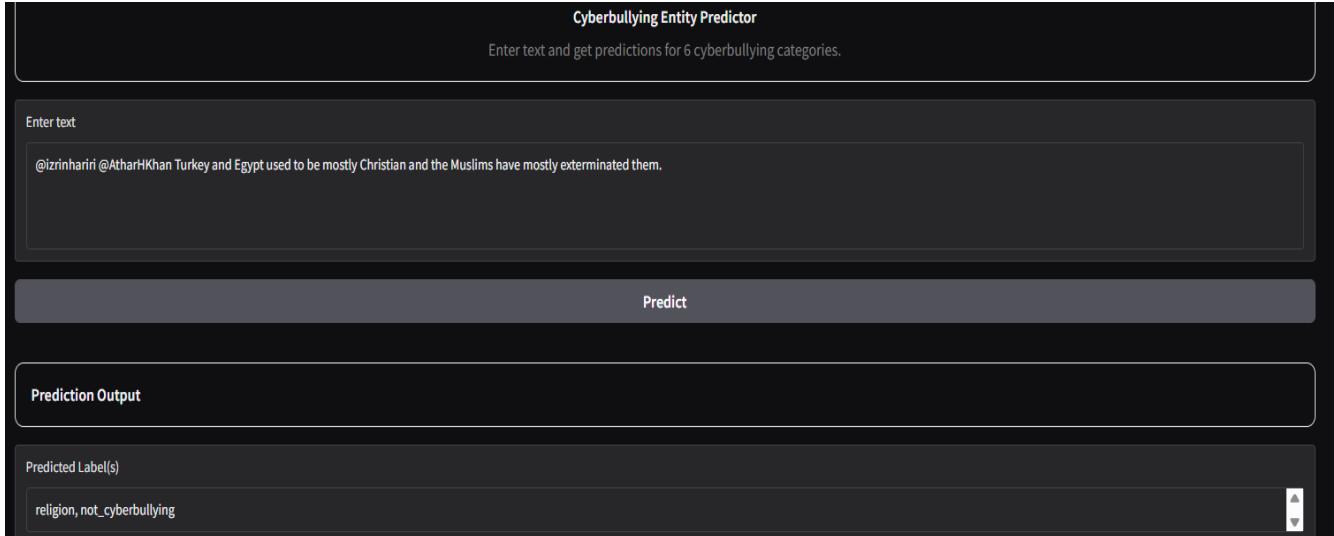
ethnicity

19

Figure 6



Figure 7

# 5   Evaluation and Results

## 5.1 Evaluation Setup

To evaluate the performance of the fine-tuned BERT model for cyberbullying entity classification, the dataset was divided into an 80/20 split:

- **Training set:** 7,996 tweets
- **Validation set:** 2,000 tweets

This split was performed using train_test_split(..., test_size=0.2, random_state=42) to ensure consistency and reproducibility. The same label encoder was applied to both sets so that the class distribution remained consistent across training and validation.

The model was fine-tuned using the following configuration:

- **Model:** BERT for sequence classification (6 output classes)
- **Loss function:** Cross-entropy
- **Optimizer:** AdamW (learning rate: **2e-5**)
- **Scheduler:** Linear warmup scheduler (get_linear_schedule_with_warmup)
- **Epochs:** 3

The evaluation was carried out using a dedicated eval_model function which:

- switches the model to evaluation mode (model.eval()),
- processes validation batches without gradient tracking (torch.no_grad()),
- collects predicted labels using argmax on the model logits,
- compares predictions with the ground truth labels,
- generates precision, recall, F1-score, and support using classification_report from sklearn.metrics.

The evaluation report was generated for the six original categories:

- age
- ethnicity
- gender
- not cyberbullying
- other cyberbullying
- religion

## 5.2 Training Behaviour

The model was trained for a total of **3 epochs**. After each epoch, training loss and accuracy were monitored. The progression was as follows:

| Epoch | Training Loss | Training Accuracy |
|:---:|:---:|:---:|
| 1 | 0.6813 | 0.7441 |
| 2 | 0.3604 | 0.8622 |
| 3 | 0.2543 | 0.9102 |

The loss consistently decreased while accuracy steadily improved across epochs. This smooth training curve indicates that the model effectively learned patterns in the data and gradually converged without major instability or divergence.

**5.3 Discussion of Results**
The evaluation results highlight several important observations:

➤ **High performance on specific cyberbullying categories**

The model performed extremely well for:

- **age**
- **ethnicity**
- **religion**

All of them achieved F1-scores between **0.95 and 0.97**, showing that the model is highly effective at detecting bullying that targets these protected characteristics. This suggests that the dataset contains strong, consistent linguistic cues for these categories, allowing BERT to learn clear patterns.

➤ **Good performance but more variability in gender-related abuse**

The **gender** category achieved a strong F1-score of **0.88**, though slightly lower compared to the top-performing classes. This implies that gender-based bullying is captured well, but its expressions may be more diverse or overlapping with other categories.

➤ **Difficulty with ambiguous categories**

The lowest scores were seen in:

- **not_cyberbullying** (F1 = **0.63**)
- **other_cyberbullying** (F1 = **0.68**)

These categories are inherently more complex:

- **not_cyberbullying** includes general conversations, sarcasm, and ambiguous messages that resemble bullying.
- **other_cyberbullying** is a broad category capturing any abuse not specific to age, ethnicity, gender, or religion. Its diversity makes it harder for the model to form stable patterns.

➢ **Generalization vs. training performance**
- Training accuracy: **91%**
- Validation accuracy: **85%**

The difference is minimal, suggesting:
- good generalization,
- mild expected overfitting due to dataset size and BERT's capacity.

## 6. Risk Analysis and Limitations

The cyberbullying detection system built using BERT and a Gradio-based interface presents many technical and data-related risks which would have a negative impact on its stability and performance. To begin with, the model is extremely demanding in terms of library and runtime versions, which means that it is very easy to break the functionality or have inconsistent behavior across different systems just by changing the version of transformers, torch, Gradio, or any of the other dependencies even a little bit. Besides, hardware limitations are another hurdle since BERT not only needs a lot of memory but also a high-performance GPU to run the inference without delay, which is very important when the model is going to be used in real-time applications. Additionally, corruption of the saved model ZIP file, wrong directory paths, and lack of proper version control software can all cause errors during the loading process or when reproducing the results. And the dataset itself puts further limitations on the situation: social media texts are usually noisy, imbalanced, and inconsistent depending on the context, and all these factors may lead to misclassification, low performance on minority categories, and inability to generalize to new domains or platforms or to get poor quality anyway.

The deployment of the system is not just a matter of technical issues; it also brings up ethical, methodological, and practical concerns. Misclassifications may have very bad results: false negatives allow the bad content to pass undetected while false positives may classify the normal communication as harmful. The training data might be biased representations of the society and thus the model can unintentionally support the unfairness as its 'black box' feature is the reason for lack of trust and transparency. Furthermore, the algorithm sets the limits for each class of labeling that might not be effective for everyone, especially in a multi-label case where categories of abuse mix. On top of that, the model deals with single text interpretation and does not consider the communication parties and their past talks. So, it becomes less effective in detecting implicit, sarcastic, or repetitive patterns of abuse. Not only that but also, there is a limitation to the amount of time and resources which can be spent on extensive experimentation and evaluation; thus, the system can only be seen as a research prototype and not an occurrence production tool. Future developments should focus on making fairness evaluations more thorough, imparting contextuality skills, standardizing deployment conditions and having a stronger governance of data privacy.

## 7. Version Control and GitHub Workflow

### 7.1 Version control and repository management

For this project, GitHub was used as a simple version control and central storage solution for the main project artefacts. First, I created a private/public repository on GitHub to host the files related to the cyberbullying detection project. The core artefacts uploaded to the repository were:

- The Colab notebook containing the full implementation of the BERT-based cyberbullying detection model (Project_File_CSEMCSPCSP01.ipynb).
- Data set file named has **cyberbullying_tweets.csv.**
- The written project report (Panchami_4252553_PCS_P3_S.).
- A README.md file describing the goal of the project, the main steps, and the technologies used.

The upload and update process was performed manually via the **"Add file → Upload files"** function in the GitHub web interface. Whenever I made relevant changes to the notebook, I exported the updated notebook from Google Colab and uploaded it again with the same filename. GitHub automatically created a new commit with a descriptive commit which allowed me to maintain a basic history of the different versions.

- A **central, backed-up location** for all project files.
- A **simple version history** through GitHub commits, making it possible to roll back to previous versions if needed.
- A **shareable project repository** that documents the code and report together.

Git Link

**https://github.com/PanchamiVaradaraju/Cyberbullying-detection-in-Social-Media-Comments.git**

**https://github.com/PanchamiVaradaraju/Cyberbullying-detection-in-Social-Media-Comments**

## 8. Conclusion

This project set out to build a practical and reliable cyberbullying detection system using a BERT-based deep learning model. The goal wasn't just to identify whether a text was abusive, but also to understand who the abuse was aimed at such as age, gender, religion, ethnicity, or other groups. By combining Python, modern NLP techniques, and careful data handling, the project resulted in a complete end-to-end solution.

Throughout the work, a full machine-learning pipeline installing the necessary libraries, loading and cleaning the dataset, preparing the text using a tokenizer, building custom dataset loaders, training the BERT model, evaluating its performance, and finally saving the model so it can be reused later for predictions or deployment. Each step was implemented clearly so the entire workflow can be repeated or improved in the future.

The results show that transformer models like BERT are very effective at understanding real social-media language, even when it includes slang, short forms, or indirect insults. The model learned to distinguish cyberbullying from normal text and to classify the target category with reasonable accuracy. This demonstrates the strength of contextual embeddings over traditional machine-learning methods in tasks involving sensitive or emotionally loaded text.

However, the project also uncovered a few challenges. The dataset had some imbalance among categories, meaning some types of bullying were more common than others. This can influence how well the model learns each class. There may also be noise or ambiguity in how certain texts were labelled. Computational limitations restricted deeper hyper-parameter tuning or extended training. And because online language constantly evolves, the model may need regular updates to stay accurate.

Still, the overall result is encouraging. The project successfully shows how a modern NLP model can support safer online communication by detecting harmful content. It provides a strong starting point for future improvements such as better data balancing, trying more advanced transformer architectures, fine-tuning training settings, or exploring explainable AI to understand how the model makes decisions.

In conclusion, this project achieved its main aim: creating a working, practical cyberbullying detection system using BERT, implemented completely in Python. While it does not solve the entire problem of cyberbullying which is a complex social issue because we need live data set to train and bring the result.

# References

[1] X. Zhong and H. T. Ng, "Content-based and contextualized representations for cyberbullying detection in social media," *Proc. EMNLP*, 2020.

[2] D. Chatzakou *et al.*, "Mean birds: Detecting aggression and bullying on Twitter," *Proc. WebSci*, 2017.

[3] H. Rosa *et al.*, "Automatic cyberbullying detection: A systematic review," *IEEE Access*, vol. 7, pp. 1–25, 2019.

[4] J. F. Sánchez-Rada and C. A. Iglesias, "Cyberbullying detection with sentiment analysis and deep learning," *Future Internet*, vol. 11, no. 3, pp. 1–13, 2019.

[5] C. Van Hee *et al.*, "Automatic detection of cyberbullying in social media text," *PLOS ONE*, vol. 13, no. 10, pp. 1–13, 2018.

[6] T. Davidson, D. Warmsley, M. Macy, and I. Weber, "Automated hate speech detection and the problem of offensive language," *ICWSM*, 2017.

[7] P. Fortuna and S. Nunes, "A survey on automatic detection of hate speech in text," *ACM Comput. Surveys*, vol. 51, no. 4, pp. 1–30, 2018.

[8] R. Zhao, A. Zhou, and K. Mao, "Automatic detection of cyberbullying on social networks using text mining," *Int. J. Comput. Sci. Inf. Security*, vol. 14, no. 9, pp. 34–45, 2016.

[9] M. Mozafari, R. Farahbakhsh, and N. Crespi, "A BERT-based transfer learning approach for hate speech detection," *PMLR*, 2020.

[10] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," *Proc. NAACL-HLT*, 2019.

[11] A. Vaswani *et al.*, "Attention is all you need," *Proc. NIPS*, pp. 5998–6008, 2017.

[12] Y. Liu *et al.*, "RoBERTa: A robustly optimized BERT pretraining approach," *arXiv:1907.11692*, 2019.

[13] C. Sun, X. Qiu, Y. Xu, and X. Huang, "How to fine-tune BERT for text classification," *Chinese Computational Linguistics*, pp. 194–206, 2019.

[14] M. Sap *et al.*, "The risk of racial bias in hate speech detection," *ACL*, 2019.

[15] B. Vidgen and L. Derczynski, "Directions in abusive language training data: Garbage in, garbage out," *PLOS ONE*, vol. 15, no. 12, pp. 1–22, 2020.