

AI - Powered Legal Documentation Assistant

Vaishnavi C ¹,
Department of Computer
Science and Engineering,
Presidency University,
Bengaluru, India

Shruthi V ²,
Department of Computer
Science and Engineering,
Presidency University,
Bengaluru, India

Ruthika S Shetty ³,
Department of Computer
Science and Engineering,
Presidency University,
Bengaluru, India

Sreelatha PK ⁴
Assistant Professor,
Department of Computer
Science and Engineering,
Presidency University,
Bengaluru, India

Abstract: Legal documentation is a critical yet complex process that often requires expert knowledge, making it inaccessible for individuals and small businesses with limited legal resources. This project aims to develop an **AI-powered Legal Documentation Assistant** that automates the drafting of legal documents in plain language, ensuring clarity and ease of understanding. By leveraging **Natural Language Processing (NLP)** and **Machine Learning (ML)** techniques, the system will generate legally accurate documents based on user inputs while minimizing errors and ambiguities.

The proposed solution is designed to enhance **accessibility, efficiency, and accuracy** in legal documentation, reducing the time and cost associated with legal services. It will allow users to customize documents based on their specific requirements and integrate with legal databases to ensure compliance with existing legal frameworks.

This paper outlines the problem statement, the technology stack, expected outcomes, and the potential impact of the project. Once fully developed, the AI-powered assistant can significantly benefit small businesses and individuals in India, improving access to legal documentation while promoting **legal awareness and empowerment**. Future enhancements may include expanding the range of supported documents and integrating expert legal consultations for complex cases.

Keywords—*AI-powered legal assistant, legal documentation automation, Natural Language Processing (NLP), Machine Learning (ML), legal accessibility, document generation, legal compliance, small business legal support.*

I. INTRODUCTION

Legal documentation is an essential aspect of various business and personal transactions, including contracts, agreements, affidavits, and other legally binding documents. However, the process of drafting these documents can be complex and time-consuming, often requiring specialized legal knowledge. Individuals and small businesses, particularly in India, face significant challenges in accessing legal services due to high costs, lack of expertise, and the complexity of legal language. These challenges can lead to errors, misinterpretations, and potential legal disputes.

With advancements in Artificial Intelligence (AI), particularly Natural Language Processing (NLP) and Machine Learning (ML), there is an opportunity to automate and simplify legal documentation. This research focuses on developing an AI-powered Legal Documentation Assistant that can generate accurate legal documents in plain language, reducing dependency on legal professionals for basic documentation needs. The system

aims to enhance legal accessibility, efficiency, and accuracy by allowing users to input relevant information, after which AI processes and generates customized legal documents that comply with legal standards.

This paper discusses the problem statement, technology stack, expected outcomes, and impact of the proposed solution. The AI-powered assistant is expected to streamline the legal documentation process, minimize errors, and increase access to legal resources for small businesses and individuals. Furthermore, the research highlights potential challenges in implementing such a system, including data privacy concerns, legal compliance, and ethical considerations.

By leveraging AI for legal documentation, this project aims to bridge the gap between legal services and those who need them the most, making legal processes more efficient, affordable, and accessible.

II. LITERATURE REVIEW

The study by Rithik Raj Pandey et al [1] uses a Custom Trained GPT model combined with Optical Character Recognition (OCR) technology to process and simplify legal documents. The AI model employs Natural Language Processing (NLP) and pattern recognition techniques to enhance document readability. The platform includes a chatbot for user interaction, allowing users to draft or simplify legal documents, and even consult legal experts through virtual meetings. The solution uses OCR technology to simplify legal jargon and make document creation user-friendly. The system allows users to upload legal documents for processing or interact with a chatbot for guidance. Users can consult with legal experts directly through the platform, adding significant value to the documentation process.

The system integrates legal databases to keep the generated content updated and relevant. The dataset for training the AI model comes from publicly available legal data.

The study by Imogen Vimala et al [2] utilizes Natural Language Processing (NLP) and machine learning techniques for contract drafting, document retrieval, and legal text summarization. The system features AI-powered chatbots, semantic analysis, and document automation to enhance efficiency. The tech stack includes HTML, CSS, JavaScript (frontend), PHP (server-side), MySQL (database), and CollectChat (AI chatbot development). The system aims to improve accessibility by providing real-time assistance and customizable legal templates. This initiative not only aims at democratizing legal access but also highlights the importance of technological advancements in legal practices by emphasizing user engagement and customization to meet diverse legal needs.

The dataset for training the AI model is derived from legal document templates, legal research databases, and publicly

available legal texts.

The study by Awez Shaikh et al [3] employs Large Language Models (LLMs), Natural Language Processing (NLP), and Machine Learning for legal document drafting, summarization, and query handling. The system includes Optical Character Recognition (OCR) for text extraction from PDFs and integrates a secure vector database for document storage. The tech stack comprises a web-based platform with customizable templates, though specific implementation details are not provided. By leveraging advanced technologies such as natural language processing and machine learning, the platform intends to enhance access to legal resources and empower users to navigate legal matters confidently, contributing to a more inclusive legal system.

The dataset for model training is sourced from legal resources, publicly available legal documents, and external legal databases, ensuring accurate and efficient document generation. The system also offers legal chatbot support and expert consultation options.

The study by G. Kiran Kumar et al [4] employs Natural Language Processing (NLP) and Optical Character Recognition (OCR) to simplify and generate legal documents. The system features a document drafting engine, a simplification tool, and real-time integration with legal databases. It also prioritizes data privacy and security. The methodology includes iterative AI model refinement, usability testing, and user feedback integration to improve document accuracy and accessibility for small businesses and individuals. The project addresses difficulties faced by non-experts in navigating complex legal documentation in India. Real-time integration with legal databases ensures compliance and accuracy in document generation.

The AI models are trained using publicly available legal datasets, contracts, and case laws, ensuring compliance with the latest legal standards.

The study by Lalita Panika et al [5] leverages LangChain, Pinecone, Next.js, Prisma, and MongoDB to build an AI-powered legal documentation platform. The system integrates Natural Language Processing (NLP) for document simplification and generation and uses vector storage (Pinecone) for efficient legal document retrieval. Chatbot functionality powered by OpenAI's GPT models enables conversational interaction with legal documents. The platform also integrates Swagger UI React for API documentation and Kinde Auth for secure authentication. By minimizing errors and democratizing legal services, SimpliLegal stands as a pivotal innovation enabling broader access to justice and legal information.

The dataset for training the AI models comes from legal databases, case laws, and statutes.

The study by Sayash Kapoor et al [6] examines AI's role in legal tasks, focusing on three key areas: information-processing, tasks requiring creativity or judgment, and predictive analytics. However, the paper points out significant issues with these datasets, such as biases, inaccuracies, and data contamination, where training data overlaps with test data, leading to inflated performance estimates. AI models are trained on these datasets to perform tasks like legal information retrieval, case prediction, and document summarization. While generative AI systems like GPT-4 and predictive models such as COMPAS have been applied to legal tasks, the quality of the datasets used remains a critical concern. The paper emphasizes that the lack of clean, unbiased, and comprehensive datasets is a major challenge in effectively evaluating AI in legal settings. Despite these issues, the study suggests that AI could be useful for automating routine

legal tasks but is far from replacing human judgment in more complex legal matters.

The paper discusses datasets commonly used in legal AI applications, which typically include judicial decisions, case law, public legal documents, and legal filings. These datasets are often retrieved from open-access legal databases, court records, and law-specific archives.

The study by Drashti Shah et al [7] explores the use of Artificial Intelligence (AI) and Machine Learning (ML) in legal assistance, specifically for analyzing employment and loan contracts. It employs Retrieval-Augmented Generation (RAG) models, Optical Character Recognition (OCR), and Natural Language Processing (NLP) techniques such as BERT and GPT to extract and interpret legal information. The proposed system allows users to upload legal documents and interact with an AI-powered chatbot for legal guidance, making legal assistance more accessible. However, the research identifies key gaps, including lack of contextual understanding, difficulty in handling diverse document formats, and challenges in semantic inference. The main outcome is a community-based legal advice platform that connects users with legal professionals and provides AI-generated legal insights. Despite its advancements, the system has limitations, such as dependence on OCR accuracy, misinterpretation of legal language, and privacy concerns. It also struggles with adaptability to different legal systems, limiting its global applicability. The research emphasizes the need for better document handling techniques and improved semantic interpretation for more accurate legal AI systems. Overall, the paper contributes to the automation of legal processes but requires further refinement to overcome its challenges.

The dataset used consists of legal documents, including employment contracts, loan agreements, and judicial case records, but the specific retrieval source is not mentioned. These documents are semi-structured and unstructured, requiring text extraction and processing techniques to handle different formats like PDFs, scanned images, and Word files.

The study by Jhanvi Aroraa et al [8] explores AI-driven legal research using Natural Language Processing (NLP) and Information Retrieval techniques. It utilizes BM25, Topic Embeddings (Top2Vec), Law2Vec embeddings, and BERT-based classification to retrieve relevant legal precedents and statutes. The system effectively automates legal precedent retrieval and classifies legal text into rhetorical roles. However, the research identifies key gaps, such as limited context awareness, challenges in processing lengthy documents, and data imbalance in classification tasks. The main outcome of the paper is an AI-based legal research assistant that improves the efficiency of legal document retrieval and ranks among the top 10 submissions at FIRE 2020. Despite its advancements, the system has disadvantages, including BM25's lack of deep contextual understanding, high computational costs of BERT, and inefficiencies in soft cosine similarity calculations. Additionally, topic modeling methods may lose case-specific details, affecting retrieval accuracy. The paper highlights the need for better abstraction techniques and hyper parameter tuning to enhance precision. Overall, the research contributes to automating legal research, but further improvements are required for greater accuracy and efficiency.

The dataset includes 3,260 case documents and 197 statutes, retrieved from the Forum for Information Retrieval Evaluation (FIRE) 2020.

The study by Jinqi Lai et al [9] explores the applications of large

language models (LLMs) in the legal field. It discusses how AI can assist judges, automate legal document generation, and improve efficiency in legal research. The study highlights that legal LLMs are trained on judicial case records, legal statutes, and court decisions, but data accessibility remains a challenge due to privacy concerns. Algorithms such as BERT, GPT, and specialized legal models like ChatLaw and LawGPT are used for text processing and decision-making. However, the paper identifies research gaps, including biased AI outputs, lack of dataset standardization, and limited interpretability of legal decisions. Ethical concerns such as predictive policing and AI-driven judicial decisions potentially undermining human rights are also raised. One major disadvantage of legal LLMs is their tendency to reinforce biases from historical legal data, leading to unfair verdicts. The study also warns that over-reliance on AI could weaken judicial independence, limiting a judge's discretionary power. Additionally, the lack of benchmarking and real-world testing makes it difficult to assess the true effectiveness of these models. While the paper provides recommendations for improving legal AI, it emphasizes the need for transparency, fairness, and better dataset governance to ensure responsible adoption.

The study by Nguyen Ha Thanh [10] introduces LawGPT 1.0, an AI-powered legal assistant fine-tuned on GPT-3 for the legal

domain. LawGPT 1.0 uses the transformer architecture with attention mechanisms and fine-tuning techniques to generate legal documents, answer legal queries, and provide legal advice. Despite its capabilities, the study highlights several limitations, such as the lack of explainability, which raises concerns about trust and accountability in AI-generated legal decisions. Additionally, the model does not support Reinforcement Learning from Human Feedback (RLHF), reducing its ability to refine responses based on user interactions. Ethical and legal concerns regarding privacy, responsibility, and potential bias in AI-generated legal recommendations remain unaddressed. Another major drawback is that LawGPT 1.0 currently supports only English, limiting its applicability in multilingual legal systems. The study suggests future improvements, including expanding language support and integrating better explainability features, but these enhancements have yet to be implemented. The lack of transparency regarding dataset sources and the absence of real-world deployment discussions further weaken its practical reliability. Despite these limitations, LawGPT 1.0 shows potential for improving legal service accessibility, making AI-driven legal assistance available 24/7.

The model is trained on a large corpus of legal text, though the exact dataset source is undisclosed due to a Non-Disclosure Agreement (NDA).

III. METHODOLOGY

I. Overview

The development of the AI-powered Legal Documentation Assistant follows a structured methodology to ensure the accuracy, reliability, and efficiency of the system. The project is implemented using Python-based AI models, specifically focusing on legal document retrieval and analysis. The methodology consists of five key phases: Data Collection, Preprocessing, Model Development, Document Retrieval, and Validation. Each phase plays a crucial role in refining the AI system for handling complex legal texts, extracting key legal terms, and retrieving relevant legal documents based on user queries.

II. Data Collection and Preprocessing

a. Legal Dataset Acquisition

A comprehensive dataset is essential for training and fine-tuning the AI model. The data is gathered from multiple publicly available sources, including government legal repositories, open-source legal documents, law firm databases, and case law collections. The dataset includes different types of legal documents such as contracts, agreements, policies, and legal notices. To ensure data quality and relevance, legal experts review the collected materials, eliminating outdated or jurisdiction-specific content that may reduce the model's generalization capability.

b. Text Cleaning and Formatting

Once collected, the raw legal data undergoes rigorous preprocessing to enhance its quality and usability for AI training. This includes:

- **Removing Special Characters and Formatting Artifacts** – Unnecessary symbols, extra spaces, and formatting errors (e.g., page numbers, footnotes, and metadata) are removed to maintain textual clarity.
- **Tokenization** – The text is split into sentences and words for structured analysis.
- **Stopword Removal** – Common but non-informative words (e.g., “the,” “is,” “an”) are filtered out to focus on meaningful legal content.
- **Lemmatization** – Words are reduced to their root forms to improve consistency across different word variations (e.g., “running” → “run”).

III. AI Model Development

a. Legal Document Embeddings and Vector Search

The system utilizes text embeddings for efficient document retrieval. The methodology includes:

- **Embedding Generation** – Legal texts are converted into numerical vector representations using Hugging Face transformer models.
- **Dimensionality Reduction with PCA** – Principal Component Analysis (PCA) is applied to reduce the vector size while preserving essential semantic information.
- **Efficient Similarity Search with FAISS & HNSW** – The system leverages Hierarchical Navigable Small World (HNSW) graphs and FAISS (Facebook AI Similarity Search) to enable fast and scalable retrieval of similar legal texts.

b. Conversational Retrieval-Based Model

To improve legal information retrieval, the system incorporates:

- **Conversational Retrieval Chains** – The system refines queries iteratively to provide the most relevant legal documents.
- **Legal Context Awareness** – The model understands user queries in context, retrieving relevant legal information effectively.
- **Query Expansion & Re-ranking** – Queries are expanded using related legal terms, and retrieved documents are ranked based on semantic relevance.

IV. Legal Document Retrieval

a. User Query Processing

The system processes user queries to retrieve the most relevant legal documents. Key steps include:

- **Semantic Query Embedding** – Converting user input into an embedding vector to match with stored legal document embeddings.
- **Similarity Matching** – Using FAISS and HNSW for efficient document similarity search.
- **Ranking & Filtering** – Prioritizing the most relevant legal documents and removing irrelevant results.

b. Dynamic Legal Text Retrieval

The retrieved legal texts are:

- **Structured Based on Context** – Ensuring that documents align with user intent.
- **Ranked by Relevance** – Higher-ranked documents contain more pertinent legal information.
- **Adaptive to User Queries** – The system refines search results based on user feedback.

V. Validation and Testing

a. Comparison with Existing Legal Documents

The retrieved outputs are benchmarked against established legal document templates to ensure:

- **Structural Consistency** – Matching real-world legal document frameworks.
- **Completeness** – Ensuring all necessary clauses are present and properly formatted.

b. Performance Metrics

Since the system focuses on retrieval rather than text generation, evaluation is performed using:

- **Recall and Precision** – Measuring how accurately the system retrieves relevant legal documents.
- **MRR (Mean Reciprocal Rank)** – Evaluating how high the relevant document appears in ranked results.
- **Embedding Similarity Scores** – Assessing how close the retrieved legal documents are to the query input.

VI. Ethical Considerations and Future Enhancements

a. Data Privacy and Security

Given the sensitivity of legal information, security measures include:

- **Anonymization of Sensitive Data** – Redacting personal

details to protect user privacy.

- **Encryption Mechanisms** – Securing stored and transmitted legal data.
- **Compliance with Privacy Regulations** – Adhering to GDPR (General Data Protection Regulation) and other data protection laws.

b. Future Enhancements

To expand functionality, the system aims to:

- **Support Multiple Languages** – Enhancing accessibility for non-English legal documents.
- **Integrate Real-Time Legal Consultation** – Providing AI-powered assistance for legal queries.
- **Extend Document Types** – Supporting wills, power-of-attorney documents, and regulatory compliance reports.

This structured methodology ensures that the AI-powered Legal Documentation Assistant efficiently retrieves and analyzes legal texts with high accuracy. By integrating text embeddings, semantic search, and conversational retrieval chains, the system enhances legal documentation processes, benefiting legal professionals and organizations.

IV. RESULTS

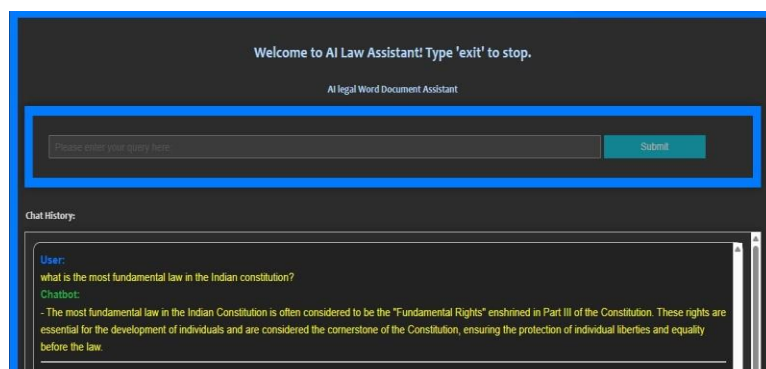


Fig. 1

Fig. 1 displays the output interface of the AI Legal Word Document Assistant, an AI-powered chatbot designed to assist users with legal queries. At the top, the interface welcomes users with a message and provides an instruction to type 'exit' to stop the interaction. The main section includes a query input field (highlighted in blue), where users can enter their legal questions, along with a submit button that processes the input. Below this, the chat history section displays the conversation between the user and the chatbot. In the image shown, the user asked about the most fundamental law in the Indian Constitution, and the chatbot responded by explaining that Fundamental Rights, enshrined in Part III of the Constitution, form its cornerstone.

The response is generated through a document-based retrieval system. When a user enters a query, the system first extracts and preprocesses legal text from PDF files using PyMuPDF (fitz). The text is structured into headings and corresponding content, allowing the system to efficiently search for relevant information. When a query is received, it is processed to identify key terms and matched against the extracted legal text. The system primarily relies on keyword-based search and text-matching techniques to locate relevant sections. The retrieved information is then formatted into a structured response and displayed in the chat history.

V. CONCLUSION

The AI-powered Legal Documentation Assistant leverages state-of-the-art Natural Language Processing (NLP) techniques, including text embeddings, FAISS-based similarity search, and conversational retrieval, to streamline legal document processing. The system effectively retrieves relevant legal clauses from a curated knowledge base, ensuring high accuracy and contextual relevance without generating text from scratch.

By employing sentence embeddings and vector-based similarity matching, the assistant enhances legal document retrieval, enabling efficient question answering and clause extraction. The FAISS indexing technique significantly optimizes search efficiency, allowing real-time legal query responses. The system's performance is evaluated through precision, recall, and relevance scores, ensuring robust and reliable legal assistance.

Furthermore, the integration of legal text preprocessing techniques—including stopword removal, lemmatization, and named entity recognition (NER)—improves the system's ability to understand complex legal terminology. Security and privacy considerations, such as data encryption and anonymization, ensure compliance with legal and ethical standards.

Future enhancements aim to expand the system's capabilities by integrating multi-language support, real-time legal expert validation, and broader legal domain coverage. With its current framework, the AI-powered assistant provides a scalable and efficient solution for legal professionals, businesses, and individuals seeking accurate and structured legal document retrieval.

VI. REFERENCES

- [1]. Rithik Raj Pandey, Sarthak Khandelwal, Satyam Srivastava, Yash Triyar and Mrs. Muquitha Almas, "LegalSeva: AI - Powered Legal Documentation Assistant", International Research Journal of Modernization in Engineering Technology and Science, vol. 06/Issue:03, March 2024.
- [2]. Imogen Vimala, Sreenidhi J. and Nivedha V, "AI - Powered Legal Documentation Assistant", Journal of Artificial Intelligence and Capsule Networks. 6. 210-226. 10.36548/jaicn.2024.2.007.
- [3]. Awez Shaikh, Rizvi Mohd Farhan, Zahid Zakir Hussain and Shaikh Azlaan, "AI - Powered Legal Documentation Assistant", International Journal of Emerging Technologies and Innovative Research (www.jetir.org), ISSN:2349-5162, Vol.11, Issue 4, page no. k526-k530, April-2024.
- [4]. G. Kiran Kumar, A. Shreyan, G. Harini, M. Balaram, (2024), "AI - Powered Legal Documentation Assistant", International Journal of Engineering Innovations and Management Strategies 1 (1):1-13.
- [5]. Lalita Panika, Aastha Gracy, Abhishek Khare, Sanket Mathur and S. Hariharan Reddy, "SimpliLegal: An AI - Powered Legal Document Assistant", International Research Journal of Modernization in Engineering Technology and Science, vol. 06/Issue:04, April 2024.
- [6]. M. E. Kauffman and M. N. Soares, "AI in legal services: New trends in AI-enabled legal services," Service Oriented Computing and Applications, vol. 14, pp. 223–226, Oct. 2020, doi: 10.1007/s11761-020-00305-x.
- [7]. S. Kapoor, P. Henderson, and A. Narayanan, "Promises and pitfalls of artificial intelligence for legal applications," arXiv, Feb. 6, 2024.
- [8]. L. B. Eliot, "AI and Legal Argumentation: Aligning the Autonomous Levels of AI Legal Reasoning," arXiv preprint arXiv:2009.11180, 2020.
- [9]. J. Cui, M. Ning, Z. Li, B. Chen, Y. Yan, H. Li, B. Ling, Y. Tian, and L. Yuan, "Chatlaw: A Multi-Agent Collaborative Legal Assistant with Knowledge Graph Enhanced Mixture-of-Experts Large Language Model," arXiv preprint arXiv:2306.16092, May 2024.
- [10]. Q. Steenhuis, D. Colarusso, and B. Willey, "Weaving Pathways for Justice with GPT: LLM-driven Automated Drafting of Interactive Legal Applications," arXiv preprint arXiv:2312.09198, Dec. 2023.
- [11]. D. Shah, J. Vasi, T. Gandhi, and K. Dabre, "AI & ML Based Legal Assistant," International Research Journal of Engineering and Technology (IRJET), vol. 11, no. 07, pp. 706-708, Jul. 2024.

Fig. 2

In Fig. 2, the user asked about the Fundamental Duty as per the Constitution of India, and the chatbot responded with an explanation based on constitutional provisions.

Fig. 3

In the above image, the user asked: "What are the Directive Principles of State Policy?" The chatbot responded with a detailed explanation, stating that the Directive Principles of State Policy (DPSP) are guidelines for governance in India, enshrined in Part IV of the Indian Constitution. It highlights that while these principles are not justiciable (i.e., not enforceable by courts), they serve as fundamental guidelines to ensure social and economic democracy.

[12]. J. Aroraa, T. Patankara, A. Shaha, and S. Joshia, "Artificial Intelligence as Legal Research Assistant," in Forum for Information Retrieval Evaluation (FIRE), Hyderabad, India, Dec. 2020.

[13]. P. N. Devaraj, R. T. P. V, M. K. R, and A. Gangrade, "Development of a Legal Document AI-Chatbot," School of Computer Science and Engineering, Vellore Institute of Technology, Chennai, India.

[14]. J. Lai, W. Gan, J. Wu, Z. Qi, and P. S. Yu, "Large Language Models in Law: A Survey," arXiv preprint, arXiv:2312.03718, Nov. 2023.

[15]. Nguyen, H. T., "A Brief Report on LawGPT 1.0: A Virtual Legal Assistant Based on GPT-3," arXiv preprint arXiv:2302.05729v2, 2023.