# AI - Powered Legal Documentation Assistant

Vaishnavi C [1], Shruthi V [2], Ruthika S Shetty [3], Sreelatha PK [4]

[1, 2, 3] Department of Computer Science and Engineering, Presidency University, Bengaluru, India

[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 Marcos Eduardo Kauffman and Marcelo Negri Soares [6] explores the transformative role of AI in the legal industry. It discusses various AI applications, including document analysis, legal research, and practice automation, which enhance efficiency and reduce costs. However, the study highlights a major challenge in the legal sector: the lack of structured and accessible legal datasets for training AI models. Public legal data, such as judicial decisions, is often scattered across different systems, making it difficult to retrieve and analyze effectively. Additionally, AI systems currently struggle with abstract reasoning and complex legal decision-making, limiting their effectiveness in nuanced cases. While predictive analytics can forecast case outcomes, biases in datasets can result in unfair or unreliable conclusions. Many law firms resist AI adoption due to business models based on billable hours, which do not incentivize automation. Ethical concerns regarding transparency and the fairness of AI decisions further hinder widespread adoption. AI also raises data privacy and cybersecurity risks, especially in handling

sensitive legal documents. Despite these challenges, AI continues to revolutionize legal services by automating repetitive tasks and improving access to justice. The paper concludes that interdisciplinary research is needed to address these limitations and ensure AI's ethical and effective integration into the legal field.

The paper does not specify a particular dataset but mentions that most law firms are "document-rich but data-poor," with legal data being either unavailable or inconsistent in format.

The study by Sayash Kapoor et al [7] 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 Dr. Lance B. Eliot [8] explores the integration of Artificial Intelligence (AI) in legal argumentation. It introduces the Levels of Autonomy (LoA) of AI Legal Reasoning (AILR), a framework that categorizes AI's role in legal decision-making from basic assistance to full autonomy. AI techniques such as Natural Language Processing (NLP), Machine Learning (ML), Deep Learning (DL), and Knowledge-Based Systems (KBS) are discussed as potential tools for legal reasoning. The paper proposes the CARE Model (Crafting, Assessing, Refining, and Engaging) to describe AI's involvement in legal argumentation. The study highlights a gap in real-world AI applications for legal reasoning, as current systems remain largely theoretical or at prototype stages. Key disadvantages include lack of structured datasets,

interpretability issues, and ethical concerns surrounding AI's role in law. The research emphasizes that AI legal reasoning must be explainable and justifiable to gain acceptance in professional practice. While AI holds promise for enhancing legal analysis, full automation remains a distant goal due to legal complexities and contextual nuances. The paper calls for further research into ethical, regulatory, and societal implications before AI can be widely adopted in legal decision-making.

The study does not use a specific structured dataset but relies on theoretical models, prior legal research, and various academic references as its foundation. Instead of retrieving data from a centralized source, the paper draws from existing legal texts, AI research papers, and conceptual frameworks.

The study by Jiaxi Cui et al [9] introduces an AI-based legal assistant designed to improve the accuracy and reliability of legal consultations. The model employs a Mixture-of-Experts (MoE) framework, integrating knowledge graphs, retrieval-augmented generation (RAG), and multi-agent collaboration to ensure accurate legal reasoning. The system features four specialized agents—Legal Assistant, Legal Researcher, Lawyer, and Legal Editor—which simulate real law firm workflows to provide structured legal services. The study demonstrates that ChatLaw outperforms GPT-4 by 7.73% in accuracy on Lawbench and by 11 points in the Unified Qualification Exam for Legal Professionals, highlighting its superior legal text understanding and reasoning capabilities. Despite its advantages, the paper identifies key research gaps, such as hallucination issues, dataset limitations, and the need for better AI explainability. Major disadvantages include high computational costs, privacy concerns, and limited generalization to legal systems outside China. Additionally, AI bias and interpretability challenges necessitate human verification in legal decision-making. The paper emphasizes that while AI can significantly enhance legal services, full automation remains challenging due to contextual complexities and ethical considerations. Future research should focus on improving dataset diversity, enhancing security, and reducing computational resource demands for practical implementation.

It utilizes a high-quality legal dataset sourced from multiple legal documents, case laws, and legal repositories, enhanced with knowledge graphs and manual refinement by legal experts.

The study by Quinten Steenhuis et al [10] explores the use of generative AI for automating the drafting of interactive legal applications. The study employs

GPT-3 and GPT-4-turbo to generate legal interview questions and assist in form automation. Three approaches are tested: a fully AI-driven method, a constrained template-based approach, and a hybrid model combining AI with human review. The findings suggest that the hybrid model is the most effective, reducing human effort while maintaining accuracy. The paper highlights a research gap in fully automated legal form generation, as AI struggles with complex conditional logic and contextual legal understanding. Key disadvantages include hallucination risks, difficulties in handling diverse legal documents, and limitations in checkbox recognition within PDFs. Additionally, AI-generated forms require significant human review to ensure compliance and usability. The study suggests further improvements in AI-assisted legal automation, particularly in refining question logic and improving PDF field recognition. Overall, the research demonstrates that AI can accelerate legal form automation but cannot replace human oversight in complex legal workflows.

It utilizes legal forms and templates from various court systems and organizations, processed through the Assembly Line Weaver tool.

The study by Drashti Shah et al [11] 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 [12] 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 Pranav Nataraj Devaraj et al [13] presents a chatbot designed to assist with legal document queries. It utilizes LangChain, an NLP framework, along with GPT-based Large Language Models (LLMs) to process and retrieve information from uploaded legal documents and the Indian Constitution. The chatbot uses Cosine Similarity to compare user queries with stored text chunks, while a Flask-based backend provides a REST API for query processing. The outcome of the research is a functional Android-based chatbot capable of answering legal queries using context-aware retrieval techniques. However, the study identifies gaps, including limited AI training capabilities, restricted query token limits, and scalability challenges. Additionally, the chatbot depends on pre-uploaded documents, lacks a real-time legal database, and struggles with complex legal reasoning beyond keyword matching. The system also faces computational inefficiencies when processing large

documents and potential security risks due to storing sensitive legal texts on a server. Despite these limitations, the research provides a solid foundation for AI-driven legal assistance, with future improvements needed in adaptive learning, document sourcing, and enhanced user experience.

The dataset consists of pre-uploaded legal texts, stored in a backend server, which are broken into vector embeddings for efficient search and retrieval.

The study by Jinqi Lai et al [14] 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 [15] 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).

## 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.

[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.