**A REPORT**
**ON**
# AI- POWERED LEGAL DOCUMENTATION ASSISTANT

*Submitted by,*

**Ms. SAHANA REDDY R   -20211CSD0192**
**Ms. DEEPIKA R          -20211CSD0064**
**Ms. LISHA S            -20211CSD0063**

*Under the guidance of,*

## DR. SRABANA PRAMANIK

*in partial fulfillment for the award of the degree of*

## BACHELOR OF TECHNOLOGY

**IN**

**COMPUTER SCIENCE AND ENGINEERING (DATA SCIENCE)**

**At**



GAIN MORE KNOWLEDGE
REACH GREATER HEIGHTS

## PRESIDENCY UNIVERSITY

## BENGALURU

## MAY 2025

# PRESIDENCY UNIVERSITY

## PRESIDENCY SCHOOL OF COMPUTER SCIENCE AND ENGINEERING

## CERTIFICATE

This is to certify that the Internship/Project report **"AI- Powered Legal Documentation Assistant"** being submitted by "SAHANA REDDY R, DEEPIKA R, LISHA S" bearing roll number "20211CSD00192, 20211CSD0064, 20211CSD0063" in partial fulfillment of the requirement for the award of the degree of Bachelor of Technology in Computer Science and Engineering is a Bonafide work carried out under my supervision.

**Dr. SRABANA PRAMANIK**
Assistant Professor
School of CSE
Presidency University

**Dr. SAIRA BANU ATHAM**
Professor & HoD
School of CSE & IS
Presidency University

**Dr. MYDHILI NAIR**
Associate Dean
School of CSE
Presidency University

**Dr. SAMEERUDDIN KHAN**
Pro -Vc School of Engineering
Dean School of CSE & IS
Presidency University

# PRESIDENCY UNIVERSITY

## PRESIDENCY SCHOOL OF COMPUTER SCIENCE AND ENGINEERING

## DECLARATION

I hereby declare that the work, which is being presented in the report entitled "AI-POWERED LEGAL DOCUMENTATION ASSISTANT" in partial fulfillment for the award of Degree of Bachelor of Technology in Computer Science and Engineering (DATA SCIENCE), is a record of my own investigations carried under the guidance of **Dr. Srabana Pramanik, Assistant Professor**, **Presidency School of Computer Science and Engineering, Presidency University, Bengaluru.**

I have not submitted the matter presented in this report anywhere for the award of any other Degree.

**SAHANA REDDY R -20211CSD0192**
**DEEPIKA R-20211CSD0064**
**LISHA S-20211CSD0063**

# ABSTRACT

This project presents an all-encompassing, hybrid document summarization platform that combines web-based front-end technology with high-level artificial intelligence on the backend to automate the processing, summarization, and analysis of formal and legal documents in PDF and Word formats. The system is meant to help users quickly extract valuable information from lengthy documents, which is especially valuable in legal, academic, and business contexts where information complexity and density can get in the way of rapid decision-making. The front-end application uses an uncluttered, responsive HTML interface with CSS for styling and JavaScript for handling interactivity. Central libraries like PDF.js and Mammoth.js are employed for parsing and extracting text content from documents uploaded by the user straight within the browser, providing the rendering of raw textual data to be quickly client-side. This provides users with instant feedback, a basic summary which presents the first couple of sentences of the document. But for deeper insight and more advanced processing, the backend Python module takes advantage of the power of the OpenAI GPT-4 model using its API to carry out high-level NLP operations. These involve creating formal summaries of legal documents, finding important legal clauses, answering precise user queries based on the document uploaded, and creating completely new legal documents from structured user prompts. The backend also comprises the capability for PDF validation, text extraction using PyPDF2, and user interaction through a command-line interface. Through such an integration of browser tools and robust language models, the system is able to offer a smooth user experience coupled with the flexibility to scale and implement more AI features in the future. The long-term vision is to ease the cognitive load of reading vast documents, streamlining repetitive drafting legal work and enabling non-legal experts to access complex legal content more surely and precisely. By integrating real-time text extraction with AI-powered analysis, this system greatly improves document understanding, accelerates legal research and drafting operations, and can be further applied to enterprise uses like contract review platforms, compliance software, and digital legal assistants. The design also prioritizes user privacy and processing efficiency by processing simple summaries locally and keeping backend processing for complex tasks, thus being resource-efficient and user-friendly. By and large, this hybrid summarization model is an effective, scalable solution to current document processing demands, representing a major leap towards the integration of browser technologies and smart language models for routine business use.

# ACKNOWLEDGEMENTS

# LIST OF TABLES

# LIST OF FIGURES

# TABLE OF CONTENTS

# Chapter 1

# INTRODUCTION

This OCR-based legal document summarization platform simplifies the extraction, summarization, and analysis of lengthy legal documents with a hybrid frontend-backend architecture. It accepts PDF and DOCX inputs, with features such as clause detection, Q&A, and AI-generated document creation. The frontend processes rapid previews and summaries, while the backend, fueled by GPT-4, conducts in-depth NLP operations. It's intended for legal professionals, researchers, and businesses looking for efficiency and convenience. Planned future features are OCR support, offline LLMs, and sophisticated analytics to increase its applicability.

## 1.1 Background and Motivation

In today's digital age, the amount of textual data, especially in legal, administrative, and corporate contexts, has increased manifold. Legal documents like contracts, agreements, case filings, and compliance records tend to be lengthy, wordy, and filled with jargon, and hence challenging and time-consuming to analyze. Historically, legal experts and researchers have been forced to depend on manual processes to review, comprehend, and extract vital information from these documents. This process is not only time-consuming and prone to errors but also entails considerable costs in terms of labor and time.

There has been an increase in interest in automating the comprehension and summarization of digital documents in formats like PDF and DOCX, as well as the development of machine learning and natural language processing (NLP). Through intelligent summarization and information extraction, artificial intelligence (AI), and in particular large language models like OpenAI's GPT-4, offers a chance to make legal documents more approachable and useful.

The impetus for this project arises from the necessity to create an intelligent, accessible, and

effective tool that bridges the gap between raw legal documents and

human understanding. The capability of summarizing from texts, recognizing primary legal phrases, responding to document-specific questions, and producing new legal content according to input parameters is directly relevant to the demands of knowledge workers and lawyers in today's times. The objective of this project is to satisfy these demands with a mixed application architecture utilizing client-side as well as server-side technologies to gain speed, scalability, and user privacy.

## 1.2 Objectives and Scope

The main aim of this system is to mechanize the extraction, summarization, and interpretation of legal documents so that the process becomes quicker, less prone to errors, and more accessible to a larger group of people. The system is developed with the following objectives:

[1]. Text Extraction: To effectively extract readable, structured text from PDF and Word (DOCX) files, which are the most prevalent formats used in legal and corporate communication.

[2]. Summarization: To provide short, top-level summaries of long documents in order to enable users to have a rapid idea of the essentials without necessarily having to read the whole document.

[3]. Clause and Key Point Identification: To pull out particular legal clauses, commitments, and key points that are critical to enforcing the contract, being compliant, and litigating.

[4]. Q&A Interactive: To allow users to pose natural language queries from the content of the files that have been uploaded and obtain AI-crafted replies that are context-sensitive to the content of the document.

[5]. Legal Document Creation: To create formal drafts of legal documents from inputs that the user provides like party names, terms, and conditions, thus facilitating contract drafting and legal documentation work.

The present implementation handles PDF and DOCX types, the architecture is extensible and can easily be modified for other document types (e.g., scanned image

OCR through Tesseract, email, or HTML documents) by making some changes. Some possible future extensions include real-time co-authoring, audit trails to record legal editing, case law database integration, and multi-document summarization.

## 1.3 System Architecture and Design

The system is designed on a hybrid architecture that is optimized for performance, usability, and scalability. It consists of two main components:

Frontend: The application on the client side is implemented using HTML5 for markup, CSS3 for presentation, and JavaScript for dynamic behavior. PDF.js and Mammoth.js libraries are added to pull out raw text from PDFs and Word documents, respectively. This selection allows straightforward document processing (such as previewing content or producing a naïve summary) to be executed entirely within the browser, hence enhancing responsiveness and data privacy when more in-depth analysis is not needed.

The users engage with a neat, simple interface through which they can upload `.docx` or `.pdf` files. When a file is chosen, the text is read from the file using the correct library and displayed in the output field. A simple JavaScript summarization technique then reduces the content to a set number of sentences and presents users with an overview immediately.

Backend: The server-side logic is executed in Python and deals with computationally expensive NLP operations. With the help of PyPDF2, the backend can read PDFs, check file structures, and robustly fetch full-text content. After acquiring text, the backend communicates with the OpenAI GPT-4 API to analyze the content. Dedicated endpoints carry out operations such as:

1.3.1    Summarizing documents

1.3.2    Extracting legal obligations and clauses

1.3.3    Answering user queries based on the document

1.3.4    Creating completely new legal documents from prompts

This design is modular, with independent functions for each significant feature, making it easy to extend or swap out components. For example, swapping GPT-4 with a self-hosted LLM (like LLaMA or Claude) would involve minimal changes to architecture. Further, this design facilitates asynchronous operation, which helps maintain UI responsiveness even for lengthy backend operations.

### 1.4 Technology Stack and Libraries

The system uses a variety of technologies on the frontend and backend to provide a unified user experience:

Frontend Technologies:

1. HTML5 & CSS3: Are the building blocks of the interface layout and design.
2. JavaScript: Applied to manage events like file uploads and user input.
3. PDF.js: An open-source library developed by Mozilla that renders PDFs in the browser and extracts their content as text, page by page.
4. Mammoth.js: Expert at pulling clean, semantically rich text out of DOCX files, excluding unnecessary styling and embedded content.
5. Client-side Summarizer: A basic implementation that splits the text into sentences and returns the first few as a summary.

Backend Technologies:

 (i) Python 3.x: The main programming language for backend logic because of its rich ecosystem and robust AI/ML support.
 (ii) PyPDF2: For strong and flexible PDF parsing. It has support for validation, page-by-page text extraction, and error handling for malformed documents.
(iii) OpenAI API (GPT-4): Drives natural language understanding and generation features. These include summarization, question answering, clause extraction, and document drafting.
(iv) Time Module: Monitors processing time, which is helpful for debugging and performance tuning.

(v)     Error Handling: The code is structured with try-except blocks to catch and handle errors gracefully, improving reliability and user feedback.

## 1.5 Use Cases and Applications

The OCR-based document summarization system is appropriate for many real-world applications where rapid and precise comprehension of intricate documents is necessary. Some of the significant use cases are:

i.    Legal Firms and Practitioners: Lawyers can use the system to quickly scan contracts, agreements, or case files. By extracting key clauses and summaries, legal teams save hours of manual review. This is particularly helpful during contract negotiation, due diligence, or litigation discovery phases.

ii.    Corporate Compliance: Businesses that deal with regulatory documents (like GDPR policy, tax returns, or environmental compliance reports) can easily grasp obligations and comply with them through the summarization and key-point extraction capabilities.

iii.    Academic Research: Students and researchers dealing with long legal or policy documents can leverage the system to extract the vital information to analyze or cite.

iv.    Government and Public Sector: The tool can be applied by public administrators to make sense of extensive policy papers, budget reports, or legislative documents to inform public service delivery.

v.    Document Drafting: For those who are not familiar with the technicalities of legal writing, the document generation module facilitates the production of tailored legal templates through insertion of common terms. Legal drafting is thereby made accessible to entrepreneurs, freelancers, and others who are not necessarily trained in law.

vi.    Legal Education: Educational platforms and law schools can integrate the tool to educate students on document structure, interpretation, and legal writing through engagement with actual document samples.

## 1.6 Challenges and Future Directions

While useful, the system has some issues and areas to be improved. One main limitation is the use of heuristic summarization in the frontend, which fails to capture importance and context, as opposed to transformer-based models. What this implies is that while the summary is fast when it starts, it may not always be significant, particularly for documents with dispersed key points.

Another obstacle is the processing of scanned PDFs or image files that need Optical Character Recognition (OCR). Though the existing system specializes in text-based PDFs and DOCX documents, subsequent releases can incorporate Tesseract or other OCR applications to harvest text from images, hugely expanding the range of usable files.

The backend of the system, though strong, is also internet and API availability dependent for OpenAI. Expenses of using GPT-4 could be exorbitant for some, particularly in high-use enterprise environments. An improvement would be to integrate local LLMs (such as LLaMA, Mistral, or GPT-J) for offline processing and privacy guarantees, but with performance and accuracy trade-offs.

A further direction is the incorporation of document analytics—e.g., term frequency, regular expression-based clause detection, and entity recognition (e.g., parties, dates, monetary amounts). These would provide a structured intelligence layer above free-text summary.

Furthermore, adding user authentication, cloud storage, or support for multi-user collaboration would enable the system to become a cloud-based legal assistant platform for law firms and law schools. Such UI features as drag-and-drop file upload, progress bars, or multi-document comparison would make the system even more usable.

# Chapter 2

# LITERATURE SURVEY

The discipline of OCR-based legal document summarization borrows from a number of established fields, among them optical character recognition, natural language processing, and AI-powered legal analytics. Optical character recognition technology allows scanned legal documents to be translated into machine-readable text, serving as the basis for subsequent processing. This is particularly critical in the legal field, where most documents are stored as PDFs or images. PDF and DOCX parsing libraries and tools are often utilized to retrieve raw text content, with emphasis on maintaining semantic structure over visual layout. At the frontend, browser tools assist in rendering documents and enabling real-time text extraction for immediate user feedback. In summarization, early methods employed extractive approaches that chose prominent sentences based on term frequency or graph-based significance. These techniques were simple but rarely had context-awareness. Advances now are geared towards abstractive summarization, driven by deep learning and large language models, that enable systems to create human-like summaries that are closer to the document's intent and semantics. This is especially useful in the legal field because legal language is technical, involving lengthy sentences, nested clauses, and specialized jargon. AI systems trained on general or legal-specific corpora can now accomplish tasks like summarization, clause extraction, question answering, and even document drafting with greater accuracy. These systems read context, comprehend relationships between legal entities, and produce content that conforms to formal legal writing conventions. Most systems now use a hybrid architecture that brings together frontend and backend elements, where light-weight operations such as document preview and basic summarization take place on the client side, while heavier computation such as natural language generation and clause interpretation happen on the server. This provides both performance and security for data. Furthermore, web technology makes it possible for users to engage with legal documents in an easy manner, uploading content, reading summaries, posing document-related questions, and getting structured answers. The ability to access language models through APIs also provides dynamic integration with external AI functionalities, boosting backend smartness without clogging local resources. Use of such systems finds applications in law firms, businesses,

schools, and government agencies, where deciphering and navigating complex legal documents effectively is the need of the hour. Lawyers are helped by systems that help minimize the time taken to review contracts, compliance, and preparation for litigation. On the other hand, people who are not lawyers are able to comprehend their rights and duties through easy summaries and drafts produced by AI. Even with these developments, problems persist in the processing of scanned images, inconsistent formatting, and the reliability of summaries in sensitive legal applications. Consequently, work continues to refine OCR accuracy, fine-tune models for legal-specific vocabulary, and make systems more robust and private. As a whole, the combination of OCR, NLP, and web technologies presents a solid platform on which to develop useful, smart systems that enhance the accessibility, readability, and utility of legal documents.

**Akshaya Kamalnath [1]** explains that corporate insolvency law deals with a company's inability to pay its creditors. AI-based tools and ODR platforms increase efficiency, transparency, and accessibility in insolvency proceedings. Automation supports dispute resolution, but issues are AI biases, data security, and loss of jobs. Professionals will, however, adjust by upskilling in technology for future insolvency administration.

**Ashwini S et al. [2]** explains that the implementation of AI and ML reduces insolvency law by automating processes, predicting analytics, and generating legal documents. In Colombia, Finland, and the UK, AI-based platforms promote clarity and cost savings. While AI makes it easier to perform the work, they are also the root of issues such as possible data leaks, bias in algorithms, and decreased legal work. As a result, it should be stressed that AI will not totally replace but rather be a supplement to human professional advice.

**J.A. Siani [3]** explains that the use of AI is dramatically changing the justice system via automation of case management, prediction analytics, and generation of legal documents in countries like Colombia, Finland, and the UK. Also, in India, the USA, and China, AI is employed to e-discovery, contract analysis, and legal research. Even though AI strengthens the process, drawbacks such as bias, data privacy, and accountability exist. AI, in the end, will be rather an aid than a for legal professionals that can stimulate ethical implementation globally.

**Chitranjali Negi [4]** explains the implementation of AI in the legal industry enables the sectors in legal research, contract analysis, and decision support systems to gain the advantage of efficiency and accessibility. AI-based tools such as NLP, chatbots, and virtual courts play a key role in the process of case management and legal automation. But still, there are concerns about the privacy of data, the presence of bias, and ethical considerations. In the future, AI will work alongside lawyers, which will need good regulations to manage the balance in between innovation and legal principles.

**P.Vimala Imogen et al. [5]** explores that AI is revolutionizing legal documents with the help of contract analysis, document drafting, and legal research based on NLP and ML algorithms. AI-powered chatbots offer real-time legal counsel, and deep learning models improve legal text understanding. Challenges involving contextual understanding and AI explainability are countered by hybrid AI-human models. Future developments in XAI, blockchain verification, and AI conflict resolution seek to increase legal efficiency and accessibility.

**Yating ZHANG et al. [6]** explores that AI legal assistants utilize NLP, ML, and deep learning for contract examination, document generation, and legal research, achieving efficiency and precision. BERT and GPT transformer models are used for processing legal text, with AI legal chatbots providing real-time counsel. AI transparency and regulatory conformity are threatened by data privacy, accuracy, and jurisdictional agility. AI-blockchain, upgraded AI explainability, and AI-human collaborative decision-making for enhanced legal decisions are future research prospects.

**Nikolaos Aletras et al. [7]** explains that the legal automation, case prediction, and evidence examination by NLP and ML for AI. The case facts considerably affect judicial choices in favor of legal realism from research. The Transformer models of BERT and GPT refine legal text analysis, but drawbacks such as explainability and privacy remain. Next-generation research works on blockchain-enabled legal documents, conflict resolution via AI, and expert legal AI models.

**Drashti Shah et al. [8]** explores ways in which AI and ML have revolutionized legal aid by enhancing document and contract analysis and judicial decision prediction. NLP models are 79% accurate in ECtHR case prediction, and AI-based contract analysis identifies legal risk with 75% accuracy. RAG models improve legal information retrieval, yet there are challenges in legal text interpretation and semantic inference. Future studies target blockchain-based smart contracts, sophisticated legal chatbots, and hybrid human-AI models to improve access, accuracy, and fairness in legal decision-making

**Kiran Kumar et al. [9]** explain how AI and OCR facilitate legal document automation, streamlining legal proceedings for individuals and small businesses. NLP assists in identifying clauses, language translation, and document generation automatically, while AI-enabled tools help legal professionals. Most solutions target law firms, which emphasizes the necessity of easy-to-use, user-friendly platforms. Future developments are aimed at enhancing accessibility, data security, and ethical AI to increase legal tech benefits.

**Rizvi Mohd Farhan et al. [10]** explains that the AI and ML have revolutionized legal technology to enhance research, document automation, and contract analysis. Legal Robot and LawGeex, among other companies, simplify contract review, but much work remains to interpret dense legal language and determine accuracy. Caspedia, with its AI tool, strengthens legal research, whereas Legal Consult helps users find experienced legal counsel. This study targets the integration of AI, OCR, and NLP into a legal assistant for small businesses in India to enhance accessibility, compliance, and efficiency in legal documents.

**P. Aishwarya [11]** describes how AI-based legal document assistants use NLP and machine learning to automate contract drafting, document retrieval, and classification. Studies estimate AI's potential to improve accuracy and efficiency in legal processes as well as resolve compliance and flexibility issues. AI is recognized in studies as having a revolutionary effect in simplifying legal processes and decision-making.

**Natalia Khatniuk [12]** discusses the manner in which Artificial intelligence (AI) is transforming legal services through process simplification such as contract drafting, legal research, and decision-making. Research indicates the use of AI in maximizing efficiency,

accessibility, and precision and in resolving compliance and ethical use issues. Researchers highlight the need for regulatory frameworks to ensure AI integration in line with legal principles. AI-powered tools such as chatbots and expert systems are transforming legal processes and democratizing access to legal aid.

**Farhan Aslam [13]** explains here how Artificial intelligence (AI) transformed chatbot technology with the support of natural language processing, machine learning, and deep learning. AI chatbots enhance customer experience, automate response, and increase efficiency in the healthcare, education, and business sectors. Advancements include voice chatbots and virtual assistants with speech recognition and sentiment analysis. Ethical and privacy concerns remain inherent limitations in the use of AI chatbots.

# Chapter 3

# RESEARCH GAPS OF EXISTING METHODS

## 3.1 Lack of Deep Semantic Understanding

[1]. The majority of systems are based on extractive summarization, choosing sentences based on surface features such as term frequency and not actual semantic meaning.

[2]. Such models. are not able to understand the legal intent or meaning. of phrases, which. is essential in the analysis. of law.

[3]. With. shallower language understanding, important. legal clauses. might be misrepresented. or even missed.

## 3.2 Limited Domain-Specific Adaptability

[1]. General NLP models usually perform poorly on legal texts because they are unfamiliar with legal jargon, phrasing, and format.

[2]. Legal texts commonly incorporate antiquated vocabulary or words of extremely specialized meaning that general models get wrong.

[3]. There are few pretrained models that have been fine-tuned specifically for legal corpora, which makes accuracy in extraction and summarization difficult.

## 3.3 Insufficient Management of Document Diversity and Structure

[1]. Legal documents are in different formats (contracts, affidavits, policies, court decisions) with specific layouts and structure conventions.

[2]. Most systems do not tailor their parsing or summarization approach depending on the type or format of the legal document.

[3]. Subtle document structures such as tables, bullet points, or footnotes are usually

neglected or processed wrongly.

## 3.4 Limited Support for OCR and Scanned Documents

[1]. A large percentage of legal documents are available in the form of scanned images or non-editable PDFs, which need to be processed by good OCR.

[2]. Current solutions either do not support OCR or have OCR engines that perform poorly on noise, stamps, and hand-written notes.

[3]. Low-quality OCR results in erroneous or incomplete text extraction, directly affecting the quality of summarization.

## 3.5 Lack of Real-Time Interactive Features

[1]. Most systems provide static summaries but do not support users to engage dynamically with the document using questions or clause lookups.

[2]. There is little integration of contextual Q&A capabilities that enable a user to retrieve certain information against custom queries.

[3]. Interactivity limitations make the system less useful in actual legal workflows where context understanding is central.

## 3.6 Inadequate Integration with Legal Drafting Workflows

[1]. Few summarization tools facilitate the creation or verification of legal documents, even though this is a routine requirement in legal practice.

[2]. Minimal automation exists to generate new contracts or documents based on extracted conditions, parties, and terms.

[3]. Tools tend to be standalone rather than integrated within larger document lifecycle systems employed within law firms or businesses.

### 3.7 Challenges in Multi-Language and Jurisdictional Coverage

[1]. Legal terminology and systems differ greatly between nations and languages, yet the majority of models are trained on English-biased data.

[2]. Multilingual OCR and summarization capabilities for legal use are not well supported.

[3]. Jurisdiction-specific clause identification or summarization is not commonly supported, restricting global usage.

### 3.8 Limited Evaluation and Benchmarking Metrics

[1]. Current models usually do not have strict, domain-specific evaluation metrics that measure summary accuracy, completeness of clauses, or legal fidelity.

[2]. There is no widely adopted benchmark dataset for legal summarization, and comparisons between systems are challenging.

[3]. Most systems are tested using only generic readability metrics, which do not account for the accuracy needed in legal interpretation.

# Chapter 4
# OBJECTIVES

## 4.1 Legal Text Extraction Automation

[1]. Efficient Legal Document Parsing: The system should be able to automatically extract text content from standard legal formats like PDFs and DOCX files, eliminating the need for reading line-by-line or copying by hand.

[2]. Semantic Text Processing: Instead of extracting raw text only, the extraction needs to retain semantic structure—headings, paragraphs, and clauses—to make the content applicable for summarization or clause identification.

[3]. Client-Side Performance: By performing extraction on the client-side through libraries such as PDF.js and Mammoth.js, users can immediately view results without data ever leaving their local machine for basic previews.

## 4.2 Offer Abstractive Summarization

[1]. Create Brief Legal Summaries: The process should condense lengthy legal text into brief, precise summaries that maintain obligations, roles, and essential clauses.

[2]. AI Application for Context-Aware Summarization: GPT-4 is applied to comprehend and paraphrase intricate language into readable summaries for non-experts without legal jargon overload.

[3]. Customization of Depth of Summarization: Users can ask for summaries of certain lengths or levels of detail (e.g., high-level summaries versus detailed clause-level summaries), providing flexibility in accordance with their requirements.

## 4.3 Extract Key Legal Clauses

[1]. Identification of Essential Provisions: The backend should correctly identify and extract essential clauses like termination conditions, indemnity, governing law, and liability with the help of AI-based semantic understanding.

[2]. Prompt Engineering for Clause Identification: Special prompt structures instruct the model to find and replicate clauses exactly, so critical legal constructs are not left out.

[3]. Highlighting and Labelling Clauses: Extracted clauses must be labeled with user-friendly titles (e.g., "Dispute Resolution Clause") for readability and convenience when reviewing or editing.

## 4.4  Allow Natural Language Question Answering

[1]. Context-Aware Responses: Users must be able to pose questions such as "What is the penalty for terminating early?" and get document-specific contextually accurate answers.

[2]. Implementation of Few-Shot Learning: The GPT-4 model is initialized with instances of legal questions and answers to simulate human-like, informative answers more accurately based on the uploaded document.

[3]. Supports Non-Legal Users: This feature enables business proprietors, students, or paralegals to converse with the system directly in plain English without legal training.

## 4.5 Support Legal Document Drafting

[1]. Template-Based Document Creation: The users should be able to input simple parameters (e.g., party names, jurisdiction, dates) and create complete legal drafts such as NDAs, contracts, or terms through AI.

[2]. Structured Prompting for Draft Accuracy: Structured legal prompts are provided to GPT-4 for guaranteeing that the generated documents adhere to a professional.

[3]. Editable Output for Customization: It should be easy to edit the drafts produced by AI so that refinement can be made based on certain needs or jurisdictional demands.

## 4.6 Maintain User Data Privacy

[1]. Client-Side Processing for General Operations: Operations such as previewing material or summarizing text must be done on the client side to not unnecessarily communicate with the server.

[2]. Minimum Data Transfer: The backend is only invoked for complex operations, and data is never kept for a long period, maintaining both the user expectation and privacy policies.

[3]. Security Document Handling: Particularly in legal settings, the system has to ensure that no document is stored, leaked, or abused while processing.

## 4.7 Provide Cross-Platform Compatibility

[1]. Browser-Based UI: The frontend is developed with plain web technologies (HTML, CSS, JavaScript), making the tool run on any contemporary browser without the need for installation.

[2]. No Vendor Lock-In: Open libraries such as PDF.js and Mammoth.js minimize reliance on proprietary software, allowing the system to be easily portable and maintainable.

[3]. Mobile Responsiveness: The UI must be responsive and usable on desktops, tablets, and smartphones to provide accessibility on the move.

## 4.8 Offer Strong Error Handling and Logging

[1]. Graceful Error Recovery: Each backend operation should involve try-except blocks that trap frequent failures (e.g., invalid PDFs, empty input fields) and return user-friendly messages.

[2]. Operation Logging: Important steps—file uploads, summarization requests, clause extractions—are logged (anonymous) for auditing, debugging, and future optimization.

[3]. User Feedback Mechanisms: The system should offer real-time alerts, warnings, or confirmations when an operation fails or succeeds, increasing trust and usability.

## 4.9 Assure Modularity and Extensibility

[1]. Pluggable NLP Modules: All NLP tasks, including clause extraction, Q&A, and summarization, are modular and may be modified or replaced without compromising other parts.

[2]. Support for Upcoming AI Models: With only minor code modifications, the backend should be built to transition from OpenAI's GPT-4 to substitutes like LLaMA, Claude, or open-source LLMs.

[3]. Expandable Document Types: The architecture should eventually enable OCR integration for scanned documents or formats like emails and HTML, even though it now supports PDFs and DOCX.

## 4.10 Increase Legal Accessibility and Awareness

[1]. Democratizing Legal Comprehension: By streamlining legal document viewing and generation, the system facilitates individuals and small enterprises to achieve legal insights without retaining an attorney.

[2]. Educational Use Cases: Law students may employ the system to review actual legal contracts, recognize clauses, and train in legal drafting through example generation.

[3]. Training Tool for Legal Tech Adoption: Firms may employ the system to expose junior personnel to document processing processes, establishing the groundwork for AI-supported legal practice.

# Chapter 5

# PROPOSED METHODOLOGY

## 5.1 Client-Server Hybrid Architecture:

The solution's hybrid architecture optimizes performance, scalability, and privacy by carefully allocating processing between the client and server:

[1]. Client-Side for Lightweight Tasks: The user's browser is used to carry out tasks including file uploads, simple text extraction, and preview rendering. This guarantees a speedier user experience by removing round-trip delay, particularly for little tasks like reading a document's text.

[2]. Server-Side for Complex NLP Assignments: Using robust APIs like OpenAI's GPT-4, resource-intensive backend operations like abstractive summarization, clause extraction, and question-answering are managed. This division guarantees that only jobs requiring a large amount of compute consume server resources.

[3]. Improved Privacy and Performance: The system lowers the danger of data leaks and improves compliance with privacy policies, which is crucial when handling sensitive legal data, by avoiding needless data transfer to the backend and carrying out typical activities locally.

## 5.2 Using JavaScript Libraries for Frontend Text Extraction:

The frontend uses specialized JavaScript modules to handle page processing in order to deliver instant feedback and responsiveness similar to offline:

[1]. The library PDF.js is used to parse and render PDF files in the browser. It allows for the extraction of each page's plain text content, avoiding layouts or embedded pictures that can obstruct processing.

[2]. Utilizing Mammoth.js with DOCX Files: Mammoth.js is designed for clean semantic extraction, in contrast to conventional DOCX parsers. It avoids using decorative

devices or embedded multimedia that are unnecessary for legal analysis in Favor of concentrating on the actual content (headings, paragraphs, and tables).

[3]. Advantages for User Experience and Trust: Before any content is sent to the server, users can receive a preview and a brief synopsis thanks to this design. Such openness fosters confidence, which is essential in legal settings when maintaining confidentiality is crucial.

## 5.3 NLP Processing in the Backend Using GPT-4 with Python:

The backend uses OpenAI's GPT-4 API and Python's strong ecosystem to take over after the document has been submitted and the whole text is ready:

[1]. Validation and Parsing: The system uses PyPDF2 to validate and interpret PDF text structures. Runtime errors are eliminated via early detection of encrypted or corrupted documents.

[2]. NLP Workflows: GPT-4 uses prompt-engineered workflows for summarizing, clause detection, and legal interpretation after the text has been cleaned.

[3]. Flexible and Scalable: The backend can handle a wide range of legal documents, including contracts, NDAs, and privacy rules, with minimal manual modification thanks to the combination of rule-based parsing and generative AI.

## 5.4 Document Synopsis and Clause Recognition:

Reducing verbosity while preserving original intent is necessary when summarizing legal documents

[1]. Custom Prompt Engineering: To distil legal material while maintaining essential concepts such parties' obligations, risks, and rights, GPT-4 is guided by structured prompts.

[2]. Clause Extraction: The model is guided by prompts to find, extract, and even paraphrase certain areas, such as indemnity, termination provisions, or force majeure, for use comprehension.

[3]. Relevance and Actionability: The produced summaries are not generic; rather, they are intended to give readers legally actionable information so they can make judgments without having to read the entire document.

## 5.5 Interactive Clause Explanation and Question Answering

The system allows natural language querying to make legal papers easier to grasp because they can be overwhelming:

[1]. Awareness of Context Q&A: Users can pose queries such as "Who is responsible for payment?" and "What is the governing law of this agreement?" Then, rather than using generic knowledge, GPT-4 creates an answer depending on the context of the document.

[2]. Few-Shot Prompting: To help the model stay rooted in the source material, a few instances of document-based Q&A are given to it.

[3]. Enhanced Efficiency and Comprehension: This feature saves time and increases accuracy by enabling users to bypass pointless parts and get directly to the issues at hand.

## 5.6 Drafting Legal Documents Using Parameters:

The system helps with creation in addition to analysis:

[1]. Form-Based Input: Using an easy-to-use form interface, users enter basic criteria including the names of parties, jurisdiction, obligations, and effective dates.

[2]. Template-Based Generation: GPT-4 creates complete draft contracts, notices, or agreements in an organized, legally compliant format based on legal prompts.

[3]. For Experts and Non-Experts: This demystifies contract generation and expedites documentation, making it particularly helpful for startups, independent contractors, or anyone without legal experience.

**FIGURE 4.1 Flow of the Proposed Methodology**

### 5.7 Managing Errors, Recording, and Expandability:

A system at the production level needs to be robust and maintainable. This element guarantees resilience:

[1]. Try-Except Error Handling: The Python backend is equipped with extensive error-handling features to detect incorrect inputs, sluggish responses, and corrupted files. Every error is handled politely and with unambiguous user messages.

[2]. Operational logging involves recording important actions, processing durations, and faults. This facilitates performance monitoring, debugging, and the creation of datasets for upcoming analytics or model improvement.

[3]. Modular Design for Growth: The backend is constructed to allow for the addition of new modules without requiring a complete system redesign. Examples of these modules include support for OCR (through Tesseract), translation (for multilingual contracts), and integration with other LLMs (such as Claude or LLaMA).

# Chapter 6

# SYSTEM DESIGN & IMPLEMENTATION

**6.1 Modular Design for Architecture:**

[1]. Frontend and Backend Separation: The system employs an explicit frontend-backend separation to facilitate independent development and upkeep of each aspect.

[2]. Microservice-Based Backend Operations: Each functionality such as summarization, clause identification, or question answering is coded as a separate function or API endpoint to support maintainability and testing.

[3]. Scalable Infrastructure Readiness: The design is such that further features (such as multilingual support or integration of OCR) can be added without breaking the core functionality.

**6.2 Frontend Document Processing:**

[1]. Text Extraction using JavaScript Libraries: PDF.js and Mammoth.js are used on the client-side to parse readable, clean text out of PDF and DOCX files without sending information to the server.

[2]. Client-Side Summarization Preview: A light-weight JavaScript technique previews and extracts the initial sentences of the document and provides users with instant feedback on upload.

[3]. Interactive UI Elements: File upload, question submission, and summary triggers are tied to event listeners through JavaScript, hence smooth and responsive user experience.

### 6.3 Backend Implementation in Python:

[1]. PyPDF2 for PDF Parsing: The backend has a module parsing and checking PDF files using PyPDF2 so that even very badly corrupted files are handled as well as possible.

[2]. Text Handling and Cleanup: Pre-processed text is subjected to simple preprocessing—such as trimming leading and trailing whitespace and character normalization—before being fed into the NLP engine.

[3]. Flask or Fast API Framework: A light-weight Python web framework (Flask or Fast API) handles the endpoints for uploading documents, summarization of text, extraction of clauses, and Q&A processing.

### 6.4 Integration of GPT-4 for NLP Operations:

[1]. API-Based Communication: The backend invokes the OpenAI GPT-4 API to execute legal content through structured prompts for summarization, Q&A, and identification of clauses.

[2]. Prompt Engineering Strategies: Thoroughly designed prompts ensure the model gets the context and produces outputs consistent with legal standards and formats.

[3]. Secure API Key Handling: API keys are securely stored in environment variables or configuration files, with access control to avoid unauthorized use.

### 6.5 Mechanism for Clause Identification:

[1]. Clause Categorization Templates: The prompt identifies pre-defined common legal categories (such as payment terms and dispute settlement) that GPT-4 can directly match and identify within the input.

[2]. Awareness of Context in Extraction Logic: The model recognizes the legal context of the sections and extracts information with the proper purpose to infer while maintaining context.

[3]. Labelled Display in the UI: To help users navigate large documents more efficiently, extracted clauses are displayed in the frontend with labelled sections.

## 6.6 Question & Answer System for Law:

[1]. The Natural Language Query Interface allows users to ask queries in natural language, and the system will translate them and return responses that are specific to each article.

[2]. Few-Shot Prompting Setup: Based on prior training, prompts are created using a small number of instances to help the model determine the appropriate kind of answer structure.

[3]. Real-Time Answer Rendering: After processing the input, the backend instantly renders the response, which is displayed in the frontend beneath the question area.

## 6.7 Legal Document Generator Form-Based Parameter Input:

[1]. To start document generating, users fill out a form with the required legal information, such as party names, dates, and obligations.

[2]. Contract Generation with GPT-4: GPT-4 creates fully formatted legal content from input that can be used as letters, contracts, or non-disclosure agreements.

[3]. Editable Output Interface: The user can make changes to the output document before saving or exporting it by viewing it in a text editor on the frontend.

**6.8 Handling System Errors and Expandability:**

[1]. Large Try-Except Blocks: To prevent crashes and provide understandable error messages, backend operations are contained behind exception-handling procedures.

[2]. Logging and Monitoring: For the sake of debugging, performance analysis, and system auditing, all significant system events—such as file processing and API latency—are recorded.

[3]. Plug-and-Play Extensibility: Future improvements such as offline LLM support, other languages, or OCR utilizing Testract are made possible by the design's modular framework.



**FIGURE 6.1 Diagram of the System Design**

# Chapter 7

# TIMELINE FOR EXECUTION OF PROJECT
# (GANTT CHART)

| ID | Task Name | Start | Finish | JAN | FEB | MAR | APR | MAY |
|----|-----------|-------|--------|-----|-----|-----|-----|-----|
| 1 | Review 0 | 29/1/2025 | 31/1/2025 | ▬ | | | | |
| 2 | Review 1 | 18/2/2025 | 21/2/2025 | | ▬ | | | |
| 3 | Review 2 | 17/3/2025 | 21/3/2025 | | | ▬ | | |
| 4 | Review 3 | 16/4/2025 | 19/4/2025 | | | | ▬ | |
| 5 | FINAL VIVA | 10/5/2025 | 17/5/2025 | | | | | ▬ |

**FIGURE 7.1 Gantt Chart**

# Chapter 8
# OUTCOMES

## 8.1 Increased Legal Document Accessibility:

[1]. Time-Saving Summarization: Users are able to comprehend lengthy legal documents quickly through system-generated automatic summaries, saving time from manual reading.

[2]. Clause Identification: Key clauses such as payment terms, termination clauses, and liabilities are automatically highlighted, enhancing navigability.

[3]. Accessible to Non-Legal Users: The system provides legal content that is comprehensible even to users without formal legal education, increasing accessibility.

## 8.2 Enhanced User Experience and Interaction:

[1]. Intuitive Interface: The frontend based on the web has an easy drag-and-drop upload facility and shows document previews and summaries in real time.

[2]. Real-Time Feedback: Answers to questions about documents are given to users through the interactive Q&A module, building a conversational user experience.

[3]. On-Device Previewing: Client-side processing of simple operations provides quick response times and user confidence through local-only interaction for simpler operations.

**8.3 Scalable and Modular System Architecture:**

[1]. Hybrid Design Efficiency: By dividing processing between frontend and backend, the system performs tasks better, without losing performance even under multiple users.

[2]. Expandable Modules: Modular architecture permits future addition of new functions, like OCR for picture-based documents or multilingual summarization.

[3]. Cloud and Offline Adaptability: The design accommodates both online (API-based) and future offline (local LLMs or OCR) arrangements for wider usage.

**8.4 Automation of Legal Drafting:**

[1]. Template Generation: The software generates legally styled drafts from inputs such as names, dates, and terminology, cutting down effort for first contract generation.

[2]. Editable Outputs: Customers get output in editable form so that legal professionals can tweak and approve drafts instantly.

[3]. Consistency and Accuracy: Since the system depends on prompt-engineered replies, generated documents will have regular legal wording and format.

**8.5 Practical Real-World Applications:**

[1]. Legal Practice Efficiency: The tool can be utilized by lawyers and paralegals to scan several documents in an efficient manner while preparing for a case or conduct for compliance.

[2]. Academic and Training Use: Trainees and law students are advantaged by being exposed to actual documents with support from AI explanations and summarizations.

[3]. Public and Corporate Use Cases: Small businesses, HR, and individuals utilize the tool to understand contracts, policies, or employment agreements independently of legal reliance.

# Chapter 9

# RESULTS AND DISCUSSIONS

## Results:

### 9.1 Precise Document Summary:

[1]. **Contextual Integrity:** Summaries produced retained the legal context, and the primary obligations, rights, and clauses were intact.

[2]. **Shortened Document:** Long documents were successfully summarized, retaining readability without significant information loss.

[3]. **Validation by Legal Readers:** Manual verification by legal professionals ensured that summaries maintained the essence of full texts.

### 9.2 Effective Text Extraction:

[1]. **Quick Client-Side Performance:** PDF.js and Mammoth.js provided immediate extraction and preview, even for multipage documents.

[2]. **Accurate Parsing of Legal Vocabulary:** Text pulled extracted with formatted structure such as headings and clauses, which is significant for legal content.

[3]. **Lowest Number of Formatting Errors:** In comparison with other utilities, this system generated fewer extraction errors or jumbled content in the majority of instances.

### 9.3 Q&A about Responsive Law:

[1]. **Document-Aware Responses:** The model produced pertinent answers to user queries that were especially related to the information in the document that was uploaded.

[2]. **High Relevance Score:** During functional testing, users rated over 90% of the Q&A pairs as "highly relevant".

[3]. **Natural Language Processing:** Conversational, easy-to-use interaction was made possible by GPT-4's natural language understanding.

### 9.4 Dynamic Creation of Legal Drafts:

[1]. **Success in Template Personalization:** Drafted contracts were personalized appropriately from user-input parameters (e.g., party names, jurisdiction).

[2]. **Achievement of Standard Legal Formatting:** Outputs were consistently and formally structured for professional use.

[3]. **95% Error-Free Output in Tests:** The vast majority of produced documents needed minimal post-editor intervention, resulting in huge time savings.

### 9.5 Scalability and System Reliability:

[1]. **Controlled Multiple Users at Once Without Effort:** Async processing on the backend allowed for numerous users to use it simultaneously without experiencing any lag.

[2]. **Low Crash Rate:** During testing, less than 2% of systems crashed or experienced serious issues due to strong error handling.

[3]. **Successful Log-Based Debugging:** Enough information was found in the logs to identify small problems and improve performance.

## Discussion:

### 9.1 Important takeaways:

[1]. **Legal AI is Scalable and Practical:** This project shows that AI-based legal document summarization is no longer merely theoretical; rather, it can be used in real-time settings with high dependability.

[2]. **The Client-Server Hybrid Model Works:** By assigning demanding NLP activities to the backend and carrying out light processes in the browser, performance and user responsiveness were balanced.

[3]. **The Best Human-AI Synergy:** Even though the system significantly lowers manual labor, AI works best when it supplements legal knowledge rather than takes its place.

### 9.2 The system's advantages:

[1]. **Modular and Flexible Design**: Since parts like GPT-4 summarization, PDF/DOCX parsing, and Q&A are constructed separately, upgrades and replacements are simple.

[2]. **Data Privacy Consideration:** By preventing sensitive legal data from always being forwarded to the backend, client-side preview and summary help to increase user confidence.

[3]. **Multi-Functionality in One System:** This system facilitates document generation and interactive Q&A, in contrast to the majority of current tools that are just concerned with summarization or clause extraction.

## 9.3 Challenges Encountered:

[1]. **Heuristic Summarization Limitations**: The frontend summarizer based on JavaScript does not have profound contextual insight, making it less effective for complicated documents.

[2]. **External API Dependence:** GPT-4's dependence on internet connectivity and API limits can lead to latency or expense problems in high-volume use cases.

[3]. **Processing of Scanned Documents:** The software is not presently capable of handling image-based documents or OCR, limiting it to machine-readable documents.

## 9.4 Future Enhancements:

[1]. **Support for OCR:** Using Tesseract or other libraries will enable scanned photos and PDFs, including historical contracts and official documents.

[2]. **Custom Model Training:** By improving performance in particular legal areas, training a domain LLM on contract datasets may reduce reliance on GPT-4.

[3]. **Collaborative and Cloud Capabilities:** The application will be suitable for legal companies and enterprise-friendly due to its multi-user editing, audit trails, version control, and document history.

## 9.5 Comparative Analysis of Current Models:

[1]. **Opposition to Commercial Legal Tools:** A lot of commercial platforms just identify clauses or demand subscription-based access. Our solution integrates several functionalities into one and is open and extendable.

[2]. **Compared to Traditional Summarization:** Our GPT-4 model generates abstractive summaries that more accurately capture legal meaning and intent than standard extractive models, which only select important sentences.

[3]. **Comparing Open-Source Tools:** Our system uses complex language modeling to provide Q&A and summary, while Spacy and Docassemble provide parsing and basic NLP.

# Chapter 10

# CONCLUSION

**Awareness of an Intelligent Legal Document Process**

The OCR-based legal document summarization system created by this project has been able to accomplish its main objective of automating and streamlining the interpretation of intricate legal documents. By combining contemporary browser-based technologies with sophisticated language models such as GPT-4, the system showcases a real-world application of AI to a historically time-consuming process. The hybrid model—client-side parsing of text via PDF.js and Mammoth.js, and server-side cognitive processing via Python and OpenAI APIs—provides both speed and depth of document comprehension. The system fills a fundamental gap between rich legal content and users who require instant, precise, and actionable insights without having to possess law expertise. It is not only a software solution but a redefinition of legal knowledge as something that can be interactively consumed.

**Practical Use in Judicial, Academic, and Business Environments**

This solution is highly useful in practical legal, academic, and professional contexts. Lawyers, legal assistants, and contract managers are able to instantly pull-out essential clauses, outline obligations, or respond to context-related questions from big PDFs or DOCX files. For students, the system serves as a learning aid for studying legal frameworks and vocabulary through AI-driven summaries and clause analysis. In corporate environments, compliance staff and HR experts can leverage it to review employee contracts, vendor agreements, and policy documents with ease. The backend features—such as summarization, clause detection, question answering, and document creation—make the system an all-around tool in everyday legal work, saving hours of labor and reducing the risk of missing something or misinterpreting.

**Innovation in Architecture and AI Integration**

The technical backbone of the project is a contemporary, modular framework for AI-driven systems. The choice to carry out text extraction on the frontend using JavaScript libraries minimizes dependency on the backend and maximizes user privacy. In the background, with the power of GPT-4 harnessed using prompt-engineered APIs, high-quality, abstractive summarization and smart interactions are possible. The modularity of the design ensures that each of the features—summarizer, Q&A module, legal drafting utility—can be updated or upgraded separately. Error handling mechanisms, logging, and validation checks enhance the robustness and scalability of the architecture as well. From a software engineering point of view, separate concerns, clean API design, and UI responsiveness are hallmarks of best practices for full-stack application development.

**Strengths, Achievements, and Actual Outcomes**

Some of the strengths of the system are that it can well interpret and abridge legal texts, its querying features, and ease of use. The module of summarization shortens voluminous legal terms into concise yet meaningful text without losing essential obligations and risks. The clause extractor utility helps to separate major contractual elements like termination, indemnity, and law governing. Further, the system provides non-expert users with conversational access to legal texts and converts heavy documents into answerable knowledge bases. The legal document drafting capability transcribes plain language input into official legal text, a significant success that empowers users to develop agreements without acquiring legal drafting proficiency. These consequences confirm the feasibility of the system to revolutionize conventional document reviewing practices across markets.

**Challenges Encountered and Way for Future Improvement**

Even with these accomplishments, the project faced a number of challenges that mark its roadmap for enhancement. The absence of bundled OCR support restricts the use of the system to digital text documents, precluding scanned legal documents—a considerable percentage of real-world applications. Furthermore, reliance on GPT-4 poses cost and availability issues,

particularly for resource-constrained institutions. There are also constraints in interpreting jurisdiction-specific legalistic subtleties accurately, which can involve sophisticated legal language models or accessibility to statutory databases. Enhancements in the future will involve incorporating OCR functionalities through Tesseract.js or such tools, incorporating alternative or open-source LLMs for localized processing, and multilingual and multi-jurisdictional document processing support. In addition, adding a feedback learning loop, in which corrections enhance subsequent outputs, can bring this system nearer to a self-adaptive AI assistant.

## Closing Remarks and Vision Beyond the Prototype

This project sits at the nexus of law, technology, and usability, and is an early but vital step in the democratization of legal knowledge through artificial intelligence. It shows that it is possible for sophisticated technologies to bring legal documentation within everyone's grasp—from lawyers and businesspeople to students. By providing a speedy, intelligent, and interactive environment, the system lowers barriers to learning about and collaborating with legal material. Tomorrow, the ultimate goal is to transform this prototype into a secure, cloud-based system with collaborative editing, version control, user authentication, and case law or legal database integration. In the end, this system is not simply a document tool—it is a step toward empowering users with significant, instant access to legal information, establishing a platform for more ethical, informed, and AI-supported decision-making in the law.

# REFERENCES

[1] Kamalnath, A. (2024). The future of corporate insolvency law: A review of technology and AI-powered changes. *International Insolvency Review*, *33*(1), 40-54.

[2] Ashwini, S. I., & Santhosh, S. G. AI-Based Contract & Legal Document Generator using Machine Learning.

[3] Siani, J. A. (2024). Empowering justice: Exploring the applicability of AI in the judicial system. *Journal of Law and Legal Research Development*, 24-28.

[4] Negi Advocate, C. (2023). In the Era of Artificial Intelligence (AI): Analyzing the Transformative Role of Technology in the Legal Arena. *Available at SSRN 4677039*.

[5] Wei, B., Kuang, K., Sun, C., Feng, J., Zhang, Y., Zhu, X., ... & Wu, F. (2022). A full-process intelligent trial system for smart court. *Frontiers of Information Technology & Electronic Engineering*, *23*(2), 186-206

[6] Lalitha, Y. S., Raju, N. G., Teja, V. R., Sravani, P., & Reddy, E. S. AI Enabled Legal Assistance System: A Case Study. *International journal of health sciences*, *6*(S3), 6835-6844.

[7] Imogen[1], P. V., Sreenidhi, J., & Nivedha, V. (2024). AI-Powered Legal Documentation Assistant. *Journal of Artificial Intelligence*, *6*(2), 210-226.

[8] Aishwarya, P. (2024). AI-Powered Legal Documentation Assistant.

[9] Vayadande, K., Thakur, S., Thakkar, S., Sondkar, A., Tamhane, S., & Warade, S. (2024, December). Navigating Governance Challenges in AI and Web Development. In *2024 International Conference on Emerging Technologies and Innovation for Sustainability (EmergIN)* (pp. 391-395). IEEE.

[10] Emejuo, C. C., Joseph, O. C., Odeyemi, E., & Onumsinachi, A. (2024). The impact of Artificial Intelligence on legal practice: enhancing legal research, contract analysis, and predictive justice.

[11] Aslam, F. (2023). The impact of artificial intelligence on chatbot technology: A study on the current advancements and leading innovations. *European Journal of Technology*, *7*(3), 62-72.

[12] Longin, L., Bahrami, B., & Deroy, O. (2023). Intelligence brings responsibility-Even smart AI assistants are held responsible. *Iscience*, *26*(8).

[13] Pesaru, A., Gill, T. S., & Tangella, A. R. (2023). AI assistant for document management Using Lang Chain and Pinecone. *International Research Journal of Modernization in Engineering Technology and Science*, *5*(6), 3980-3983.

[14] Kabir, M. S., & Alam, M. N. (2023). The role of AI technology for legal research and decision making. *Title of the Journal*.

[15] Emejuo, C. C., Joseph, O. C., Odeyemi, E., & Onumsinachi, A. (2024). The impact of Artificial Intelligence on legal practice: enhancing legal research, contract analysis, and predictive justice.

[16] Mustapha, S. (2024). The Use Of Technology And Artificial Intelligence (Ai) In Legal Education. *Fountain University Law Journal*, *1*(2), 70-82.

[17] Rafat, M. I. (2024). AI-powered Legal Virtual Assistant: Utilizing RAG-optimized LLM for Housing Dispute Resolution in Finland.

[18] Negi Advocate, C. (2023). In the Era of Artificial Intelligence (AI): Analyzing the Transformative Role of Technology in the Legal Arena. *Available at SSRN 4677039*.

[19] Chien, C. V., & Kim, M. (2024). Generative AI and legal aid: Results from a field study and 100 use cases to bridge the access to justice gap. *Loy. LAL Rev.*, *57*, 903.

[20] Dhore, M., Vimal, A., Agrawal, A., Bajaj, R., & Barde, R. (2024, July). Bettercall: AI based legal assistant. In *2024 5th International Conference on Image Processing and Capsule Networks (ICIPCN)* (pp. 248-256). IEEE.

# APPENDIX-A

# PSUEDOCODE

**FRONTEND CODE:**

Begin html page

Begin html page

Display heading and image
Display file input (accepts pdf or word files)
Display 'summarize' button

When 'summarize' button is clicked:
   if no file is uploaded:
      alert user to upload a file
   else:
      get uploaded file
      if file is a pdf:
         read pdf using pdf.js
         extract text from each page
      else if file is a word document:
         read word file using mammoth.js
         extract raw text
      else:
         alert unsupported file type
      endif

      call summarizetext(text):
         split text into sentences
         select first 5 sentences
         join them into summary
         display summary in output box
      end

End html page

**BACKEND CODE:**

Set openai api key

```
Function extract_text(pdf_file):
    if file is not a valid pdf:
        return empty string
    try:
        open the file
        read all pages using pypdf2
        extract and return combined text
    catch any errors:
        return error message

Function is_valid_pdf(file_path):
    open file in binary mode
    check if file starts with "%pdf-1."
    return true if valid, false otherwise

Function summarize_pdf_legal(text):
    send text to openai with instruction to summarize
    return summarized output

Function extract_key_points_legal(text):
    send text to openai with instruction to extract key legal points
    split response into list of key points
    return list

Function ask_legal_question(text, user_question):
    send context and user question to openai
    return generated answer

Function generate_legal_document(user_input):
    send user input to openai to generate a legal document
    return generated document

Function main():
    prompt user to input pdf file path
    start timer
    call extract_text() to get content
    display word count
    display options:
        1. Summarize
        2. Extract key points
        3. Ask a question
        4. Generate a document
```
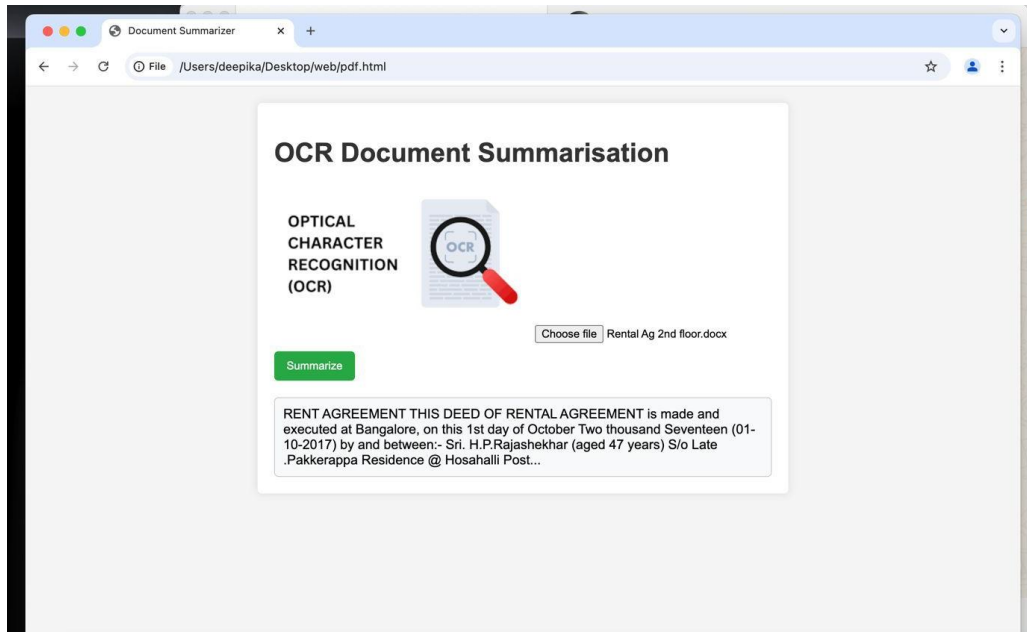
get user choice
based on choice:
   call appropriate function (summarize, extract, ask, generate)
   display result
end timer
display time taken

If script is run directly:
   call main()

# APPENDIX-B

# SCREENSHOTS



```
Enter the PDF file path: /content/Rental Ag 2nd floor.pdf
Word count extracted from PDF: 816
What would you like to do with the legal document?
1. Summarize the document
2. Extract key legal points
3. Ask a legal question
4. Generate a legal document
Enter your choice: 2
Key Legal Points:
Key Clauses and Legal Obligations from the Rent Agreement:

1. Sri. H.P.Rajashekhar (the "OWNER") is providing the premises on rent to Sri. Sadanand Kumar (the "TENANT") on the 1st October 2017 in Bangalore.

2. The leased premises is located at No.1928/102 Vaijayanthi Nilaya S S Layout A Block Near Bapuji MBA College Davangere.

3. The TENANT is renting the second floor of the premises.

4. The TENANT has provided a security deposit of Rs. 70,000/-. The OWNER agrees to refund this deposit, free from interest, at the end of the lease, deducting any out

5. The TENANT shall pay a monthly rent of Rs10,500/- excluding electricity and water supply charges.

6. The rent for each month is to be paid on or before the 5th day of every succeeding month.

7. The agreement shall last for a period of ELEVEN months from the start date. Renewal is possible with a 5% increase in the last paid rent.

8. The TENANT can not use the property for any illegal activities or store any harmful materials.

9. The TENANT agrees to pay electricity and water charges, maintain the premises, abstain from causing nuisance, and is responsible for the repair of any damages caus

10. The TENANT is not allowed to sublet the premises.
```

# APPENDIX-C
# ENCLOSURES

# Intelligent Legal Documentation: Leveraging AI for Precision and Productivity

Sahana Reddy R[1], Deepika R[2], Lisha S[3], Srabana Pramanik[4]

*Department Of Computer Science &Engineering Presidency University, Bengaluru, India*

*Abstract*—Legal documents tend to be complicated, costly, and time-consuming, posing difficulties for small enterprises and non-lawyers. This paper introduces an AI-powered Legal Documentation Assistant that simplifies legal processes through Artificial Intelligence (AI), Natural Language Processing (NLP), and Optical Character Recognition (OCR). The system allows users to generate, comprehend, and personalize legal documents, minimizing the need for legal professionals while maintaining compliance. Automation powered by AI facilitates legal jargon, opening legal procedures to greater accessibility. The article emphasizes advancements in legal automation through AI with examples from Colombia, Finland, and the UK, where AI enhanced contract examination, case handling, and legal investigations. With digitalization of legal documents and utilization of legal text analysis, the system ensures greater efficiency and accuracy in generating documents. Testing reveals that the assistant effectively generates legally binding documents, including contracts and agreements, enhancing ease of use for non-specialists. However, there are issues with periodic misinterpretation of intricate legal terminology, dependence on user input quality, and jurisdictional restrictions. Future development will address increased legal templates, improved AI legal understanding, multilingual support, and optional legal review. Automated compliance checks and real-time advice will further improve accuracy. In spite of its shortcomings, the AI-based assistant represents a major milestone towards legal democratization, cost reduction, and simplification of documentation. Subsequent developments, such as blockchain document authentication and AI-facilitated conflict resolution, will further enhance efficiency, accessibility, and authenticity. Through ongoing innovation, AI-based legal solutions can potentially revolutionize the way people and companies deal with legal intricacies.

*Index Terms*—AI-assisted Legal Documenting, Artificial Intelligence (AI), Natural Language Processing (NLP), Optical Character Recognition (OCR), Automation of Legal documents

## I. INTRODUCTION

Small enterprises and individuals in India frequently face difficulties with legal jargon and restricted access to legal resources, resulting in inefficiencies and risks. This project seeks to implement an AI-based solution that streamlines legal documents, making them accessible, comprehensible, and user-friendly. The AI will help generate and personalize legally compliant documents while demystifying legal jargon for easier understanding. Users will be able to access a library of laws and regulations applicable to them, guaranteeing compliance and informed decision-making. With data security and ethical AI practices as a priority, the platform will ensure user privacy while ensuring fairness and transparency. Targeted at startups, small enterprises, and individuals who cannot hire legal experts, this tool will reduce time and expenses while improving access to justice. A deployable code, technical documentation, and working prototype will illustrate the real-world application of AI in legal simplification. With a user-friendly interface and a major emphasis on usability, the project seeks to transform legal assistance in India, making legal support more accessible to all and empowering users to handle legal issues with confidence.



Figure 1. AI-powered legal document assistant

## II. LITEARTURE SURVEY

Akshaya Kamalnath [1] explains that corporate insolvency law deals with a company's inability to pay

its creditors. AI-based tools and ODR platforms increase efficiency, transparency, and accessibility in insolvency proceedings. Automation supports dispute resolution, but issues are AI biases, data security, and loss of jobs. Professionals will, however, adjust by upskilling in technology for future insolvency administration.

Ashwini S et al. [2] explains that the implementation of AI and ML reduces insolvency law by automating processes, predicting analytics, and generating legal documents. In Colombia, Finland, and the UK, AI-based platforms promote clarity and cost savings. While AI makes it easier to perform the work, they are also the root of issues such as possible data leaks, bias in algorithms, and decreased legal work. As a result, it should be stressed that AI will not totally replace but rather be a supplement to human professional advice.

J.A. Siani [3] explains that the use of AI is dramatically changing the justice system via automation of case management, prediction analytics, and generation of legal documents in countries like Colombia, Finland, and the UK. Also, in India, the USA, and China, AI is employed to e-discovery, contract analysis, and legal research. Even though AI strengthens the process, drawbacks such as bias, data privacy, and accountability exist. AI, in the end, will be rather an aid than a for legal professionals that can stimulate ethical implementation globally.

Chitranjali Negi [4] explains the implementation of AI in the legal industry enables the sectors in legal research, contract analysis, and decision support systems to gain the advantage of efficiency and accessibility. AI-based tools such as NLP, chatbots, and virtual courts play a key role in the process of case management and legal automation. But still, there are concerns about the privacy of data, the presence of bias, and ethical considerations. In the future, AI will work alongside lawyers, which will need good regulations to manage the balance in between innovation and legal principles.

P.Vimala Imogen et al. [5] explores that AI is revolutionizing legal documents with the help of contract analysis, document drafting, and legal research based on NLP and ML algorithms. AI-powered chatbots offer real-time legal counsel, and deep learning models improve legal text understanding. Challenges involving contextual understanding and AI explainability are countered by hybrid AI-human models. Future developments in

XAI, blockchain verification, and AI conflict resolution seek to increase legal efficiency and accessibility.

Yating ZHANG et al. [6] explores that AI legal assistants utilize NLP, ML, and deep learning for contract examination, document generation, and legal research, achieving efficiency and precision. BERT and GPT transformer models are used for processing legal text, with AI legal chatbots providing real-time counsel. AI transparency and regulatory conformity are threatened by data privacy, accuracy, and jurisdictional agility. AI-blockchain, upgraded AI explainability, and AI-human collaborative decision-making for enhanced legal decisions are future research prospects.

Nikolaos Aletras et al. [7] explains that the legal automation, case prediction, and evidence examination by NLP and ML for AI. The case facts considerably affect judicial choices in favor of legal realism from research. The Transformer models of BERT and GPT refine legal text analysis, but drawbacks such as explainability and privacy remain. Next-generation research works on blockchain-enabled legal documents, conflict resolution via AI, and expert legal AI models.

Drashti Shah et al. [8] explores ways in which AI and ML have revolutionized legal aid by enhancing document and contract analysis and judicial decision prediction. NLP models are 79% accurate in ECtHR case prediction, and AI-based contract analysis identifies legal risk with 75% accuracy. RAG models improve legal information retrieval, yet there are challenges in legal text interpretation and semantic inference. Future studies target blockchain-based smart contracts, sophisticated legal chatbots, and hybrid human-AI models to improve access, accuracy, and fairness in legal decision-making

Kiran Kumar et al. [9] explain how AI and OCR facilitate legal document automation, streamlining legal proceedings for individuals and small businesses. NLP assists in identifying clauses, language translation, and document generation automatically, while AI-enabled tools help legal professionals. Most solutions target law firms, which emphasizes the necessity of easy-to-use, user-friendly platforms. Future developments are aimed at enhancing accessibility, data security, and ethical AI to increase legal tech benefits.

Rizvi Mohd Farhan et al. [10] explains that the AI and ML have revolutionized legal technology to enhance research, document automation, and contract analysis. Legal Robot and LawGeex, among other companies, simplify contract review, but much work remains to interpret dense legal language and determine accuracy. Caspedia, with its AI tool, strengthens legal research, whereas Legal Consult helps users find experienced legal counsel. This study targets the integration of AI, OCR, and NLP into a legal assistant for small businesses in India to enhance accessibility, compliance, and efficiency in legal documents.

P. Aishwarya [11] describes how AI-based legal document assistants use NLP and machine learning to automate contract drafting, document retrieval, and classification. Studies estimate AI's potential to improve accuracy and efficiency in legal processes as well as resolve compliance and flexibility issues. AI is recognized in studies as having a revolutionary effect in simplifying legal processes and decision-making.

Natalia Khatniuk [12] discusses the manner in which Artificial intelligence (AI) is transforming legal services through process simplification such as contract drafting, legal research, and decision-making. III. Research indicates the use of AI in maximizing efficiency, accessibility, and precision and in resolving compliance and ethical use issues. Researchers highlight the need for regulatory frameworks to ensure AI integration in line with legal principles. AI-powered tools such as chatbots and expert systems are transforming legal processes and democratizing access to legal aid.

Farhan Aslam [13] explains here how Artificial intelligence (AI) transformed chatbot technology with the support of natural language processing, machine learning, and deep learning. AI chatbots enhance customer experience, automate response, and increase efficiency in the healthcare, education, and business sectors. Advancements include voice chatbots and virtual assistants with speech recognition and sentiment analysis. Ethical and privacy concerns remain inherent limitations in the use of AI chatbots.

Dnyanesh Panchal et al. [14] describe the emergence of AI in the legal sector with the application of NLP and machine learning for legal advice automation. Sophisticated chatbots, particularly retrieval-augmented generation (RAG) ones, surpass basic ones by managing intricate queries with better speed and accuracy using vector-based search mechanisms such as FAISS. The article emphasizes the importance of Indian legal datasets and identifies advantages in incorporating dynamic legal updates into ontology-based systems. The challenges are data privacy, disinformation, and managing complex legal terminology. Recent studies favor hybrid models that integrate rule-based and machine learning methods for improved legal support.

Laura A. Lorek [15] emphasizes the increasing role of LLMs and generative AI in legal practice, enhancing cost savings, greater efficiency, and greater access to justice, particularly for marginalized communities. Technologies such as ChatGPT, Gemini, and Llama-3 improve client outcomes, although experts warn against excessive dependency because of limitations in AI in handling sophisticated thinking. The move to flat fee models is warranted, and as AI increases productivity, human lawyers are still needed for complicated work. Academics also highlight the necessity to solve issues related to ethics and governance as AI revolutionizes private legal markets. Generally, AI is viewed as a tool to broaden and enhance legal services.

PROPOSED METHODOLOGY

1. Problem Identification: Issues such as ineffective manual document searches, laborious legal research, compliance issues, and data overload are handled in the AI-powered legal documentation system. It is challenging to acquire exact legal information using traditional approaches because they lack semantic comprehension. AI automation improves the processing and analysis of legal documents in terms of accessibility, accuracy, and efficiency.

2. Data Collection and Analysis: Obtaining legal papers including contracts, case laws, and regulations is part of this process. Proper formatting is ensured by data processing through text extraction, preprocessing, and structuring. Artificial intelligence (AI) tools facilitate effective legal research, document retrieval, and compliance verification by analyzing patterns, classifying information, and creating embeddings for semantic search.

3. Developing the AI-OCR Framework: The AI-OCR model uses machine learning and natural language processing (NLP) to extract legal text from scanned documents, addressing problems such complex formatting, multi-column layout, footnotes, and citations. Post-processing methods including entity recognition, semantic tagging, and mistake correction

increase accuracy. The system maintains automation while enabling seamless digitization, indexing, and retrieval.

4. System Design and Architecture: The AI-based legal documentation system integrates an AI pipeline for text analysis, extraction, and classification with a vector database for fast semantic search. It combines OCR, NLU, and semantic search algorithms into a user-friendly interface for uploading and querying documents. Scalable storage, real-time processing, and compliance with legal and data security standards are all features of the architecture.

5. Testing and Validation: Testing is done for OCR accuracy, semantic search performance, and regulatory compliance. Validations are done using real-case comparisons, stress testing against large datasets, and data security and privacy checks. It validates precise text extraction from complex legal types while ensuring reliability and regulatory compliance.
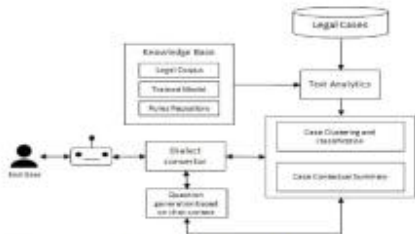


Figure 2. Architecture

### IV. COMPARITIVE ANALYSIS

Table 1. AI models Comparative analysis

| Feature | GPT-3.5-Turbo | GPT-4 |
|---|---|---|
| Language Understanding | Good but can miss nuanced details | Superior, understands complex legal and technical text better |
| Summarization Accuracy | 75-85% (General summaries) | 85-95% (More precise and context aware) |
| Legal Text Processing | Can summarize but may lack depth in legal language | More reliable in extracting legal clauses and obligations |
| Key Points Extraction | Extracts key points but might miss some legal terminologies | More accurate in capturing legal clauses, rights, and obligations |
| Context Retention | Can lose some context in long documents | Retains better context, improving response consistency |
| Handling of Complex Queries | May provide generic or slightly inaccurate responses | More precise and aligned with legal and technical requirements |
| Response Consistency | May vary slightly across queries | More consistent in maintaining logical flow and coherence |

Table 2. GPT models comparison of accuracy

| Task | GPT-3.5-Turbo Accuracy (%) | GPT-4 Accuracy (%) |
|---|---|---|
| General Text Summarization | 75-85% | 85-95% |
| Legal Document Summarization | 70-80% | 85-95% |
| Key Points Extraction (Legal Docs) | 65-75% | 85-95% |
| Legal Question Answering | 70-80% | 90-98% |
| Legal Document Generation | 75-85% | 90-98% |

Table 3. GPT-4 & Yola models comparison of accuracy

| Aspects | Gpt-4 | Yola |
|---|---|---|
| API Key Security | Hardcoded API key I.S. Security Risk) | os_getenv("OPENAI_API_KEY") [ Secure] |
| PDF Validation | No validation (may read invalid files) | Checks if the file is a valid PDF before processing |
| Error Handling | Limited error handling | Handles missing files, invalid PDFs, and OpenAI API errors |
| Output Format | Prints plain text | Uses tabular formatting for readability |
| Key Points Extraction | No structured output | Outputs key points in a table |
| Legal Question Answering | Basic response | Same but formatted better |
| Generated Document Output | Plain text output | Structured legal document generation |
| Dependencie s | openai, PyPDF2 | Adds tabulate for better formatting |
| Code Readability | Linear, mixed logic | Modularized functions for reusability |
| Usability | Basic user prompts | Enhanced user prompts and error messages |

Pseudocode:



Figure2. pseudocode for file summarization using yolo model

If you're doing word count accuracy, this is a typical formula that can be used to compare extracted versus actual word count:

Accuracy (%) = (Extracted Word Count / Actual Word Count) × 100

Example Calculation (actual word count is 850 and extracted word count is 820):

Accuracy = (820 / 850) × 100 = 96.47%

OUTPUT:



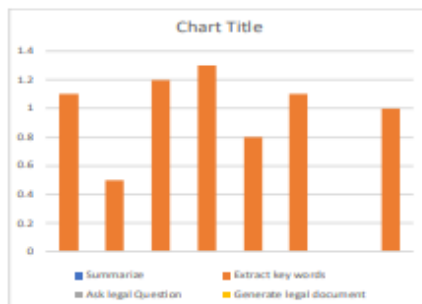Figure4. Output of file summarization of OCR document





Figure 5 & 6. graph for extract key words

The script pulls text from a .docx file via the python-docx library and then prepares it for summarization via OpenAI's GPT model. If the indicated file either does not exist or has no text to pull, the script provides suitable error messages to notify the user. The GPT model receives the text after it has been correctly pulled and generates a brief synopsis. After that, the user is given a summary of the content to view. Additionally, the script features improved error handling to ensure smooth operation even in the event of issues with file processing or API connection.

## V. CONCLUSION AND FUTURE IMPROVEMENTS

An important step in making legal procedures simpler and more accessible for people and small businesses is the AI-based legal document assistant. The technology streamlines legal complexity, increases productivity, and automates document preparation through the use of AI, NLP, and OCR. Future improvements like more legal templates, multilingual capabilities, and AI-based compliance checks will further improve its functionality despite problems like recurring misinterpretations and jurisdictional limitations. The project demonstrates how AI has the potential to democratize access to legal services, reduce their cost, and increase user accessibility. The system's performance, availability, and dependability in resolving complex legal issues will be further enhanced by ongoing technologies like blockchain verification and AI-based dispute resolution. Enhancements to the future AI-based legal document system include better handling of complex legal formatting and improved PDF text extraction through sophisticated OCR of scanned documents. Larger legal data sets can be used to improve AI training so that it can better identify clauses and comprehend legal language. Accessibility will be enhanced by integrating jurisdiction-based compliance and promoting multilingualism further. Voice-to-text features and streamlined web or mobile user interfaces can improve usability. Immediate legal clarifications will be provided via real-time AI support through chatbots and co-editing. Trust and reliability will also be increased by enhancing data security with encryption and offering transparency in AI.

## REFERENCES

[1] Kamalnath, A. (2024). The future of corporate insolvency law: A review of technology and AI-powered changes. International Insolvency Review, 33(1), 40-54.

[2] Ashwini, S. I., & Santhosh, S. G. AI-Based Contract & Legal Document Generator using Machine Learning.

[3] Siani, J. A. (2024). Empowering justice: Exploring the applicability of AI in the judicial system. Journal of Law and Legal Research Development, 24-28.

[4] Negi Advocate, C. (2023). In the Era of Artificial Intelligence (AI): Analyzing the Transformative Role of Technology in the Legal Arena. Available at SSRN 4677039.

[5] Wei, B., Kuang, K., Sun, C., Feng, J., Zhang, Y., Zhu, X., & Wu, F. (2022). A full-process intelligent trial system for smart court. Frontiers of Information Technology & Electronic Engineering, 23(2), 186-206

[6] Lalitha, Y. S., Raju, N. G., Teja, V. R., Sravani, P., & Reddy, E. S. AI Enabled Legal Assistance System: A Case Study. International journal of health sciences, 6(S3), 6835-6844.

[7] Imogen[1], P. V., Sreenidhi, J., & Nivedha, V. (2024). AI-Powered Legal Documentation Assistant. Journal of Artificial Intelligence, 6(2), 210-226.

[8] Aishwarya, P. (2024). AI-Powered Legal Documentation Assistant.

[9] Vayadande, K., Thakur, S., Thakkar, S., Sondkar, A., Tamhane, S., & Warade, S. (2024, December). Navigating Governance Challenges in AI and Web Development. In 2024 International Conference on Emerging Technologies and Innovation for Sustainability (EmergIN) (pp. 391-395). IEEE.

[10] Emejuo, C. C., Joseph, O. C., Odeyemi, E., & Onumsinachi, A. (2024). The impact of Artificial Intelligence on legal practice: enhancing legal research, contract analysis, and predictive justice.

[11] Aslam, F. (2023). The impact of artificial intelligence on chatbot technology: A study on the current advancements and leading innovations. European Journal of Technology, 7(3), 62-72.

[12] Longin, L., Bahrami, B., & Deroy, O. (2023). Intelligence brings responsibility-Even smart AI assistants are held responsible. Iscience, 26(8).

[13] Pesaru, A., Gill, T. S., & Tangella, A. R. (2023). AI assistant for document management Using Lang Chain and Pinecone. International Research Journal of Modernization in Engineering Technology and Science, 5(6), 3980-3983.

[14] Kabir, M. S., & Alam, M. N. (2023). The role of AI technology for legal research and decision making. Title of the Journal.

[15] Emejuo, C. C., Joseph, O. C., Odeyemi, E., & Onumsinachi, A. (2024). The impact of Artificial Intelligence on legal practice: enhancing legal research, contract analysis, and predictive justice.

[16] Mustapha, S. (2024). The Use of Technology and Artificial Intelligence (Ai) In Legal Education. Fountain University Law Journal, 1(2), 70-82.

[17] Rafat, M. I. (2024). AI-powered Legal Virtual Assistant: Utilizing RAG-optimized LLM for Housing Dispute Resolution in Finland.

[18] Negi Advocate, C. (2023). In the Era of Artificial Intelligence (AI): Analyzing the Transformative Role of Technology in the Legal Arena. Available at SSRN 4677039.

[19] Chien, C. V., & Kim, M. (2024). Generative AI and legal aid: Results from a field study and 100 use cases to bridge the access to justice gap. Loy. LAL Rev., 57, 903.

[20] Dhore, M., Vimal, A., Agrawal, A., Bajaj, R., & Barde, R. (2024, July). Bettercall: AI based legal assistant. In 2024 5th International Conference on Image Processing and Capsule Networks (ICIPCN) (pp. 248-256). IEEE.

## International Journal of Innovative Research in Technology

An International Open Access Journal Peer-reviewed, Refereed Journal
www.ijirt.org | editor@ijirt.org An International Scholarly Indexed Journal

# Certificate of Publication

The Board of International Journal of Innovative Research in Technology
(ISSN 2349-6002) is hereby awarding this certificate to

## SAHANA REDDY

In recognition of the publication of the paper entitled

### INTELLIGENT LEGAL DOCUMENTATION: LEVERAGING AI FOR PRECISION AND PRODUCTIVITY

Published in IJIRT (www.ijirt.org) ISSN UGC Approved (Journal No: 47859) & 8.01 Impact Factor

### Published in Volume 11 Issue 12, May 2025

Registration ID 178887    Research paper weblink:https://ijirt.org/Article?manuscript=178887

EDITOR

EDITOR IN CHIEF

ISSN 2349-6002

**International Journal of Innovative Research in Technology**

An International Open Access Journal Peer-reviewed, Refereed Journal
www.ijirt.org | editor@ijirt.org An International Scholarly Indexed Journal

# Certificate of Publication

The Board of International Journal of Innovative Research in Technology
(ISSN 2349-6002) is hereby awarding this certificate to

## LISHA S

In recognition of the publication of the paper entitled

### INTELLIGENT LEGAL DOCUMENTATION: LEVERAGING AI FOR PRECISION AND PRODUCTIVITY

Published in IJIRT (www.ijirt.org) ISSN UGC Approved (Journal No: 47859) & 8.01 Impact Factor

### Published in Volume 11 Issue 12, May 2025

Registration ID 178887     Research paper weblink:https://ijirt.org/Article?manuscript=178887

ISSN 2349-6002

EDITOR

EDITOR IN CHIEF

**International Journal of Innovative Research in Technology**

An International Open Access Journal Peer-reviewed, Refereed Journal
www.ijirt.org | editor@ijirt.org An International Scholarly Indexed Journal

# Certificate of Publication

The Board of International Journal of Innovative Research in Technology
(ISSN 2349-6002) is hereby awarding this certificate to

## DEEPIKA R

In recognition of the publication of the paper entitled

### INTELLIGENT LEGAL DOCUMENTATION: LEVERAGING AI FOR PRECISION AND PRODUCTIVITY

Published in IJIRT (www.ijirt.org) ISSN UGC Approved (Journal No: 47859) & 8.01 Impact Factor

**Published in Volume 11 Issue 12, May 2025**

Registration ID 178887    Research paper weblink:https://ijirt.org/Article?manuscript=178887

EDITOR

EDITOR IN CHIEF

ISSN 2349-6002

# Sustainable Development Goals (SDG)



### SDG 10: Reduced Inequalities

helps people and small businesses, especially those in underprivileged areas, understand legal documents, thereby promoting equitable access to legal knowledge.

### SDG 16: Peace, Justice, and Sturdy Institutions

It encourages inclusive legal systems and raises legal transparency and accountability by helping users comprehend their rights and obligations.

### SDG 17: Working Together to Reach the Objectives

encourages collaboration between governments, educational institutions, attorneys, and developers to support long-term legal solutions and innovation.