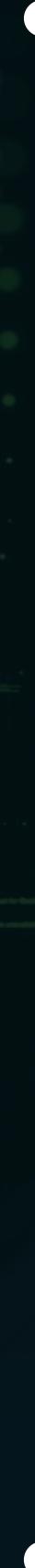


Unmasking the AI chatbot

Detecting AI generated content

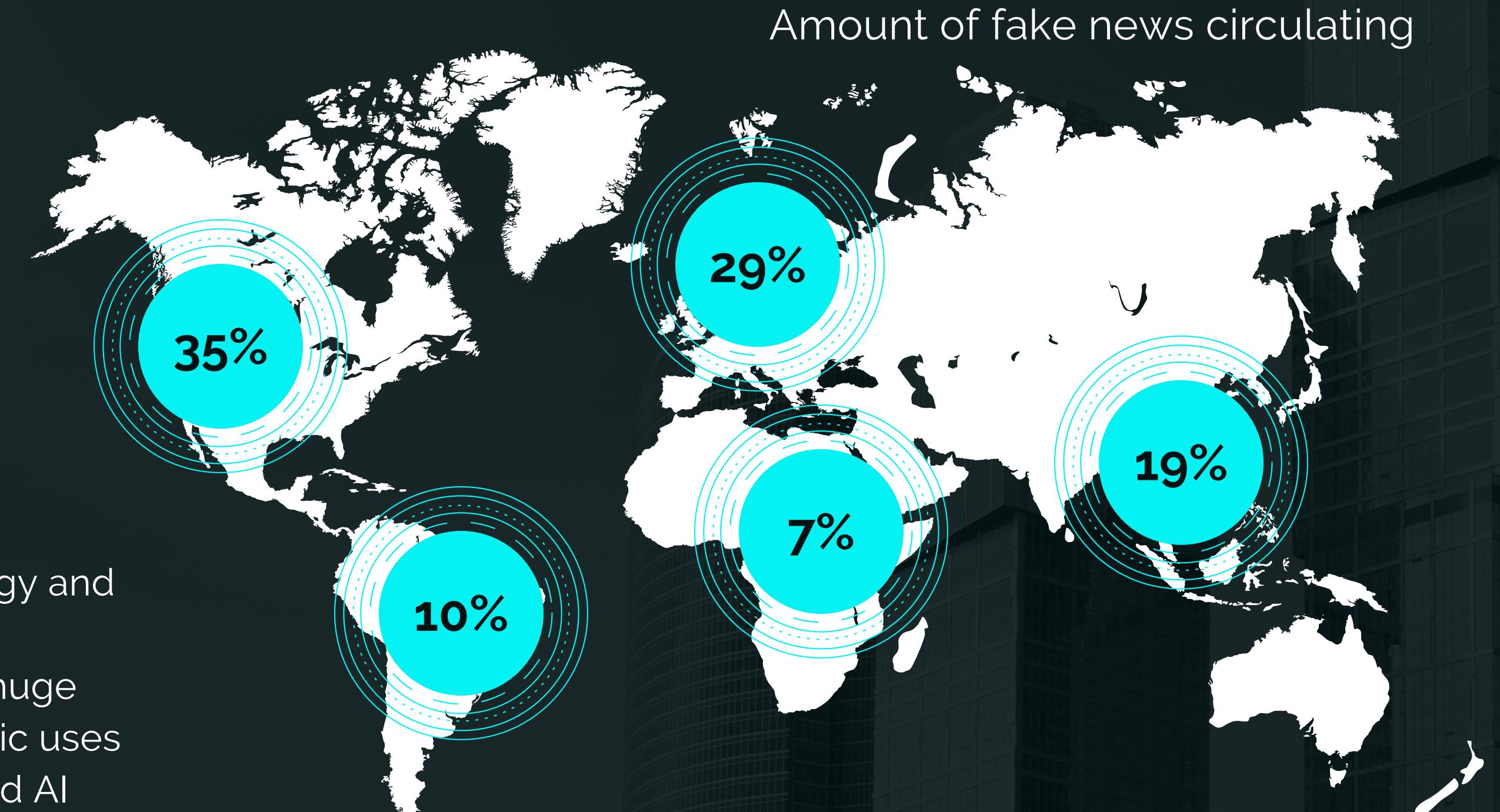
Saswat Susmoy Sahoo
Dev Mishra

Table Of Content

- 
- 01** Need of the hour
 - 02** About the Tech
 - 03** Key Features
 - 04** Training and Evaluation
 - 05** Problems faced
 - 06** Potential Applications
 - 07** Future Work
 - 08** Codebase
 - 09** Questions

Need of the hour

With the advancing technology and mass adoption of the same, misutilization can become a huge problem. As the general public uses AI chatbots like ChatGPT, Bard AI etc, the chances of misutilising the resources increases.



North A
36%

South A
24%

Africa
16%

Asia
28%

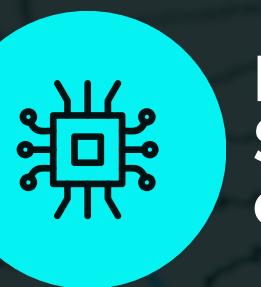
Europe
32%

Our AI model

We started out designing this model using various Python Libraries like Numpy, Pandas, Scikit Learn, Tensorflow, Keras etc.

We used Numpy and Pandas for Data Handling purposes and Tensorflow with Keras for importing LSTM / Logistic Regression models. Scikit Learn was used to train and test our model and also to check our accuracy.

To implement concepts of OOPs we designed our GUI using PyQt5 which is based on Python Programming Language.



**Numpy, Pandas,
Scikit Learn, Keras
etc**



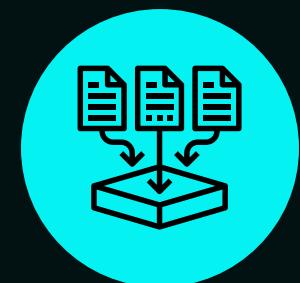
**Model Used - LSTM /
Logistic Regression****

** Logistic Regression isn't used in the working model rather it was used to compare accuracies in the early stage of development of this project

Key Features

Our Model uses Deep Learning concepts of LSTM which produces quite accurate results in the fields of NLP and Cognitive Learning.

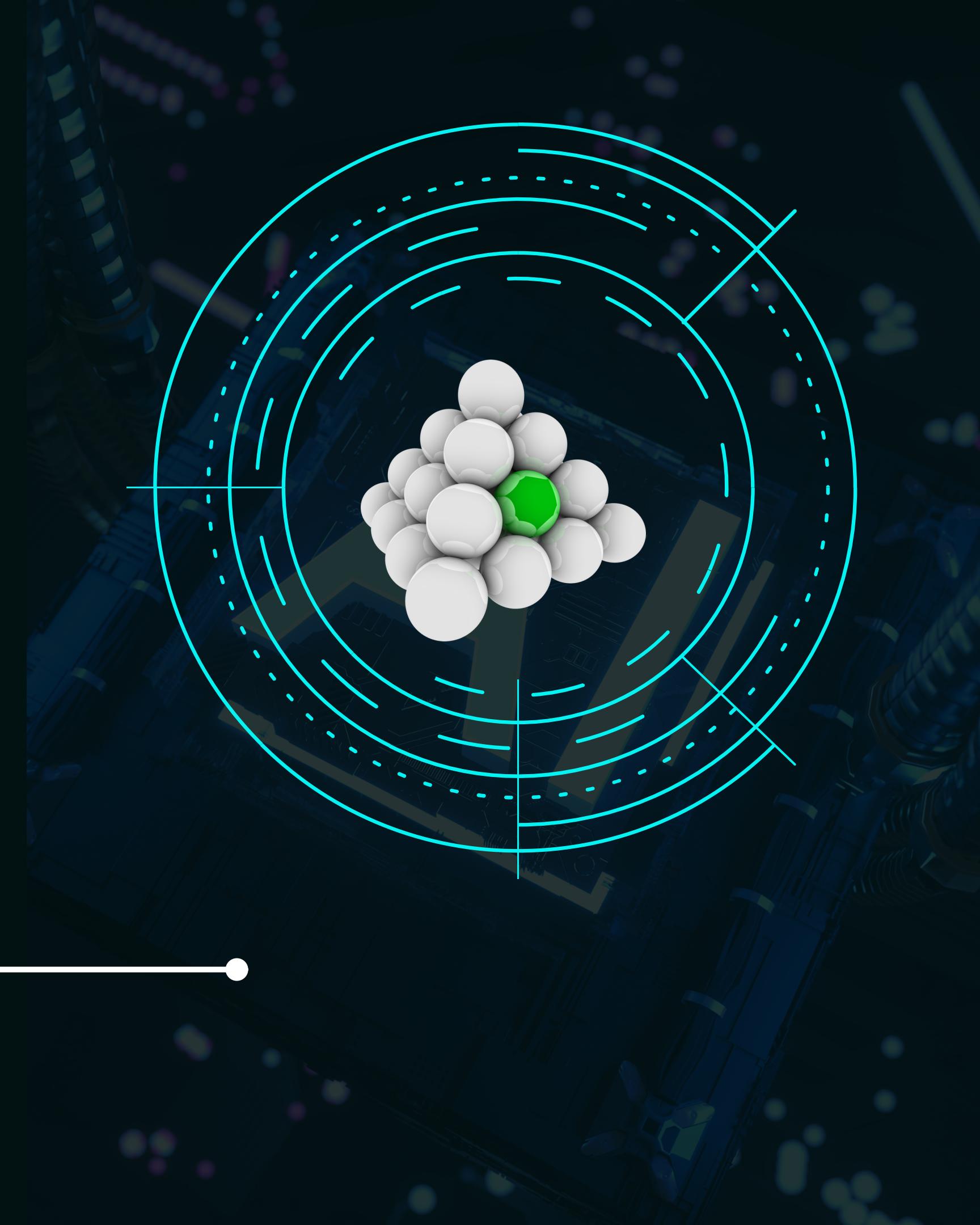
As our model is completely based out on Python, it will be easier and convenient for us to optimize our model further in future and also add more functionalities.



**Data Used
8-10k**



**Accuracy
85-95** %**



** Accuracy above 90% sometimes denotes that the model is overfitting which might be in our case. So Accuracy statistics should not be considered as the sole basis to judge our project

Training and Evaluation

AI Generated

Apply page animations and transitions to your Canva presentation to emphasize

51%

Content Used

AI Generated Essays - 20%
ManaGPT prompt responses - 35%
Miscellaneous - 45%



Human generated

Apply page animations and transitions to your Canva presentation to emphasize

49%

Content Used

- Covid research abstracts 30%
- MacD consumer reviews 30%
- Medium Articles 10%
- US Prez speeches >1%
- BBC News RSS feed >1%
- Miscellaneous ~ 30%

Problems Faced



Data Handling

Due to our less exposure and knowledge in the field of Data Analytics, we found it difficult to handle such huge amount of data that we could access.



Cases of Overfitting and Underfitting

As we were starting out, we faced problems of Overfitting and Underfitting as we couldn't set our parameters to optimize our model and sometimes ample/shortage of data created the problem.



Less Knowledge

Due to less exposure and knowledge from our side in the ML field, we at times found it difficult to solve problems and understand concepts. We could not optimize the LSTM model to match our requirements properly.

Potential Application

We initiated this idea with the aim to detect AI generated content but as we entered the development phase we realised the potential this small idea holds in various streams and fields.



Detecting Plagiarism
Versions of this model can be used to detect plagiarism in the content creating industry.



Detecting Fake News
This model can be trained on Real & Fake News data to enable it to detect fake news which can be implemented in social media sites to prevent misinformation

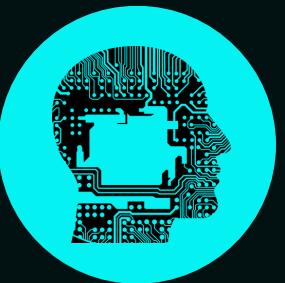


College Applications
A very accurate and precise version of this model can be used to detect fake/AI generated college essays or Letters of intent



Forensics
The model aids in identifying manipulated or tampered media, assisting in investigations and legal proceedings.

Future Work



Latest Technology

As AI algorithms improve, it becomes increasingly difficult to distinguish between AI-generated and human-created content.

We would like to explore various options to implement methods and integrate other technologies to keep our model updated as per industry standards.



Image Detection

We would like to explore the option to detect AI generated images in the future.



Codebase



```
1 from PyQt5.QtWidgets import QApplication, QWidget, QVBoxLayout, QPushButton, QTextEdit, QLabel  
2 from PyQt5.QtCore import Qt  
3 import sys  
4 import numpy as np  
5 import pickle  
6 import pandas as pd  
7 from sklearn.model_selection import train_test_split  
8 from sklearn.metrics import accuracy_score  
9 from tensorflow.keras.preprocessing.text import Tokenizer  
10 from tensorflow.keras.preprocessing.sequence import pad_sequences  
11 from tensorflow.keras.models import Sequential  
12 from tensorflow.keras.layers import Embedding, LSTM, Dense, Dropout  
13 from tensorflow.keras.models import load_model
```

This contains all Libraries and Modules imported



```
1 df = pd.DataFrame(data)
2 loaded_model = load_model('model3_extended.h5')
3 # preprocessing
4 tokenizer = Tokenizer()
5 tokenizer.fit_on_texts(df['text'])
6 vocab_size = len(tokenizer.word_index) + 1
7
8 sequences = tokenizer.texts_to_sequences(df['text'])
9 max_sequence_length = max([len(seq) for seq in sequences])
10 padded_sequences = pad_sequences(sequences, maxlen=max_sequence_length)
```

This loads the data and load the trained model saved as "model3_extended". This also processes the data before training the model

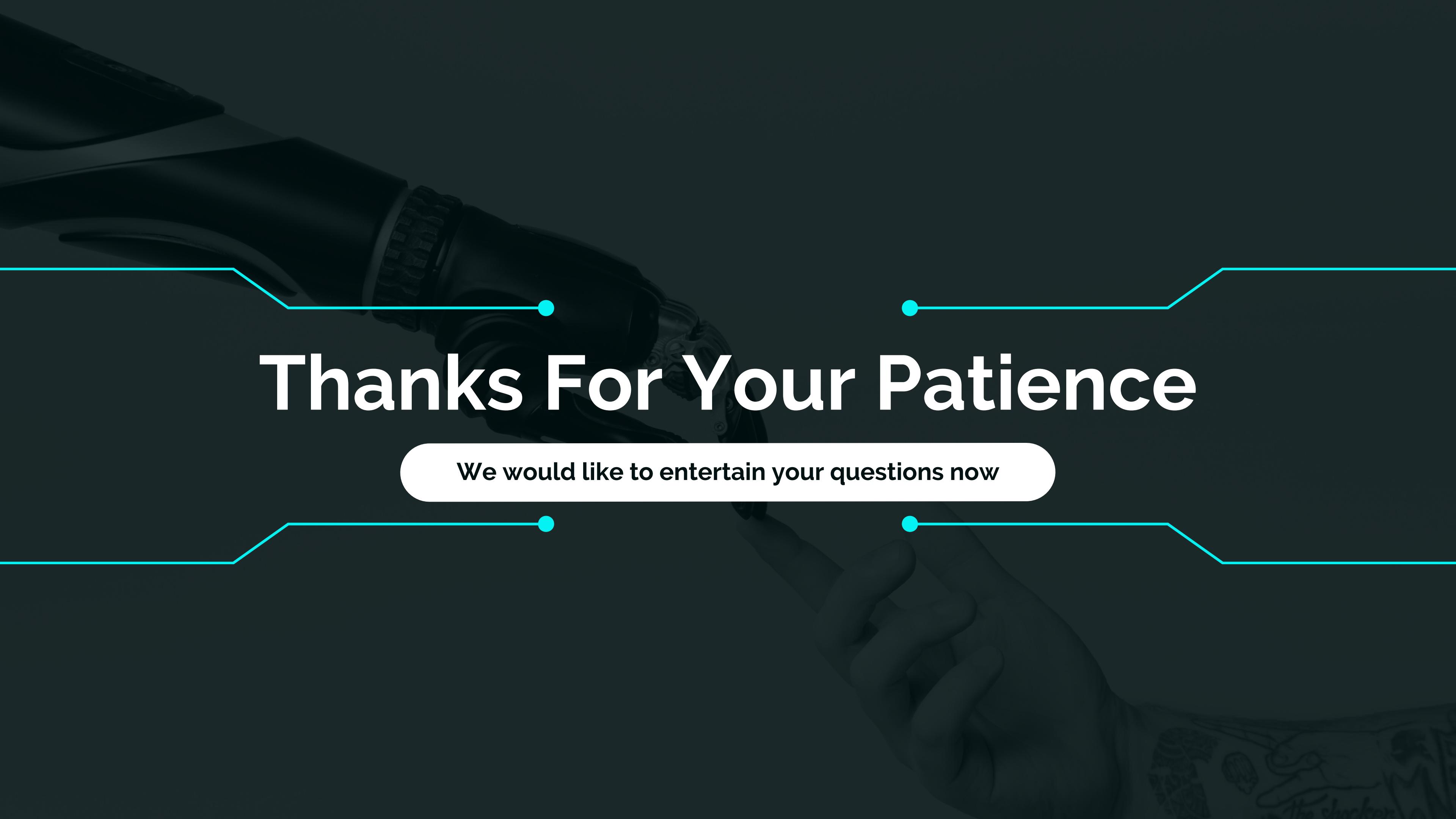


```
1 #dl model
2 model = Sequential()
3 model.add(Embedding(vocab_size, 100, input_length=max_sequence_length))
4 model.add(LSTM(128))
5 model.add(Dropout(0.5))
6 model.add(Dense(1, activation='sigmoid'))
7
8 model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
```

This loads the Deep Learning model (LSTM) and gives various parameters to tune the same.

```
1  class AIChatbotPredictor(QWidget):
2      def __init__(self):
3          super().__init__()
4
5          self.initUI()
6
7      def initUI(self):
8          layout = QVBoxLayout()
9
10         self.label = QLabel("Enter Content:")
11         layout.addWidget(self.label)
12
13         self.text_edit = QTextEdit()
14         layout.addWidget(self.text_edit)
15
16         self.button = QPushButton("Predict")
17         self.button.clicked.connect(self.predict_content)
18         layout.addWidget(self.button)
19
20         self.setLayout(layout)
21
22     def predict_content(self):
23         content = self.text_edit.toPlainText()
24         content_sequence = tokenizer.texts_to_sequences([content])
25         content_padded = pad_sequences(content_sequence, maxlen=max_sequence_length)
26         prediction = model.predict(content_padded)
27
28         if prediction[0] >= 0.5:
29             result = "The content is produced by an AI chatbot."
30         else:
31             result = "The content is not produced by an AI chatbot."
32
33         self.label.setText(result)
34
35     def main():
36         app = QApplication(sys.argv)
37
38         window = AIChatbotPredictor()
39         window.show()
40
41         sys.exit(app.exec_())
42
43 if __name__ == '__main__':
44     main()
```

This is the GUI for the program



Thanks For Your Patience

We would like to entertain your questions now