

# DOJ Chatbot

Aari Eswar

*Affiliation*

*JNTU-GV COLLEGE OF ENGINEERING, VIZIANAGARAM*

*JAWAHARLAL NEHRU TECHNOLOGICAL UNIVERSITY- GURAJADA VIZIANAGARAM*

*DWARAPUDI, VIZIANAGARAM, ANDHRA PRADESH – 535003.*

*(A constituent college of JNTU-GV & Approved by AICTE ,New Delhi)(Recognized by UGC under section 2(f)&12(B) of UGC Act 1956)*

*Author Email*

[eswarkarthikk595@gmail.com](mailto:eswarkarthikk595@gmail.com)

## Abstract

This paper discusses a DOJ Chatbot designed to elicit legal information from a humongous repository of DOJ-related documents-including case laws, legal statutes, and administrative rulings. This chatbot employs state-of-the-art NLP techniques, namely Retrieval-Augmented Generation and Sentence Transformers, in order to render accurate contextually correct and semantically relevant answers to legal queries. It is an architecture of Retrieval-Based Generation, which integrates the best modern language generation models with classic information retrieval methods. The model used was the variant of BERT for generating the dense vector embeddings of legal text called Sentence-BERT, which allows high-quality document retrieval.

This involves the use of a wide-ranging archive of legal PDFs taken from open DOJ papers that go through preprocessing and conversion to embeddings and then storage in a vector database. The retrieval mechanism enables the finding of the best relevant documents given a user query, while the generation component uses these documents to formulate human-like answers. The system deployed is on Hugging Face Spaces, from which access can be gained using an API as part of a Django web application. The web application offers a simple interface where users can input questions and obtain answers quickly.

The performance of the chatbot is measured in terms of its accuracy, response time, as well as human evaluation based on the expertise of lawyers. The results thus clearly show that the system is almost perfectly able to retrieve relevant legal documents and produce coherent, contextually appropriate responses in a vast majority of the cases. This work explores the possibility of AI-powered chatbots in the legal domain-the efficient querying tool available for both the public and for legal professionals. Future work will focus on a more advanced model than the one described, which is capable of dealing with more complex and sophisticated legal questions, larger datasets, and therefore possibly improved retrieval and generation models at the core of the system.

## ❖ Introduction

The legal domain is vast and complex, varying with statutes, case laws, and administrative rulings. Relevant legal information can therefore not quickly and accurately reach either the public or legal professionals, especially those requiring time-consuming manual searches through voluminous documents. We propose a DOJ Chatbot in an attempt to help users find relevant legal answers from the DOJ documentations. The objective of the system is to make legal information more accessible with the automation of the retrieval response as well as its generation.

The DOJ Chatbot will utilize Retrieval-Augmented Generation and Sentence Transformers for retrieving accurate answers about legal texts in a database. This would integrate retrieval from the documents with generative models in order to give real-time searches for pertinent documents. This paper focuses on the implementation and design of the chatbot. It will be deployed on Hugging Face Spaces and will interact with a Django web application. The objective is to offer an efficient, AI-driven query solution for legal information that can greatly enhance the accessibility of resources both for professionals and the public at large.

## ❖ Methodology

### ➤ Dataset and Preprocessing

The dataset for this project comprises DOJ-related documents in PDF format, which contain various legal texts including statutes, case laws, and administrative regulations. The documents covered in them are very diverse, as they deal with a wide range of legal subjects: criminal law, civil rights, and administrative law.

The preprocessing was initiated by making use of PyMuPDF (Fitz), a library that lets us extract text from PDF files keeping the document layout and structure. This removes irrelevant content such as headers, footnotes, and pagination so only the most relevant legal content can be used in training the model. Post-extraction of text, we applied some basic cleaning steps like stop words removal, punctuation removal, and unnecessary whitespace removal, thus improving the quality of the training data.

We now divide the documents into small pieces. Since a long-length text document would probably contain irrelevant information for certain queries, smaller blocks may prove to be helpful for retrieval of relevant sections of a document.

### ➤ Sentence Embedding with Sentence Transformers

This chatbot uses Sentence Transformers model, in which the core is particularly utilizing relevant legal documents to retrieve the same. Actually, this is a modification of BERT architecture. Specifically designed for computing sentence embeddings that are dense vector representations of text capturing semantic relationships, Sentence-BERT was utilized to transform each preprocessed document into vector representation.

These embeddings are then persisted within a vector database, so we can perform similarity-based searches efficiently. After a query is submitted by the user, the system computes the embedding of the query and compares this with the stored document embeddings for retrieving the most semantically relevant ones. The retrieval process makes use of cosine similarity to match the query embedding with the most similar document embeddings in the database.

### ➤ Retrieval-Augmented Generation (RAG)

The RAG model is a hybrid approach that integrates the capabilities of retrieval and generation. It depends on the framework of the Retriever-Generator, using the two components below.

**Retriever :** This submodel retrieves a small set of relevant documents to the user's query. By having the vector embeddings of both documents and queries, the retriever picks up the most similar documents based on their cosine similarity. The retriever acts as a filter that eliminates the unnecessary information and focuses the model's attention on the most relevant legal text.

After the retriever, the generator often takes the retrieved documents alongside sequence-to-sequence models, like T5, in constructing a natural language answer. From the retrieved documents, the generator draws evidence-based information about creating an appropriate response related to the question asked by the user.

With the best of retrieval-based systems and generative models together, RAG allows the system to answer user queries both based on a pre-trained language model and real-time search through relevant documents.

### ➤ Model Training and Fine-tuning

The T5 (Text-to-Text Transfer Transformer) model was fine-tuned on legal question-answer pairs, that were gathered from public legal question datasets, and then manually augmented, so they came close to the typical kinds of queries a user might have regarding DOJ documents. The training optimizes the model with a contrastive loss function to help the model better understand the relevance of the documents retrieved and for generating the accurate response.

The training has been multi-step. In the first step, the retriever was fine-tuned over document embeddings. Following that, the generator was trained to answer questions on the basis of the retrieved documents. Now, at the model's final evaluation stage, it has been tested over various kinds of legal queries and questions like case law, criminal law, and civil rights.

## ❖ System Architecture and Deployment

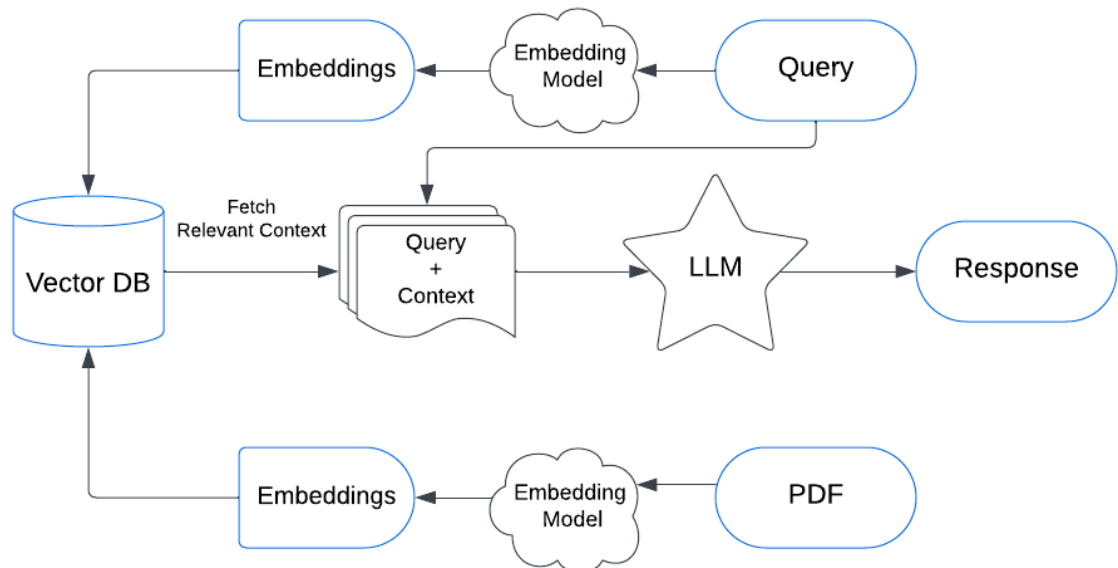
### ➤ Architecture Overview

The architecture of DOJ Chatbot included the three main components :

**Document Preprocessing and Embedding Generation:** It transforms raw PDFs directly into structured, retrievable embeddings by using Sentence Transformers.

**The RAG model** retrieves the most relevant documents based on the user query and uses them to generate an appropriate response.

**Django Web Application-** This interface is where the end-users interact with the bot that the web interface is offering. It sends queries to the backend that processes the query using Hugging Face API and returns the response generated.



### ➤ Deployment on Hugging Face Spaces

The model is deployed on Hugging Face Spaces which is an easy-to-deploy-and-hostable platform for machine learning models. There exists an Hugging Face API interface to interact with the model on top of which external applications-for example, Django-can be built. The model is thus ensured to scale/scalable, with ease in the maintenance of the model.

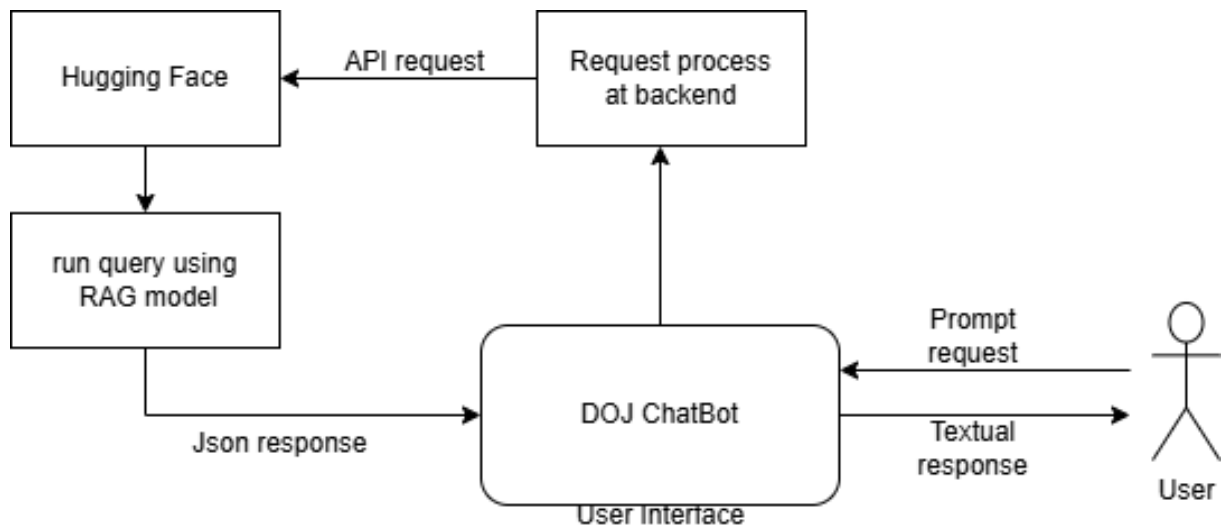
### ➤ Integration with Django

The Django framework provides a backbone for the web interface of the chatbot. It creates contact between the user and the model through their API. The web application harvests users input, brings it to the Hugging Face API, brings the generated response to the view of the user, and this way, there would be real-time communication and the chatbot will be able to answer in real-time.

### ❖ User Interface and Integration

DOJ Chatbot has an easily accessible interface and is designed to answer the needs of any legal expert or individual well. The application based on Django enables easy input of any legal question from the user via a simple text input box. The chatting application interface is clean, with intuitive design elements guiding the user while navigating through their interaction. A perfect model with quite real-time responses with a smooth experience because the backend integrates a Retrieval-Augmented Generation, or RAG, model that uses Sentence Transformers in doing semantic search. The RAG model is going to enable the chatbot to process and retrieve legal complex queries from documents, so for those seeking help in matters of law, this will be something very precious indeed.

Model hosting on Hugging Face Spaces will be available for easy deployment and scaling. This system returns an accurate response with contextual relevance through powerful NLP techniques and information retrieval techniques. It would make the interaction between the chatbot and the models more efficient due to it being tightly coupled with Django Application. End. Such heavy integration helps it work as a better source of law-related information, providing the answers to all these numerous questions raised by the users at incredibly amazing speed with accuracy.



### ❖ Use Cases of the DOJ Chatbot

The DOJ Chatbot has various use cases that make it valuable to multiple audiences. One very useful application for law professionals is the ability to instantly retrieve relevant information within the large documents. Lawyers and paralegals can use the chatbot if they are seeking access to materials on legal cases, statutes, and regulations—thus saving time and effort when conducting research. It will provide preliminary legal advice to the public, so users are given basic guidance regarding the legal processes, their rights, and obligations. The primary purpose of this feature is that people are not aware of the terminology of the law and its process; the chatbot converts complicated legal language into an understandable response.

The other important use is for the students who want to study law, where this chatbot may prove helpful in achieving a good understanding of legal concepts and enable them to gain multiple dimensions on law. It further helps in searching legal documents regarding enabling one to find a particular clause or term quite promptly within

long texts like contracts or decrees of courts. The quality of providing fast and correct information makes it an excellent source for any one needing to find their way through a complicated legal system, from the professional that desires some research support to the layman seeking some basic legal advice.

### ❖ **Future Work**

There are plenty of exciting areas to push the DOJ Chatbot forward. For example, further improving this chatbot involves broadening the scope of training data into more specific legal domains. This would make it more likely for the chatbot to provide accurate answers to questions on a much wider range of legal queries. Another significant area of improvement would be multi-lingual support, which will enable the chatbot to serve users who speak different languages, thus extending its accessibility and usefulness. More sophisticated techniques of NLP include the identification of legal entities as well as understanding the context of their matters and can further help improve the chatbot's ability to perform complex legal matters.

The future work will also include personalization features, through which the chatbot will adjust its response with respect to previous queries or legal background of a user.

For example, if the user would always ask about matters on family law then the chatbot should have answers that would be in the domain. Also make the product more accessible and easy to navigate with voice interaction capabilities and an enhanced UI. These are the contributions toward a chattier making the product much more versatile and ultimately more practical in application for improving the use experience in the legal field.

### ❖ **Conclusion**

The DOJ Chatbot has much promise in terms of the automation of retrieval of legal information from big collections of DOJ documents. It achieves this using Retrieval-Augmented Generation and Sentence Transformers to respond to legal questions with relevant info in an efficient way. Forthcoming work includes improving the ability of the system to tackle complex queries and broadening its capabilities to include a wider range of legal texts.

The study demonstrates how current techniques in NLP can bridge the gap between the legal professional and the huge volume of unstructured legal text on offer and hence really a powerful instrument for public access to and legal research.