

Project Report: AI-Generated Content Detection Platform

Contents

1	Introduction	3
2	Model Description	3
2.1	Chosen AI Model	3
2.2	Design and Development Process	3
2.2.1	Data Collection	3
2.2.2	Data Preprocessing	3
3	User Interface Development Tools	3
3.1	Front-end Development	3
3.2	Back-end Development	4
4	Front-end UI	4
5	Back-end and model	4
6	Workflow diagram	5
7	Comparative Analysis Table	6

1 Introduction

As fraud methods evolve, from financial fraud, identity theft, credit card misuse, and even plagiarism, modern technology uses advanced techniques such as machine learning and artificial intelligence to detect these frauds.

This project involves the development of a website or mobile application capable of identifying whether a given text or script is generated by AI tools. It combines natural language processing (NLP) and web development.

2 Model Description

2.1 Chosen AI Model

The core of the detection system is based on transformer-based models such as BERT. BERT is a pre-trained deep learning model designed to process and understand text by capturing context and meaning in a bidirectional manner. This model is trained to recognize patterns in text that are often indicative of AI-generated content.

2.2 Design and Development Process

2.2.1 Data Collection

- **Gather Text Data:**
 - *Sources for AI-generated texts:* GPT models (e.g., ChatGPT), open-source datasets (Kaggle `AI_Human.csv`), GitHub.
 - *Sources for human-written texts:* Books (e.g., *Rich Dad Poor Dad*), articles (e.g., news from News in Levels).
- **Label Data:** Ensure that the collected texts are properly labeled as *AI-generated* or *human-written*.

2.2.2 Data Preprocessing

- **Text Cleaning:** Remove any irrelevant information (e.g., HTML tags, special characters).
- **Tokenization:** Split the text into words or sentences for easier analysis.
- **Normalization:** Convert text to a consistent format (e.g., lowercasing, stemming).

3 User Interface Development Tools

3.1 Front-end Development

- **HTML/CSS/JavaScript:** For creating the web UI.

3.2 Back-end Development

- Python: For the backend logic and model implementation.
- Flask: For building APIs and managing model requests.

4 Front-end UI

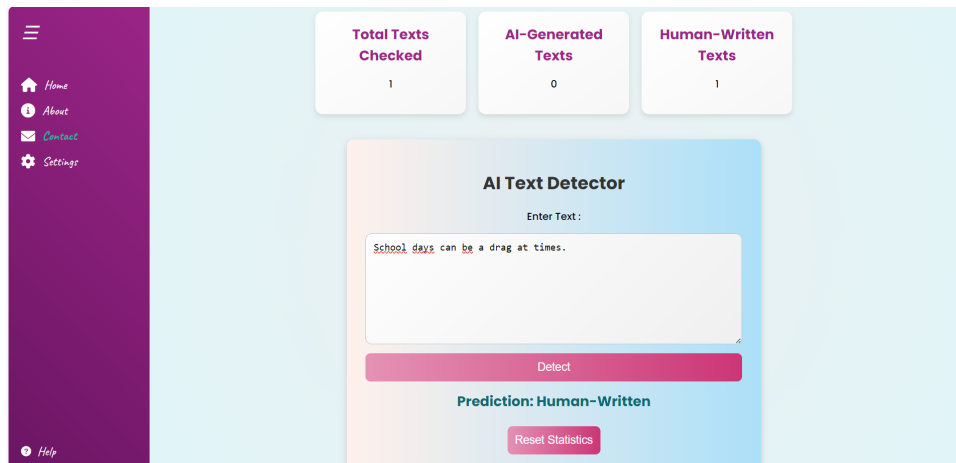


Figure 1: Home page UI design

5 Back-end and model

```

Run Terminal Help
app.py x text_data.csv
app.py > predict
32 def predict():
33     try:
34         data = request.get_json()
35         if not data or 'content' not in data:
36             return jsonify({"error": "No content provided"}), 400
37
38         text = data['content']
39
40         # Check if the input text is in the dataset
41         dataset_result = is_in_dataset(text)
42         if dataset_result:
43             return jsonify({"result": dataset_result})
44
45         # If not in dataset, use the model to predict
46         inputs = tokenizer(text, return_tensors="pt", padding=True, truncation=True)
47         with torch.no_grad():
48             outputs = model(**inputs)
49             prediction = torch.softmax(outputs.logits, dim=1)
50
51         # Get the predicted label
52         ai_generated_prob = prediction[0][1].item()
53         result = "AI-Generated" if ai_generated_prob > 0.5 else "Human-Written"
54
55         return jsonify({"result": result})
56
57     except Exception as e:
58         print("An error occurred:", e)
59         return jsonify({"error": "An internal error occurred"}), 500
60
61 if __name__ == '__main__':
62     app.run(host='0.0.0.0', port=8000, debug=True)
63

```

Figure 2: flask

```

Run Terminal Help
app.py x train_model.py
model > train_model.py > ...
29
30 # Define Dataset class
31 class TextDataset(torch.utils.data.Dataset):
32     def __init__(self, encodings, labels):
33         self.encodings = encodings
34         self.labels = labels
35
36     def __len__(self):
37         return len(self.labels)
38
39     def __getitem__(self, idx):
40         item = {key: val[idx] for key, val in self.encodings.items()}
41         item['labels'] = torch.tensor(self.labels[idx], dtype=torch.long)
42         return item
43
44 # Prepare datasets
45 train_dataset = TextDataset(train_encodings, train_labels)
46 val_dataset = TextDataset(val_encodings, val_labels)
47
48 # Load pre-trained BERT model
49 model = BertForSequenceClassification.from_pretrained("bert-base-uncased", num_labels=2)
50
51 # Set up training arguments
52 training_args = TrainingArguments(
53     output_dir="./model",
54     num_train_epochs=3,
55     per_device_train_batch_size=8,
56     per_device_eval_batch_size=8,
57     warmup_steps=500,
58     weight_decay=0.01,
59     logging_dir="./logs",
60     logging_steps=10,
61     evaluation_strategy="epoch",
62     save_strategy="epoch", # Ensure models are saved after each epoch
63     save_total_limit=1, # Keep only the latest checkpoint
64 )
65

```

Figure 3: Model training

6 Workflow diagram

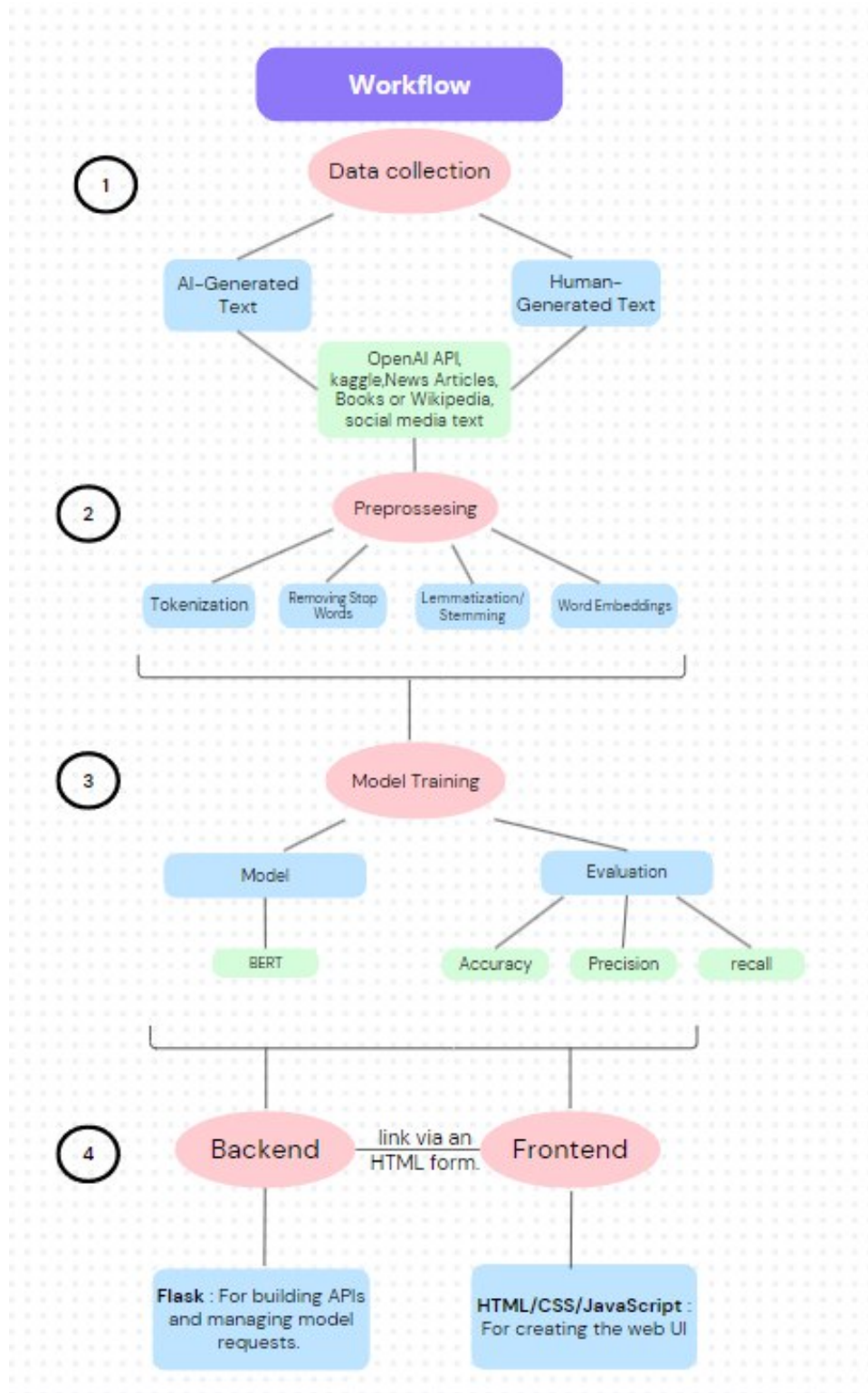


Figure 4: Workflow diagram modeling

7 Comparative Analysis Table

Study	Approach/Method	Model/Algorithm	Dataset	Key Metrics	Findings
GPT-2 Output Detector	Uses fine-tuned transformers	GPT-2 fine-tuned model	OpenAI dataset (Reddit)	Accuracy, Precision, Recall	Fine-tuning GPT-2 improves detection of its own outputs, achieving high accuracy on similar models
RoBERTa for Detection	Fine-tunes RoBERTa on labeled text data	RoBERTa	Custom dataset of AI-generated text	F1 Score, Accuracy	Outperforms traditional methods with a balanced accuracy on longer texts
GLTR (Giant Language Model Test Room)	Statistical anomaly detection based on token likelihoods	BERT-based model	Custom dataset from diverse models	Probability distribution, PPL	Highlights statistical anomalies in AI text, visual aid for human-assisted verification
Turing-NLG Detection	Uses Turing Natural Language Generation model	Turing-NLG-based detection	OpenAI GPT-2 and GPT-3 datasets	Accuracy, Specificity	Provides high specificity in distinguishing text generated by large language models
Grover Detector	Fine-tuned language model specifically for AI-generated news	Grover	AI-generated news articles (Grover dataset)	Accuracy, F1 Score	High accuracy on news detection, emphasizes importance of training on AI-generated news
DetectGPT	Utilizes log probability-based features	Log probability comparison method	GPT-3 generated texts	F1 Score, Accuracy	Employs probability mismatches to flag AI text, works well across different GPT models
OpenAI Classifier	Fine-tuned on diverse text corpora	Transformer-based classifier	Custom dataset from AI and human texts	Precision, Recall, Accuracy	Effective in general scenarios, but less accurate on short and complex texts

Table 1: Comparative Table of Previous Work in Detecting AI-Generated vs. Human-Written Text