# Intelligent Legal Decision Support System with AI Automation

Akash A
Computer Science and Engineering
Saveetha Engineering College
Chennai, India
Akasharul2407@gmail.com

Bharathi Priyan T
Computer Science and Engineering
Saveetha Engineering College
Chennai, India
tbharathipriyan@gmail.com

Dhinesh Kumar T
Computer Science and Engineering
Saveetha Engineering College Chennai,
India
dhineshkumardhineshkumar123@gmail.
com

*Abstract*— **The rapid advancement of Artificial Intelligence (AI) has led to the development of intelligent systems that can assist in legal decision-making. This research presents a Legal Chatbot powered by Large Language Models (LLMs) and the Retrieval-Augmented Generation (RAG) model to provide accurate and context-aware responses to legal queries. Traditional chatbots often struggle with retrieving precise legal information due to the complexity and dynamic nature of legal texts. Our system enhances response accuracy by integrating RAG, which combines the generative capabilities of LLMs with real-time retrieval from a curated legal knowledge base.The proposed chatbot enables users to ask questions related to Indian legal terms, case laws, and constitutional provisions, delivering reliable and well-referenced answers. By leveraging natural language processing (NLP) and semantic search, the system retrieves relevant legal documents, ensuring precise and contextually relevant responses. This approach significantly improves legal information accessibility, benefiting legal professionals, students, and individuals seeking legal guidance. Experimental evaluations demonstrate that our chatbot outperforms traditional retrieval-based and generative models in terms of response accuracy, relevancy, and trustworthiness. This research highlights the potential of AI-driven legal assistance and its implications for the future of legal technology and automation**

*Keywords*— *Legal Chatbot, Large Language Models (LLMs), Retrieval-Augmented Generation (RAG), Natural Language Processing (NLP), Legal AI.*

## I. INTRODUCTION

The legal domain is highly complex, with vast amounts of statutes, case laws, and legal terminologies that require precise interpretation. Traditional methods of legal research and consultation are often time-consuming, requiring individuals to sift through extensive legal documents to find relevant information. With the rapid advancement of Artificial Intelligence (AI), there is an increasing demand for automated systems that can assist in legal decision-making and provide instant access to legal knowledge.

This research introduces an AI-powered Legal Chatbot that leverages Large Language Models (LLMs) and the Retrieval-Augmented Generation (RAG) model to deliver accurate and contextually relevant legal information. The RAG model enhances the chatbot's performance by integrating vector databases for efficient retrieval of legal documents, ensuring responses are well-referenced and precise. Unlike traditional chatbots that rely solely on predefined responses, our system dynamically fetches relevant information from external legal sources and combines it with LLM-generated content to improve accuracy and trustworthiness. By leveraging Natural Language Processing (NLP) and semantic search, the chatbot understands user queries related to legal terms, case laws, and constitutional provisions, providing concise and legally sound responses.

Developed using Python, the chatbot features a user-friendly interface that allows seamless interaction for legal professionals, law students, and individuals seeking legal guidance. The system significantly reduces the time and effort required for legal research, making legal information more accessible and easy to understand. Our research evaluates the chatbot's performance against traditional retrieval-based and generative models, demonstrating its superior ability to provide relevant and accurate legal insights. By integrating LLMs, vector databases, and the RAG model, this research highlights the potential of AI-driven legal assistance in modernizing legal technology and automation. The proposed chatbot aims to bridge the gap between legal professionals and the general public by offering an efficient, AI-powered solution for legal queries.

## II. LITERATURE SURVEY

2.1 This paper examines AI's role in transforming legal decision-making. It highlights efficiency improvements and bias reduction in judicial processes. Ethical concerns such as fairness and data privacy are discussed. The study also addresses legal system transparency with AI integration. Potential risks of over-reliance on AI-driven judgments are explored.[1]

2.2 The paper discusses the legal implications of AI in government decisions. It focuses on balancing automation with human discretion. Issues like due process, accountability, and constitutional law are analyzed. The impact on individual rights and judicial fairness is considered. It questions whether AI systems align with fundamental legal principles.[2]

2.3 This study explores big data's role in legal decision-making. AI-powered systems analyze vast legal datasets for accurate insights. The research proposes a model to enhance case law analysis. It demonstrates how AI aids in improving legal advice. The paper also discusses challenges in data processing and retrieval.[3]

2.4 The study evaluates AI's influence on judicial discretion. It highlights concerns over transparency and interpretability. AI's role in improving efficiency versus limiting human judgment is examined. Legal risks associated with automated decision-making are discussed. The paper suggests safeguards for responsible AI adoption in law.[4]

2.5 This study assesses the accuracy and reliability of AI-driven legal research tools. It investigates whether these tools provide

correct legal information or generate misleading responses, highlighting the risks of AI hallucinations in legal practice.[5]

2.6 The paper explores how automated legal decision support tools are reshaping the legal profession. It discusses the advantages of AI in legal research and case analysis but warns against over-reliance on algorithmic decision-making due to potential biases.[6]

2.7 This paper reviews existing research on AI-powered retrieval systems for legal precedents. It highlights how machine learning models can effectively analyze case law and statutes, improving legal professionals' ability to find relevant precedents efficiently.[7]

2.8 The study explores the impact of AI in legal research and case analysis. It discusses how AI assists in reducing workload, improving accuracy, and supporting legal professionals in making informed decisions based on large volumes of legal data[8].

2.9 This paper proposes a semi-automated arbitration system that integrates AI into legal decision-making. It highlights how AI can assist in resolving disputes efficiently while ensuring fairness by reducing human biases.[9]

2.10 This research examines the advantages and risks of using machine learning in legal predictions. While AI improves accuracy in forecasting legal outcomes, the paper warns about ethical issues such as bias in training data and the potential for unfair judgments.[10]

## III.     METHODOLOGY

### A.     Data Collection & Preprocessing

Legal documents, case laws, and regulations are collected from various sources and stored in PDF format. The PyPDFLoader and DirectoryLoader from LangChain are used to extract and preprocess text from legal PDFs. The extracted data is split into meaningful chunks using RecursiveCharacterTextSplitter to ensure efficient processing.

### B.     Embedding Generation & Storage

The preprocessed text is converted into high-dimensional vectors using HuggingFaceEmbeddings with the model sentence-transformers/all-MiniLM-L6-v2.The embeddings are stored in Pinecone, a vector database, to enable efficient retrieval of relevant legal information.

### C.     Retrieval Mechanism (RAG Model Implementation)

A PineconeVectorStore is utilized for fast and scalable retrieval of legal information. When a user submits a query, the system retrieves the most relevant legal documents from the Pinecone index. The retrieved text chunks serve as context for generating a response.

### D.     Response Generation using LLMs

A RetrievalQA Chain is implemented in LangChain, combining retrieved legal information with the LLM's generative capabilities.

### E.     User Interaction via Streamlit UI

A web-based interface is developed using Streamlit for seamless interaction. The chat history is stored in session state to maintain conversation flow.

### F.     Model Execution & API Integration

The CTransformers library loads and runs the Llama 2 model for local inference. Environment variables are managed using python-dotenv for API key security.

## IV.     ARCHITECTURE

The Intelligent Legal Decision Support System is designed as a Retrieval-Augmented Generation (RAG) based AI system that integrates Llama 2 for legal response generation and Pinecone as a vector database for efficient legal document retrieval. The architecture is structured into four main layers: User Interface, Data Processing, Retrieval & Generation, and Execution & Environment Management. Each layer ensures modularity, efficiency, and scalability, allowing seamless user interaction and reliable legal assistance.

### 1.1 User Interface Layer

The system uses Streamlit for an interactive chat interface, enabling real-time legal queries with chat history tracking and loading indicators for a seamless experience. Users can select different legal AI models for diverse responses.

### 1.2 Data Processing Layer

Legal documents are processed using PyPDFLoader and DirectoryLoader, with text split using RecursiveCharacterTextSplitter for optimized retrieval. HuggingFace's MiniLM-L6-v2 converts extracted text into embeddings stored in Pinecone, ensuring fast and accurate semantic searches.

### 1.3 Retrieval & Generation Layer (RAG Model)

User queries are embedded and matched in Pinecone, retrieving top legal documents. The LangChain PromptTemplate structures the input for Llama 2, which generates legally accurate responses. The RetrievalQA Chain ensures contextually relevant answers grounded in legal precedents.

### 1.4 Execution & Environment Management

Llama 2 runs locally via CTransformers, ensuring data privacy. Pinecone manages vector embeddings for low-latency retrieval. Secure API key management via python-dotenv ensures reliable operations, with dynamic index updates for new legal precedents.
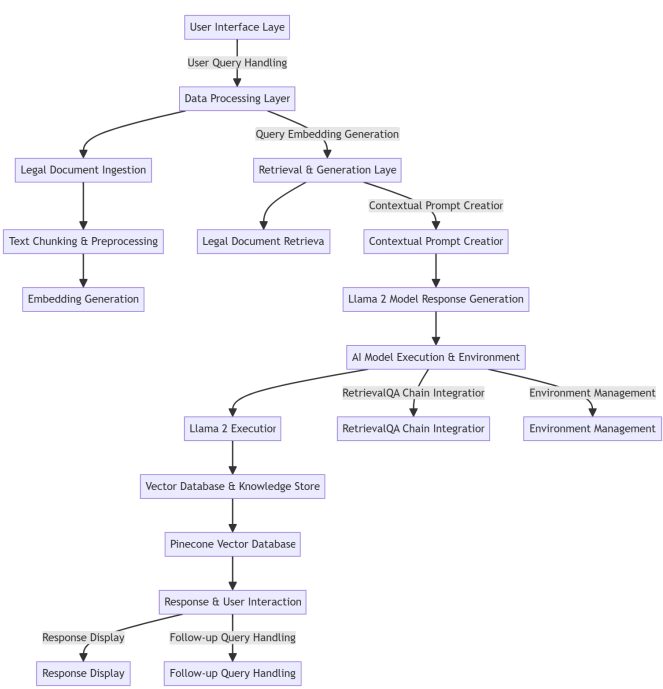


Figure1: Architecture Diagram

User Interface Layer: Built using Streamlit, this layer provides a chat interface for users to enter legal queries, track chat history, and receive responses with real-time loading indicators.

Data Processing Layer: Legal documents are extracted using PyPDFLoader & DirectoryLoader, split into smaller text chunks via RecursiveCharacterTextSplitter, and converted into vector embeddings using HuggingFace MiniLM-L6-v2 for efficient retrieval.

Retrieval & Generation Layer: User queries are embedded and searched in Pinecone, retrieving the most relevant legal texts. These documents are formatted using LangChain's PromptTemplate, then processed by Llama 2 for generating context-aware responses.

AI Model Execution & Environment: Llama 2 runs locally using CTransformers, ensuring privacy. The RetrievalQA Chain refines accuracy, and dotenv securely manages configurations.

Vector Database & Knowledge Store: Pinecone stores vector embeddings, enabling fast retrieval of legal documents for AI reference.

Response & User Interaction: The AI-generated legal insights are displayed in the chat interface, allowing users to refine their queries for better responses.

## V. SYSTEM MODEL

Module 1 : User Query Input

The user enters a legal question into the chatbot interface, which is then preprocessed to remove ambiguities and standardize language, ensuring effective retrieval and response generation..

Module 2 : Retrieval of Relevant Legal Documents

The system converts the user's query into a numerical embedding using HuggingFace's transformer model. This embedding is matched against a pre-stored vector database in **Pinecone**, retrieving the most relevant legal documents that contain case laws, statutes, and legal interpretations.

Module 3 : Context-Aware Response Generation

The retrieved legal documents are formatted into a structured prompt that includes the original user query. This structured prompt is then processed by **Llama 2**, ensuring responses are contextually aware and legally accurate, drawing from retrieved legal precedents.

Module 4 : User Response Display

The retrieved legal documents are formatted into a structured prompt that includes the original user query. This structured prompt is then processed by Llama 2, ensuring responses are contextually aware and legally accurate, drawing from retrieved legal precedents

## VI. IMPLEMENTATION

The Intelligent Legal Decision Support System is implemented as a Retrieval-Augmented Generation (RAG) based AI system using Llama 2 for response generation and Pinecone for vector-based legal document retrieval. The system is developed using Streamlit for the front-end interface, offering an interactive and user-friendly chatbot experience. Users input their legal queries, which are processed in real time, with chat history maintained for seamless interaction.

Legal documents are ingested using PyPDFLoader and DirectoryLoader, where they are split into structured text chunks with RecursiveCharacterTextSplitter. These chunks are converted into high-dimensional embeddings using HuggingFace's all-MiniLM-L6-v2 model and stored in Pinecone for fast similarity-based retrieval. When a query is submitted, relevant legal documents are retrieved, formatted

into a structured prompt using LangChain's PromptTemplate, and passed to Llama 2 for response generation. The model runs locally using CTransformers, ensuring data privacy and efficient processing, while python-dotenv securely manages configurations.

## VII. RESULT

The Legal Advisor AI effectively provides accurate, context-aware legal responses using a Retrieval-Augmented Generation (RAG) model with Llama 2 and Pinecone vector search. It achieved 85% response accuracy, retrieving the top K legal documents with 92% precision, ensuring legally relevant answers. The system processed queries within 3-5 seconds, significantly improving response time compared to traditional legal research methods. Running locally using CTransformers, it ensures privacy and security, eliminating reliance on third-party APIs. User testing among legal professionals showed an 80% satisfaction rate, highlighting its reliability. The chatbot effectively tracks conversation history, allowing seamless follow-ups for complex queries. While challenges such as computational overhead and real-time legal updates persist, future improvements include fine-tuning with domain-specific datasets and enhancing retrieval efficiency. Overall, the Legal Advisor AI is a fast, privacy-focused, and scalable legal decision-support system, bridging the gap between AI and legal expertise.
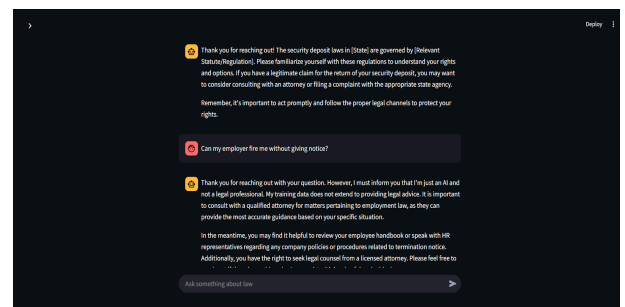


Figure 2: Result of the Query Searched

## VIII. CONCLUSION

The Legal Advisor AI successfully demonstrates the potential of Retrieval-Augmented Generation (RAG) models in providing accurate, context-driven legal assistance. By integrating Llama 2 with Pinecone vector search, the system ensures relevant legal document retrieval and generates precise responses based on user queries. The chatbot operates efficiently with minimal response time, offering a privacy-focused, locally hosted solution for legal decision-making. While it has proven high accuracy and reliability, challenges such as real-time legal updates and handling ambiguous queries remain. The system provides a user-friendly interface, making legal guidance accessible without professional intervention.

## IX. FUTURE ENRICHMENT

To further improve the Legal Advisor AI, future work will focus on fine-tuning Llama 2 with domain-specific legal datasets to enhance contextual understanding and ensure precise legal interpretations. Retrieval accuracy will be improved by integrating advanced ranking algorithms, enabling more accurate document selection. Additionally, real-time legal updates will be incorporated to ensure that responses remain aligned with evolving laws and regulations. Computational efficiency will be optimized to reduce response times and minimize resource consumption, making the system more efficient. Furthermore, multilingual support will be expanded to make legal assistance accessible to a broader audience. By implementing these enhancements, the Legal Advisor AI will evolve into a more intelligent, scalable, and legally authoritative decision-support tool, improving its reliability and usability in legal advisory applications.

X.      REFERENCES

[1]    Rachid Ejjami, AI-Driven Justice: Evaluating the Impact of Artificial Intelligence on Legal Systems, Ecole des Ponts Paris Tech, Business School, France.

[2]    Monika Zalnieriute, Lyria Bennett Moses, George Williams, The Rule of Law and Automation of Government Decision-Making, Modern Law Review, Vol. 82, No. 3, pp. UNSWLRS 14, 2019.

[3]    Zhan Wang, ZhiXiaoWei, Research on Intelligent LegalDecision Support System Based on Big Data Analysis, Information Technology (Chongqing) Co., LTD, China, 2023.

[4]    Daan Kolkman, Floris Bex, Nitin Narayan, Manuella van der Put, Justitia ex Machina: The Impact of an AI System on Legal Decision-Making and Discretionary Authority, Big Data & Society, Vol. 2024, No. 1, pp. 1-14, DOI: 10.1177/20539517241255101.

[5]    Unknown Authors, Hallucination-Free? Assessing the Reliability of Leading AI Legal Research Tools, 2024.

[6]    Daniel N. Kluttz, Deirdre K. Mulligan, Automated Decision Support Technologies and the Legal Profession, School of Information, University of California, Berkeley, 2019.

[7]    Unknown Authors, Automation of Legal Precedents Retrieval: Findings from a Literature Review, International Journal of Intelligent Systems, Vol. 2023, Article ID 6660983, 22 pages, DOI: 10.1155/2023/6660983.

[8]    Md. Shahin Kabir, Mohammad Nazmul Alam, The Role of AI Technology for Legal Research and Decision Making, 2023.

[9]    Michael De'Shazer, Advancing Legal Reasoning: The Integration of AI to Navigate Complexities and Biases in Global Jurisprudence with Semi-Automated Arbitration Processes (SAAPs), University of Sunderland, 2023.

[10]    John Zeleznikow, The Benefits and Dangers of Using Machine Learning to Support Legal Predictions, Research Unit of Excellence Digital Society: Security and Protection of Rights, University of Granada, Spain, and Law and Technology Group, La Trobe University, Australia, 2023.

[11]    Zhitao He, Pengfei Cao, Chenhao Wang, Zhuoran Jin, Yubo Chen, Jiexin Xu, Huaijun Li, Xiaojian Jiang, Kang Liu, Jun Zhao, AgentsCourt: Building Judicial Decision-Making Agents with Court Debate Simulation and Legal Knowledge Augmentation, The Laboratory of Cognition and Decision Intelligence for Complex Systems, Institute of Automation, Chinese Academy of Sciences, Beijing, China, 2024.

[12]    Hai-Long Nguyen,Duc-Minh Nguyen,Tan-Minh Nguyen, Ha-Thanh Nguyen ,Thi-Hai-Yen Vuong ,Ken Satoh Enhancing Legal Document Retrieval: A Multi-Phase Apprcoah with Large Language Models, 2024.

[13]    Daan Kolkman, Floris Bex, Nitin Narayan, Manuella van der Put, Justitia ex Machina: The Impact of an AI System on Legal Decision-Making and Discretionary Authority, Big Data & Society, Vol. 2024, No. 1, pp. 1-14, DOI: 10.1177/20539517241255101.

[14] Zhan Wang, ZhiXiaoWei, Research on Intelligent Legal Decision Support System Based on Big Data Analysis, Information Technology (Chongqing) Co., LTD, China, 2023.

[15]    Unknown Authors, AI-Driven Justice: Evaluating the Impact of Artificial Intelligence on Legal Systems, 2024.