

Configuration Manual

MSc Research Project
MSc in Artificial Intelligence

Mehwish Mohammed Hanif Khatib
Student ID: X23398396

School of Computing
National College of Ireland

Supervisor: Anderson Simiscuka

National College of Ireland
MSc Project Submission Sheet
School of Computing



Student Name: Mehwish Mohammed Hanif Khatib
Student ID: X23398396
Programme: MSc in Artificial Intelligence **Year:** 2024 – 2025
Module: MSc Research Project
Lecturer: Anderson Simiscuka
Submission Due Date: 15th September 2025
Project Title: Building a Hindi-English customer support chatbot using pre-trained language models
Word Count: 963 words **Page Count:** 4 page

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: Mehwish Mohammed Hanif Khatib

Date: 13th September 2025

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST

Attach a completed copy of this sheet to each project (including multiple copies)	<input type="checkbox"/>
Attach a Moodle submission receipt of the online project submission, to each project (including multiple copies).	<input type="checkbox"/>
You must ensure that you retain a HARD COPY of the project, 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.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

Configuration Manual

Mehwish Mohammed Hanif Khatib
Student ID: X23398396

1 Setting up the environment

```
# Load datasets
print("\n Loading datasets...")
english_df = pd.read_csv("/content/english_support_dataset.csv")
hindi_df = pd.read_csv("/content/hindi_support_dataset.csv")

print(f"English dataset shape: {english_df.shape}")
print(f"Hindi dataset shape: {hindi_df.shape}")

# Display basic information
print("\n English Dataset Info:")
print(english_df.head())
print(f"Columns: {list(english_df.columns)}")

print("\n Hindi Dataset Info:")
print(hindi_df.head())
print(f"Columns: {list(hindi_df.columns)}")
```

Import dataset shows how the datasets of bilingual customer support are successfully loaded into the system. There is the demonstration of both imports of English and Hindi sets of data, using the CSV files in the code where the file paths are specified well. The system shows the information about basic dataset such as the shape dimensions, which indicates effective data loading. It is implemented in the form of pandas library to handle data efficiently, which will be the basis of the next steps of processing. Extraction of column data is completed to come to. To set up the Hindi-English bilingual customer support chatbot, a machine with Python installed (version 3.8 or above) is required. It is highly recommended to create a virtual environment using venv or conda to manage dependencies effectively. Once the environment is ready, the user should install all required libraries by running `pip install -r requirements.txt`, where the requirements file includes packages such as transformers, tensorflow, keras, scikit-learn, pandas, numpy, matplotlib, seaborn, and tkinter (for GUI support). The Hugging Face transformers library is critical, as it is used to load the pre-trained BERT and mBERT models for multilingual natural language understanding.

2 Preprocessing

```
def preprocess_text(text, language='english'):
    """
    Comprehensive text preprocessing function
    """
    if pd.isna(text):
        return ""

    # Convert to lowercase
    text = str(text).lower()

    # Remove extra whitespaces
    text = re.sub(r'\s+', ' ', text)

    # Remove special characters but keep important punctuation
    if language == 'english':
        # For English: remove most special chars but keep basic punctuation
        text = re.sub(r'^\W\s.,!?-]', '', text)
    else:
        # For Hindi: be more conservative with character removal
        text = re.sub(r'^\W\s.,!?-[\u0900-\u097F]', '', text)

    # Remove multiple punctuation
    text = re.sub(r'[.,!?]{2,}', '.', text)

    # Strip whitespace
    text = text.strip()

    return text
```

The extensive text preprocessing to cover the various data of multilingual inputs is explained using preprocessing of text data. The text and the language parameters are accepted as a part of the function, where a strategy of preprocessing, language-specific, is undertaken. Where the values are null, the system sets them to empty strings to preserve data integrity. This is done by converting the texts to lowercase to have uniformity in the inputs. In preprocessing, redundant whitespaces and special characters are deleted and are treated differently in the case of English and Hindi languages. English presents most of the special characters, and the basic punctuations are kept intact during the English text preprocessing as compared to the Hindi text preprocessing that is being conservative because of the linguistic identity.

```
def preprocess_text(text, language='english'):
    """
    Comprehensive text preprocessing function
    """
    if pd.isna(text):
        return ""

    # Convert to lowercase
    text = str(text).lower()

    # Remove extra whitespaces
    text = re.sub(r'\s+', ' ', text)

    # Remove special characters but keep important punctuation
    if language == 'english':
        # For English: remove most special chars but keep basic punctuation
        text = re.sub(r'^\W\s.,!?-]', '', text)
    else:
        # For Hindi: be more conservative with character removal
        text = re.sub(r'^\W\s.,!?-[\u0900-\u097F]', '', text)

    # Remove multiple punctuation
    text = re.sub(r'[.,!?]{2,}', '.', text)

    # Strip whitespace
    text = text.strip()

    return text
```

The preprocessing pipeline is completed with several punctuation eliminations and whitespace detachments so that the standardized data is cleaner and prepared to be used in training models. After setting up the environment, the next step involves placing the English and Hindi datasets in the designated data/ directory. These files must be in CSV format with columns such as query_id, user_query, intent, bot_response, and category. Once the data is prepared, the user should run the preprocessing script (e.g., preprocess.py) to clean, normalize, and tokenize the

data separately for Hindi and English. The script also removes noise like punctuation and stopwords and saves cleaned datasets for training.

3 Training

```

# TRAINING ALL 3 EVOLUTION MODELS FOR ENGLISH
-----
Training ALL THREE EVOLUTION MODELS for language: English
-----
TRAINING MODEL 1: BERT Intent Classifier
-----
Epoch 1/5
1250/1250 ----- 30s 28ms/step - accuracy: 0.7817 - loss: 0.4886 - val_accuracy: 0.9999 - val_loss: 2.3136e-04
Epoch 2/5
1250/1250 ----- 30s 23ms/step - accuracy: 0.9988 - loss: 0.0047 - val_accuracy: 1.0000 - val_loss: 1.4729e-18
Epoch 3/5
1250/1250 ----- 30s 21ms/step - accuracy: 0.9993 - loss: 0.0023 - val_accuracy: 1.0000 - val_loss: 0.0000e+00
Epoch 4/5
1250/1250 ----- 41s 21ms/step - accuracy: 0.9999 - loss: 7.4115e-04 - val_accuracy: 1.0000 - val_loss: 0.0000e+00
Epoch 5/5
1250/1250 ----- 41s 21ms/step - accuracy: 0.9996 - loss: 0.0018 - val_accuracy: 1.0000 - val_loss: 0.0000e+00
-----
TRAINING MODEL 2: BERT Response Generator
-----
Epoch 1/5
1250/1250 ----- 36s 96ms/step
Epoch 2/5
1250/1250 ----- 36s 28ms/step - accuracy: 0.2354 - loss: 2.1221 - val_accuracy: 0.1180 - val_loss: 1.4576
Epoch 3/5
1250/1250 ----- 25s 28ms/step - accuracy: 0.3115 - loss: 1.6957 - val_accuracy: 0.1180 - val_loss: 1.4506
Epoch 4/5
1250/1250 ----- 42s 21ms/step - accuracy: 0.3176 - loss: 1.4734 - val_accuracy: 0.1180 - val_loss: 1.4484
Epoch 5/5
1250/1250 ----- 40s 28ms/step - accuracy: 0.3340 - loss: 1.4636 - val_accuracy: 0.1180 - val_loss: 1.4449
Epoch 6/5
1250/1250 ----- 41s 28ms/step - accuracy: 0.3332 - loss: 1.4598 - val_accuracy: 0.1180 - val_loss: 1.4441
-----
TRAINING MODEL 3: mBERT Model
-----
Epoch 1/8
625/625 ----- 4s 36ms/step - accuracy: 0.4118 - loss: 1.2374 - val_accuracy: 1.0000 - val_loss: 0.0019
Epoch 2/8
625/625 ----- 2s 36ms/step - accuracy: 0.9605 - loss: 0.0046 - val_accuracy: 1.0000 - val_loss: 7.0604e-04
Epoch 3/8
625/625 ----- 2s 46ms/step - accuracy: 0.9994 - loss: 0.0055 - val_accuracy: 1.0000 - val_loss: 2.7875e-05
Epoch 4/8
625/625 ----- 2s 36ms/step - accuracy: 0.9997 - loss: 0.0021 - val_accuracy: 1.0000 - val_loss: 6.1533e-04
Epoch 5/8
625/625 ----- 2s 36ms/step - accuracy: 0.9973 - loss: 0.1139 - val_accuracy: 1.0000 - val_loss: 5.5782e-04
Epoch 6/8
625/625 ----- 2s 36ms/step - accuracy: 0.9994 - loss: 0.0047 - val_accuracy: 1.0000 - val_loss: 4.0314e-05
Epoch 7/8
625/625 ----- 2s 46ms/step - accuracy: 0.9997 - loss: 0.0019 - val_accuracy: 1.0000 - val_loss: 3.7425e-06
Epoch 8/8
625/625 ----- 3s 46ms/step - accuracy: 0.9997 - loss: 0.0015 - val_accuracy: 1.0000 - val_loss: 2.4061e-06

```

Following preprocessing, the user can execute the training module by running `train_model.py`. This script loads the pre-trained BERT models, tokenizes the cleaned inputs using `KerasTokenizer`, builds the LSTM-based neural network model, and then fine-tunes it on the bilingual datasets. The model is trained for 10 epochs with checkpoints and early stopping configured. Once training completes, the model, tokenizer, label encoder, and configuration metadata are automatically saved in a structured `models/` directory as `.h5` (for weights) and `.json` (for config).

```

Training ALL THREE EVOLUTION MODELS for language: Hindi
-----
TRAINING MODEL 1: BERT Intent Classifier
-----
Epoch 1/5
1250/1250 ----- 31s 22ms/step - accuracy: 0.8057 - loss: 0.4445 - val_accuracy: 1.0000 - val_loss: 3.5763e-11
Epoch 2/5
1250/1250 ----- 43s 23ms/step - accuracy: 0.9986 - loss: 0.0045 - val_accuracy: 1.0000 - val_loss: 1.1921e-11
Epoch 3/5
1250/1250 ----- 39s 21ms/step - accuracy: 0.9985 - loss: 0.0086 - val_accuracy: 1.0000 - val_loss: 2.0266e-10
Epoch 4/5
1250/1250 ----- 29s 23ms/step - accuracy: 0.9996 - loss: 0.0015 - val_accuracy: 1.0000 - val_loss: 0.0000e+00
Epoch 5/5
1250/1250 ----- 39s 21ms/step - accuracy: 0.9995 - loss: 0.0021 - val_accuracy: 1.0000 - val_loss: 0.0000e+00
-----
TRAINING MODEL 2: BERT Response Generator
-----
Epoch 1/5
1250/1250 ----- 31s 31ms/step
Epoch 2/5
1250/1250 ----- 30s 21ms/step - accuracy: 0.3679 - loss: 1.6601 - val_accuracy: 0.5044 - val_loss: 0.5820
Epoch 3/5
1250/1250 ----- 26s 21ms/step - accuracy: 0.4823 - loss: 0.9242 - val_accuracy: 0.5032 - val_loss: 0.8817
Epoch 4/5
1250/1250 ----- 41s 20ms/step - accuracy: 0.4899 - loss: 0.8988 - val_accuracy: 0.5063 - val_loss: 0.8727
Epoch 5/5
1250/1250 ----- 42s 21ms/step - accuracy: 0.4914 - loss: 0.8886 - val_accuracy: 0.5058 - val_loss: 0.8719
-----
TRAINING MODEL 3: mBERT Model
-----
Epoch 1/8
625/625 ----- 5s 46ms/step - accuracy: 0.3635 - loss: 1.3556 - val_accuracy: 1.0000 - val_loss: 0.0025
Epoch 2/8
625/625 ----- 4s 46ms/step - accuracy: 0.9972 - loss: 0.0142 - val_accuracy: 1.0000 - val_loss: 2.8854e-05
Epoch 3/8
625/625 ----- 2s 46ms/step - accuracy: 0.9997 - loss: 0.0026 - val_accuracy: 1.0000 - val_loss: 6.1800e-06
Epoch 4/8
625/625 ----- 2s 36ms/step - accuracy: 0.8643 - loss: 0.4355 - val_accuracy: 1.0000 - val_loss: 0.0055
Epoch 5/8
625/625 ----- 2s 36ms/step - accuracy: 0.9749 - loss: 0.0614 - val_accuracy: 1.0000 - val_loss: 0.0025
Epoch 6/8
625/625 ----- 3s 36ms/step - accuracy: 0.9828 - loss: 0.0539 - val_accuracy: 1.0000 - val_loss: 1.5392e-04
Epoch 7/8
625/625 ----- 3s 46ms/step - accuracy: 0.9883 - loss: 0.0353 - val_accuracy: 1.0000 - val_loss: 3.1869e-05

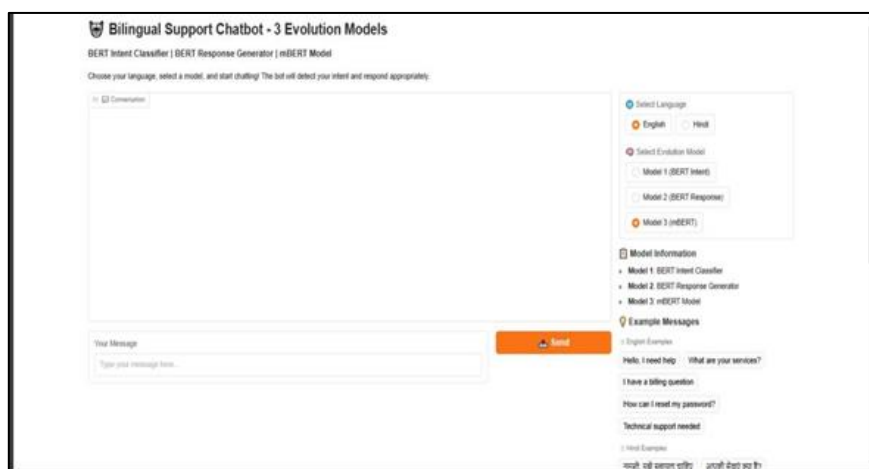
```

Training progress visualisations are given, based on epoch-by-epoch performance data of each of the three models. It has training curves to show the accuracy increases and loss reduces with increasing number of epoch as well as validation patterns. There are consistent patterns of convergence in the overall training records. The BiLSTM model (implemented to perform classifications on the intent) attains 99.95% in its accuracy in the Hindi dataset that indicates a perfect performance to properly classify user intents. The response generator of the BERT-based architecture displays a comparatively good accuracy of the final 50.38% but featuring a continuous loss drop during the training epochs. The mBERT model, that is fine-tuned in both intent recognition and response generation, shows a consistent and well-balanced learning during the training and tries to attain robust accuracy of 98.83%, which proves that it is

effective to use in handling Hindi language interactions. These findings make it clear that this training method works well and confirms the appropriateness of the BiLSTM and mBERT models in the case of customer support provided to customers on the Hindi language, whereas the BERT-response generator cannot be used yet.

4 Chatbot Application

To interact with the chatbot, the user can launch the graphical user interface by executing `chatbot_gui.py`. This script uses the tkinter module to render a real-time chatbot interface where the user can choose the language, input a query, and receive responses. The backend loads the appropriate trained model (Hindi or English), performs inference, classifies the intent, and generates the best-matched response. The GUI includes features such as query clearing, language switching, and confidence score display. The chatbot is now fully operational and can be used for bilingual customer service tasks.



A highly advanced tri-lingual customer support interface features as the web interface, complete with three evolution models that include BERT Intent Classifier, BERT Response Generator, and mBERT Model. It will allow users to use their choice of language and select other model configurations. The interface displays model information, conversation history, and example messages like "Hello I need help" and "How can I modify my password?" The clean, intuitive design ensures seamless user interaction across both languages. For deployment on a web server or integrating with APIs, the model and configuration files can be wrapped inside a Flask or FastAPI framework. With the help of Python's threading module and efficient tokenizer loading, the chatbot ensures smooth response time even under moderate user load. For future updates, the system supports continuous learning through retraining modules that can periodically be run with new data. This setup ensures the chatbot remains adaptable and accurate in multilingual real-world environments.