

# Medical Q&A Chatbot Using Machine Learning

## 1. Introduction

The rapid expansion of healthcare information on the internet has made it challenging for individuals to access reliable medical knowledge. The **Medical Q&A Chatbot** aims to provide users with **accurate, relevant, and fast** responses to their medical queries. The chatbot is built using **Streamlit, spaCy, TF-IDF vectorization, cosine similarity, and Deep Translator** to support multilingual interactions. It leverages the **MedQuAD dataset**, a collection of medical question-answer pairs, to retrieve the most appropriate responses based on user input.

## 2. Objectives

- Develop an **interactive chatbot** for answering medical-related questions.
- Utilize **NLP techniques** (such as entity recognition and text vectorization) for question processing.
- Implement **TF-IDF-based retrieval** for fetching relevant answers.
- Support **multiple languages** using language detection and translation.
- Provide a **user-friendly interface** using **Streamlit**.

## 3. System Architecture

### 3.1 Components

1. **User Interface (UI)**: Built with **Streamlit**, enabling seamless user interaction.
2. **Data Processing Module**:
  - **Parsing MedQuAD XML Data**
  - **Preprocessing text using spaCy**
3. **Retrieval Mechanism**:
  - **TF-IDF Vectorization** for encoding questions.
  - **Cosine Similarity** to find the closest match.
4. **Language Support**:
  - **LangDetect** for automatic language identification.
  - **Deep Translator** for question and answer translation.
5. **Session Management**:
  - **Maintains conversation history** using Streamlit's session state.

## 4. Methodology

### 4.1 Data Collection

The **MedQuAD dataset** contains medical question-answer pairs extracted from trusted sources like **NIH, NCI, and Genetics Home Reference**. These XML files were parsed and converted into a structured **CSV format**.

### 4.2 Text Preprocessing

- **Tokenization & Lemmatization:** Using **spaCy’s en\_core\_web\_sm** model.
- **Stopword Removal:** Eliminates unimportant words to enhance accuracy.
- **Lowercasing:** Ensures uniform text comparison.

### 4.3 TF-IDF Model for Answer Retrieval

- **Vectorization:** Converts questions into numerical representations.
- **Cosine Similarity:** Measures similarity between user queries and dataset questions.
- **Best Match Selection:** The highest similarity score determines the best response.

### 4.4 Multilingual Support

- **LangDetect** identifies the input language.
- **Deep Translator** translates questions to English for processing.
- **Answers are translated back** to the original language before display.

## 5. Implementation

### 5.1 Technologies Used

Component	Technology
Programming Language	Python
Frontend Framework	Streamlit
NLP Library	spaCy
Machine Learning	Scikit-learn (TF-IDF)
Language Detection	LangDetect
Translation API	Deep Translator
Dataset	MedQuAD (XML)

## 5.2 Installation & Setup

To install all required libraries, create a requirements.txt file consisting:

```
streamlit==1.31.1
pandas==2.2.0
spacy==3.7.2
scikit-learn==1.4.0
langdetect==1.0.9
deep-translator==1.11.4
```

**Install using:** pip install -r requirements.txt

**Download spaCy model:** python -m spacy download en\_core\_web\_sm

## 6. Results & Discussion

### 6.1 Accuracy & Performance

- The chatbot provides **highly relevant answers** with an **efficient retrieval mechanism**.
- Multilingual support enhances usability across different user demographics.
- Entity recognition helps identify symptoms, diseases, and treatments.

### 6.2 Limitations

- The bot lacks **deep contextual understanding** and relies on pre-existing questions.
- Responses are **limited to the dataset**; no new knowledge is learned.
- **Translation errors** may occasionally affect response accuracy.

---

## 7. Conclusion & Future Scope

### 7.1 Conclusion

The Medical Q&A Chatbot successfully **retrieves relevant medical answers**, supports multiple languages, and provides an **interactive user experience**. By leveraging **machine learning and NLP**, it simplifies access to healthcare information.

## 7.2 Future Enhancements

- **Fine-tune a transformer-based model (e.g., BERT, BioBERT)** for improved contextual responses.
- **Expand the dataset** to include more diverse medical topics.
- **Enable voice-based interaction** for accessibility.
- **Integrate with a medical knowledge graph** for more precise answers.

## 8. References

1. MedQuAD Dataset: <https://github.com/abachaa/MedQuAD>
2. spaCy NLP Documentation: <https://spacy.io/>
3. Streamlit Framework: <https://streamlit.io/>
4. TF-IDF & Cosine Similarity: <https://scikit-learn.org/>