# End-Term Report

## Inter-IIT Technical Meet 13.0
## Pathway

# PATHWAY LEGAL NAVIGATOR

**Submitted By:** Team-91

# Contents

# 1   INTRODUCTION

Due to the complexity and amount of legal precedents that must be analyzed in order to make educated choices, an Agentic Retrieval-Augmented Generation (RAG) system designed for legal case review is the need of the hour. Existing legal research tools lack full contextual knowledge and dynamic interaction, necessitating manual intervention to link similar cases or comprehend nuances. A well-designed Agentic RAG system can close this gap by automatically finding relevant legal precedents, analyzing case specifics, and providing accurate, contextual summaries. The lack of a stable, automated solution on the market emphasizes the need for innovation, which will save legal professionals important time and improve decision-making efficiency.

Our solution is designed to address a critical gap in the industry by providing legal professionals with a rapid, accurate, and robust Agentic-RAG system. This system aims to significantly enhance the efficiency of case analysis and improvement, reducing the time required for such tasks. By incorporating innovative techniques such as Guidance and DiskANN, our solution effectively addresses the challenges faced by current systems, ensuring superior performance and more reliable outcomes.

# 2   UNIQUENESS

## 2.1   LIMITATIONS OF CURRENT SYSTEMS

Current Legal RAG systems are far from perfect,they deal with a lot of issues and hence hold a large room for improvement. Some of them are listed below:

- Lack of Contextual Understanding: Many RAG systems struggle to grasp the nuanced legal context of cases, leading to irrelevant or shallow document retrieval and flawed summaries.

- Insufficient Precedent Mapping: Existing solutions often fail to effectively map legal input cases to precise and relevant precedents, limiting their practical utility for complex legal analysis.

- Hallucination of Facts: Legal RAG systems sometimes generate inaccurate or non-existent information, which can be detrimental in high-stakes scenarios.

- Limited Scalability: The large volume of legal documents and the complexity of case interconnections pose challenges for scaling existing systems while maintaining accuracy.

- Lack of Explanability:Many systems do not provide clear reasoning or evidence for their output, making it difficult for legal professionals to trust or verify their findings.

- Inflexibility in Input Formats:Many systems are limited in the types of legal queries or document formats they can handle, reducing usability.

## 2.2   PROPOSED SOLUTION

To address these challenges, we implemented a suite of methods specifically designed to mitigate the aforementioned limitations, thereby improving the accuracy, efficiency, and relevance of the results.

- Enhanced Precedent Mapping: Our solution improves the mapping of legal cases to relevant precedents through a dual-stage retrieval process. Initially, similar legal cases are identified and retrieved from "Indian-Kanoon" and stored in the Pathway live vectorstore. Subsequently, a retrieval agent performs a more targeted search within the vectorstore to ensure precise precedent mapping, optimizing legal research efficiency.

- Hallucination Agent: This agent evaluates the accuracy of LLM-generated outputs by cross-referencing the input query, retrieved documents, and the generated response. The output is classified as fully accurate, partially accurate, or hallucinated. In cases of unsatisfactory results, the retrieved documents are reintroduced to the LLM to produce a more accurate and refined output.

- Hypothetical Document Embeddings (HyDE): HyDE uses a Language Learning Model, like ChatGPT, to create a theoretical document when responding to a query, as opposed to using the query and its computed vector to directly seek in the vector database. It goes a step further by using an unsupervised encoder learned through contrastive methods. This encoder changes the theoretical document into an embedding vector to locate similar documents in a vector database. Rather than seeking embedding similarity for questions or queries, it focuses on answer-to-answer embedding similarity. Its performance is robust, matching well-tuned retrievers in various tasks such as web search, QA, and fact verification.

- Retrieval-Augmented Reasoning Enhancement(RARE) RARE brings forth the Retrieval-Augmented Factuality Scorer (RAFS) which uses external evidence for reasoning and grants logical consistency. Using Monte Carlo Tree Search, it explores multi-step reasoning paths focused on medical and commonsense tasks. Experiments with RARE show better results than base-line methods and match proprietary models like GPT-4. It demonstrates scalability, working effectively with open-source models without requiring additional fine-tuning.

- Pathway Vectorstore: Our model is capable of operating on a diverse array of document and query formats due to the dynamic Pathway vectorstore. By doing so, our methodology becomes considerably more user-friendly.

.

## 3 USE CASE SELECTION AND NOVELTY

### 3.1 USE CASES

Our use case centers around legal data. The following are the main use cases of our project:

- Legal Advisor: Given a legal document/ string and using pathway dynamic vectorstore, we create a dynamic vector storage for similar cases and relevant sections of Indian Penal Code(IPC) and constitution. A user query based on the given case is then answered by fetching relevant cases and law sections from the vectorstore, while simultaneously checking for hallucinations and quality of response through a hallucination agent. It is designed for legal advisors/ lawyers to identify strengths, weaknesses, loopholes and flaws in any given case.

- Legal Form Filling: By leveraging Pathway ETL and user-provided information, we developed a system capable of automatically populating legal forms with user details.It is also capable of generating customized legal forms, contracts,etc. in any custom context.

### 3.2 NOVELTY

We integrate the following novelties in our project:

- Explicit Reasoning Chains implemented with python: Using Python scripting, we developed multiple pipelines utilizing reasoning chains. Traditional methods often yield suboptimal results when dealing with large context sizes. To address this, we segmented the context into smaller sub-chains and applied distinct pipelines to each segment, enhancing overall performance.

- **Interpretation of Legislative Changes with Contextual Impact Analysis:** Our use case focuses on developing an advanced RAG (Retrieval-Augmented Generation) system designed to keep legal databases updated with new court rulings and judgments. Leveraging existing laws, historical verdicts, and ongoing case data, the system provides in-depth insights into user queries related to specific cases. Its innovation lies in its capability to map interconnected legal principles and identify analogous cases across diverse legal domains. Achieving this requires building extensive legal databases, implementing sophisticated contextual algorithms, and integrating cross-jurisdictional data to empower legal professionals with proactive strategy refinement.

- **Case Preparation:** Our solution leverages RAG to enhance case preparation by synthesizing legal precedents, client-specific contexts, and emerging legal theories. It dynamically adapts to each case, offering insights and trend analysis to predict potential regulations and discrepancies. The innovation lies in real-time integration of structured and unstructured data using open-source models, providing tailored, data-driven strategic support. This requires advanced data integration, contextual adaptation, and customizable frameworks to address specific client needs.

- **Microsoft Guidance AI:** Microsoft Guidance AI optimizes token usage and streamlines the orchestration of agent interactions, thereby facilitating the generation of contextually accurate and concise responses to intricate legal enquiries. In comparison to popular open-source models such as LlamaIndex, Guidance substantially reduces API token consumption, saving an average of 200 tokens per API call. It achieves this by further regulating the LLM's token generation process. It employs an innovative approach to constraint decoding that significantly reduces the utilization of our LLM and other resources, thus giving much better control over the format and structure over the output along with faster execution time.

## 4   SOLUTION OVERVIEW

The solution leverages the Pathway dynamic vector store to store relevant Indian Penal Code (IPC) articles and similar cases retrieved in real time through API calls to IndianKanoon, based on the user-provided case input. The user query is subsequently answered using context derived from the retrieved documents. Additionally, a hallucination-checking agent evaluates the quality of the generated output to ensure accuracy and reliability. This approach facilitates the generation of a detailed case review, providing valuable insights into the strengths and weaknesses of the case.

### 4.1   CORE CONCEPTS

- **Similar Cases:** Examination of precedent cases with comparable circumstances to gain insights from past judicial decisions, interpretations, and persuasive authority, thus enhancing case understanding.

- **Relevant laws:** Identification of applicable sections of the Indian Penal Code laws and regulations, and legal principles that form the foundational framework for legal interpretation in the current case.

- **Case Summarization:** Generating a summary of the input case is vital for retaining only the most relevant information, thereby enhancing the efficiency of the search process. This summary helps in constructing a more targeted and effective Pathway dynamic vector store, ensuring that only the most pertinent data is utilized, which in turn improves the quality of search results.

- **Query Enhancement:** Refining the input query to be more specific to the legal domain is crucial for obtaining more accurate results during similarity searches within the Pathway vector database. By

tailoring the query to the nuances of legal terminology and context, the system can more effectively identify relevant documents and improve the quality of the retrieval process.

- Hallucination Check: Hallucination checks are crucial for assessing the quality of the language model's responses and the relevance of retrieved documents. By providing actionable feedback, this process helps to improve the overall accuracy and reliability of the system's outputs, significantly enhancing the quality of analysis.

## 4.2 EXPLANATION OF APPROACH

- Pathway dynamic vector store is utilized to retrieve relevant laws and similar cases based on the input provided by the user.

- AI agents are deployed to perform query enhancement, case summarization and fetching similar cases.

- A comprehensive case review is generated providing a detailed analysis of the case.

- An AI agent is employed to evaluate the quality of the response and provide feedback to the language model regarding its accuracy.

# 5 SYSTEM ARCHITECTURE

## 5.1 COMPONENTS IN THE SYSTEM

Our legal assistance system was constructed such that it delivers the most relevant and accurate answers as far as the user's queries about his legal case are concerned. It employs a variety of agents that work together to provide the user with the most contextually accurate answers regarding their legal cases. Below is an in-depth analysis of each agent in our system along with their respective roles and sub-components for which they work together, thus providing holistic legal support:

- Case Summarizer Agent: The Case Summarizer Agent is the agent which takes care of parsing and compressing the admissible true data of the user's legal case into actionable data. This agent comprises two sub-agents specifically targeting the two most crucial aspect features of information extraction directed towards the rest of the data retrieval process.

  - Lexicographer Sub-Agent: This sub-agent carefully analyses each user's input to extract most informative key terms actually needed for the understanding of the case in a legal sense. Using natural language techniques, it could make sense of crucial words and phrases that speak to the very core elements of the case. It accepts the user's raw legal case, then using advanced NLP algorithms processes raw text recognizing the legal terminology, entities, and situational message cues. On the basis of this, it generates keywords that best shape the nature of the case to ensure relevance and specificity and identifies the most information-rich keyword from the case and passes it to the KanoonIQ Agent for mapping to legal precedents.

  - Case Summary Generation Sub-Agent: This sub-agent generates a comprehensive and accurate summary of a user's legal case encapsulating all the essential facts of the legal issues as well as of the relevant contexts and enables other agents to analyze the case deeply. It takes the context provided by the user as input. Then, the input gets synthesized to be a coherent and concise summary that focuses on the critical aspects and legal questions. It gets relevant background information that will make this description have a more complete picture of the case.
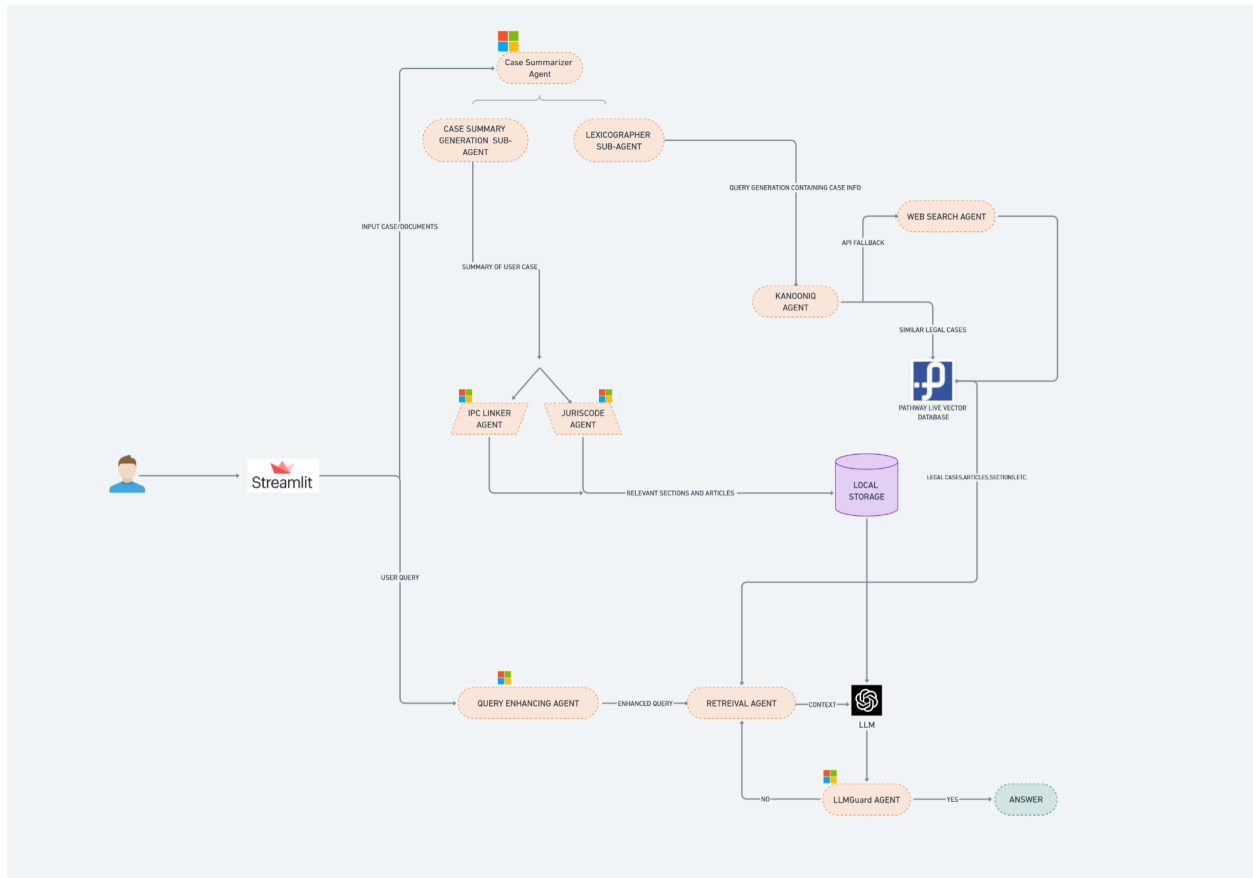
Figure 1: System Architecture

Then it provides a detailed summary to both the IPC Linker Agent and JurisCode Agent for further retrieval. This will enable IPC Linker Agent and JurisCode Agent to reach a well-rounded summary which is far better than individual retrieval for fetching applicable legal sections and constitutional articles.

- KanoonIQ Agent: The KanoonIQ Agent heps in matching the user's case with existing legal precedents, ensuring that the advice and information provided are grounded in relevant case law. This agent interfaces with the Indian Kanoon API to perform accurate and efficient mappings. It uses the keyword provided by the Case Summarizer Agent's Keyword Extraction Sub-Agent as input and queries the Indian Kanoon API, to identify matching or similar legal precedents. It then gets the necessary and relevant case information, consisting of judgments, legal reasoning, and results. It stores the similar legal cases on Pathway's Vector Store and create a repository of analogous legal cases that could be referenced for the present legal situation to improve the accuracy and reliability of the system's responses.

- Web-Search AGENT: The Web Search Agent is designed to enhance the Indian Kannon Agent by integrating error resilience and optimizing query handling. It utilizes a robust web search function that processes user queries, retrieving both recent legal cases and historical cases (from 5-10 years ago) to provide comprehensive, structured output. The agent ensures that all data is relevant and formatted correctly for further analysis. By leveraging an assistant and system framework, it adapts to varied query scenarios, improving the robustness and accuracy of case retrieval, even in the face of system

errors or ambiguous inputs.

- **IPC Linker Agent:** The IPC Linker Agent specializes in identifying and retrieving pertinent sections from the Indian Penal Code (IPC) that are applicable to the user's legal case. This ensures that the user will get information about provisions directly applicable in law to their situation. It works by extracting the specific IPC sections that correspond with the particulars of the user's case, as determined by the case summary provided by the Case Summarizer Agent's Case Summary Generation Sub-Agent, looks at the full summary of the case to determine its legal context and issues, searches IPC for relevant sections specific to the situation of the case. It then gets the text of the found sections, with their relevant subsections or clauses and saves the retrieved IPC section in the local storage of the system for accessible and safe reference. This will help the user get a clearer understanding of the statutory framework which applies to their case.

- **JurisCode Agent:** The JurisCode Agent helps in identifying and retrieving applicable articles from the Constitution that relate to the user's legal case. This ensures that constitutional provisions are considered alongside statutory laws to provide a comprehensive legal perspective. It analyzes the detailed summary provided by the case summarization agent to understand the constitutional issues involved, searches the Constitution for articles corresponding to the identified legal themes and concerns of the case, extracts the full text of the relevant constitutional articles and stores the retrieved constitutional articles in the system's local storage, thereby ensuring that the retrieved constitutional articles will be securely stored in the local storage of the system, making them easily available for referencing and analysis.

- **Pathway Live Vector Database:** The application integrates a live vector store using Pathway by loading raw data (`binary`) with metadata via `pw.io.fs.read`. Data is parsed with `ParseUnstructured` and split into chunks using `RecursiveCharacterTextSplitter`. These chunks are embedded into dense vectors using `SentenceTransformerEmbedder` (`"all-MiniLM-L6-v2"`), representing their semantic meaning. The embeddings are stored in a KNN-based dynamic index for similarity-based retrieval, and Pathway updates the vector store in real-time with new data or embeddings as needed.

- **Query Enhancing Agent:** The Query Enhancing Agent helps in refining and optimizing the user's initial query to ensure that the data retrieval process is both accurate and contextually relevant. This optimized and enhanced query helps retrieve far better information by capturing the nuances and specifics of the user's legal case. The user query is routed through sophisticated algorithms and templates for rephrasing and expansion with relevant legal context and information from Case Summarizer Agent to ensure that the enhanced query incorporates the specific characteristics of the user's case to generate an optimized query to be sent to the Retriever Agent for retrieval. The quality of queries thereby improves the chances of retrieving high-quality legal information through specific queries.

- **Retriever Agent:** The Retriever agent helps in getting all the most relevant legal precedents and information based on the enhanced queries provided by the Query Enhancing Agent. It uses Pathway's live VectorStore for legal precedents that match closely to the enhanced query based on their semantic similarity and relevance and selects the top-ranking legal cases and documents that best correspond to the user's query. It then gathers additional relevant information from local storage, including IPC sections and constitutional articles retrieved by the IPC Linker Agent and JurisCode Agent and builds the retrievable precedents along with the contextual information into a consolidated dataset and provides the aggregated data to the Large Language Model for final response generation. This process provides the LLM with extensive amplifying context of statutory laws, legal precedents and relevant past legal cases for more informed and nuanced answers.

- **Large Language Model:** We use the GPT3.5T model to generate our final answer. It receives all the context it needs from the retriever agent and local storage and uses it as the base to generate the output.

- **LLMGuard Agent:** The LLMGuard Agent serves as a safeguard against inaccuracies and hallucinations within the response output of the large language model. The agent makes available the much-needed credibility and factual authenticity in the response to the users. This agent reads the answer produced by the Large Language Model, utilizes advanced algorithms and predetermined standards of identification of inconsistencies and unsupported claims in the answer to compare it to the pre-retrieved data and its legal references for factual accuracy. If no error is detected, it shows the response to the user. But if any hallucinations or inaccuracies are present, the agent flags the response for revision, and the cycle resets to the Retriever Agent for another data-retrieval and response process.

## 6 RESULTS AND METRICS

Our legal RAG model delivers highly reliable results in retrieving and analyzing legal precedents, constitutional articles, and relevant sections to address queries related to input cases. By leveraging advanced techniques such as a dynamic vector store and multi-agent architecture, the model ensures accurate mapping of user queries to similar cases and applicable legal provisions. Its robust retrieval capabilities, combined with query enhancement and hallucination checks, provide contextually accurate and dependable outputs. This ensures comprehensive legal insights, making the system a valuable tool for efficient case preparation and analysis.

## 7 CHALLENGES FACED AND MITIGATION STRATEGIES

### 7.1 CHALLENGES FACED

Our project face the following challenges during the development of the advanced legal pipeline

- **Absence of Knowledge Graphs in Pathway's Vectorstore:** Pathway's vectorstore, known for its cutting-edge capabilities in vector-based storage and retrieval, currently lacks the support for knowledge graphs. This necessitated a shift from our initial strategy of implementing Microsoft's LightRAG, which relies on graph-based contextual data representation. Instead, we leveraged Pathway's optimized query-matching and metadata capabilities to create a seamless system that replicates the benefits of knowledge graphs while aligning with the strengths of the vectorstore.

- **Handling Large Legal Documents in Indian Contexts:** Legal cases in India often involve extensive documents filled with intricate details and precedents, presenting a challenge for efficient data extraction and embedding. Addressing this required developing methods to break down large documents into semantically coherent segments and embedding them in a way that reduced token count while preserving critical information. Additionally, handling multilingual content added complexity, requiring robust pre-processing pipelines.

- **Efficient Multi-Agent Pipeline Implementation:** Our pipeline's multi-agent design, with dedicated agents for tasks like legal clause extraction and semantic analysis, posed challenges in terms of token management and contextual consistency. Each agent's independent yet interconnected operations necessitated heavy computational resources and an advanced orchestration system to ensure minimal token usage while maintaining high accuracy.

## 7.2  MITIGATION STRATEGIES

Several innovative solutions were employed to overcome these challenges:

- **Disk Accelerated Nearest Neighbors (DiskANN)**: DiskANN mitigates the issue of embedding and efficiently utilizing large legal documents. It segments extensive data into smaller, semantically relevant chunks and performs similarity searches through a highly optimized disk-based architecture. By focusing retrieval on clustered data with similar features, DiskANN ensures rapid and precise query resolution. Using the Pathway framework, we further refined this process, enhancing both speed and scalability for handling large datasets in the Indian legal context.

- **Microsoft Guidance AI:** Guidance AI addresses the challenge of token optimization and effective contextual orchestration of multi-agent interactions. By employing policy-based execution, it dynamically allocates token budgets and prioritizes embedding tasks for each agent. Additionally, Guidance AI's innovative constraint decoding mechanisms regulate the LLM's token generation process, ensuring concise and accurate outputs. This approach not only reduces token consumption but also enhances response quality by harmonizing agent outputs into a cohesive response, making the pipeline efficient and reliable for complex legal queries.

## 8  RESILIENCE TO ERROR HANDLING

The Web Search Agent acts as a cornerstone of error resilience by augmenting the Indian Kanoon Agent's functionality. It efficiently handles ambiguous inputs and system errors by retrieving both contemporary and historical legal cases, ensuring comprehensive case analysis. Utilizing a robust web search framework, the agent validates and cross-checks retrieved data to eliminate inconsistencies, enhancing the relevance and accuracy of outputs. It dynamically structures the information for further processing, adapting to diverse query scenarios. This layered approach strengthens the system's robustness, maintaining reliable performance and providing precise legal insights even under challenging conditions or incomplete query inputs.

## 9  RESPONSIBLE AI PRACTICES

### 9.1  GUARDRAILS

While specific guardrails are not explicitly integrated into our system, the inherent design of the pipeline ensures strong safeguards against errors or misuse. Leveraging the capabilities of OpenAI's GPT-3.5 Turbo as our LLM, the system benefits from its robust content moderation, ethical query handling, and advanced safety protocols. This LLM is adept at avoiding hallucinations, maintaining the integrity of sensitive legal content, and ensuring that generated responses adhere to ethical and factual standards. Combined with the precision of Pathway's vectorstore, which prevents irrelevant or spurious data from influencing query results, the system achieves a harmonious balance between performance and safety. This partnership between state-of-the-art retrieval and generation technologies ensures that users receive accurate, contextually appropriate, and reliable outputs, solidifying our pipeline as a trusted tool for navigating complex legal landscapes.
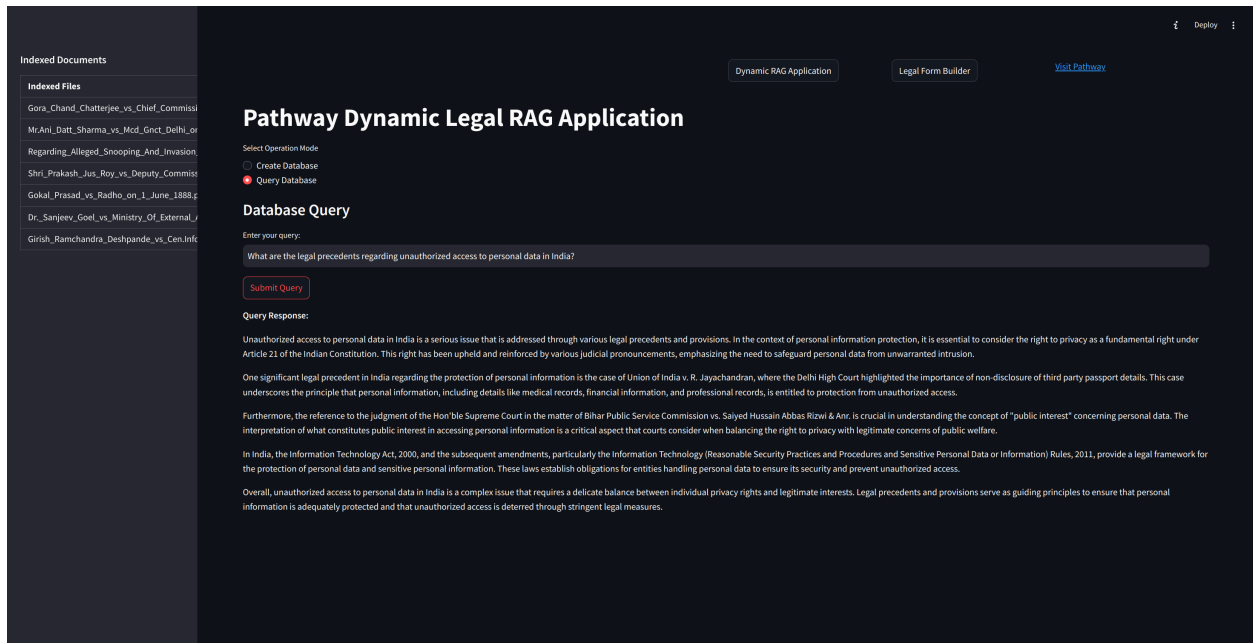
Figure 2: User Interface

# 10 USER INTERFACE

## 10.1 COMPONENTS IN THE SYSTEM

The system is designed to facilitate efficient document management and information retrieval through its key components. These components are described below:

- Header and Navigation: The header prominently displays the title of the application, *Pathway Dynamic RAG Application*, clearly indicating its purpose. Navigation options include:

  - **Legal Form Builder:** Provides access to a tool for creating legal documents.
  - **Dynamic RAG Application:** Allows seamless switching between different modules of the system.
  - **Visit Pathway:** A link to external resources or additional information related to the application.

- Indexed Documents Sidebar: This section displays a list of documents that have been indexed and are ready for querying. It is organized and scrollable, ensuring users can easily reference and identify available data sources. The documents present here are dynamically fetched from Pathway's vectorstore. The sidebar is labeled as *Indexed Files* for clarity.

- Query Interface: The core interaction area where users can perform database operations. Key features include:

  - **Operation Modes:**
    * **Create Database:** Enables users to upload or index new documents into the database.
    * **Query Database:** Allows users to retrieve information from the indexed documents.
  - **Database Query Input Field:** Provides an input area for users to enter queries in natural language or structured formats. Designed for precise query execution and information retrieval.

- **Submit Query Button:** Executes the user's input query and initiates the search process in the indexed database.

This system's design ensures a streamlined approach to document management and querying, enabling users to interact with their data efficiently while maintaining simplicity in its operation.

## Lessons Learned

1. **Importance of Contextual Understanding in Legal AI Systems:** The project highlighted the critical need for nuanced contextual understanding in the legal domain, especially to map precedents effectively and avoid shallow or irrelevant outputs.

2. **Challenges of Large-scale Data Handling:** Dealing with extensive legal documents and multilingual content demanded advanced techniques for segmenting and embedding data while maintaining accuracy and efficiency.

3. **Power of Modular Multi-agent Systems:** A modular, multi-agent design proved effective for dividing complex tasks like clause extraction, semantic analysis, and error checking, enabling targeted improvements and scalability.

4. **Value of Token Optimization:** The integration of Microsoft Guidance AI demonstrated how optimizing token usage can significantly enhance computational efficiency without sacrificing response quality.

5. **Balancing Innovation with Usability:** Ensuring a user-friendly interface while incorporating advanced technical features, such as vector-based retrieval and real-time updates, proved essential for practical application.

6. **Avoiding Hallucination in AI Responses:** Implementing safeguards, like the LLMGuard Agent, emphasized the importance of verifying AI-generated outputs to maintain trust and reliability in high-stakes applications.

7. **Future-readiness Through Scalability:** Designing the system for scalability, with plans to integrate hierarchical knowledge frameworks and constitutional indexing, reinforced the importance of anticipating future needs in system design.

## 11  CONCLUSION

This project presents a dynamic, agent-based RAG system tailored for the legal domain, offering contextual precision, robust error handling, and scalability. Core applications include analyzing legal cases with the help of legal precedents and the constitution. The system employs advanced chunking techniques for structured data processing, optimized retrieval methods, and corrective measures to tackle challenges such as hallucinations, retrieval inefficiencies, and reasoning gaps in legal research. Future enhancements aim to abstract legal concepts, build a hierarchical knowledge framework, and index constitutional clauses, further improving its interpretive and reasoning capabilities for addressing complex legal queries.

## References

[1] *CBR-RAG*: https://arxiv.org/html/2404.04302v1

[2] *Guidance-AI*: https://github.com/guidance-ai/guidance

[3] *Legal Bench RAG*: `https://arxiv.org/html/2408.10343v1`

[4] *Tree RAG*: `https://ipchimp.co.uk/2024/02/16/rag-for-legal-documents/`

[5] *DiskANN: Fast Accurate Billion-point Nearest Neighbor Search on a Single Node*: `https://suhasjs.github.io/files/diskann_neurips19.pdf`

[6] *FreshDiskANN for Similarity Search*: `https://arxiv.org/pdf/2105.09613`

[7] *SAILER: Structure-aware Pre-trained Language Model for Legal Case Retrieval*: `https://arxiv.org/abs/2304.11370`

[8] *Legal Case Document Similarity*: `https://arxiv.org/pdf/2209.12474`

[9] *Financial Report Chunking for Effective RAG*: `https://arxiv.org/pdf/2402.05131`

[10] *LawGPT 1.0: A Virtual Legal Assistant Based on GPT-3*: `https://arxiv.org/pdf/2302.05729`

[11] *Legal Case Document Summarization*: `https://arxiv.org/pdf/2210.07544`

[12] *Lawyer LLaMA: Enhancing LLMs with Legal Knowledge*: `https://arxiv.org/pdf/2305.15062`

[13] *Corrective RAG*: `https://arxiv.org/pdf/2401.15884`

[14] *Adaptive RAG*: `https://arxiv.org/abs/2403.14403`

[15] *Light RAG*: `https://arxiv.org/abs/2410.05779`