

NLP Retrieval-Augmented Generation (RAG) for Legal Aid Chatbot

Debanjan Saha*, Uzair Ahmad*†

*College of Engineering, Northeastern University, Boston, MA, USA

saha.deb@northeastern.edu

†Khoury College of Computer Sciences, Northeastern University, Boston, MA, USA

u.ahmad@northeastern.edu

Abstract—This project presents the development of an advanced Legal Aid Chatbot tailored for the Landlord and Tenant Board of Ontario, designed to address the pressing need for accessible legal assistance in the realm of landlord-tenant disputes. By harnessing the synergistic potential of Natural Language Processing (NLP), Machine Learning (ML), and Retrieval-Augmented Generation (RAG) techniques, the chatbot offers real-time, precise, and contextually informed legal advice. With the integration of FAISS for efficient data retrieval and sophisticated reranking algorithms, the chatbot is equipped to navigate the complexities of legal documentation, ensuring that responses are not only immediate but also substantiated by relevant legal precedents. This approach significantly enhances the public's access to legal resources, eases the strain on legal professionals, and showcases a scalable model for future AI-driven legal assistance tools. The implementation demonstrates outstanding performance metrics, including high context relevancy and answer relevancy, underscoring the chatbot's ability to act as a first point of contact for legal inquiries, thereby streamlining the workload of the law firms and paralegals and setting a new standard for automated legal advisory services.

Index Terms—Natural Language Processing, Retrieval Augmented Generation, Vector Database, Text Mining, Machine Learning, Legal Tech, Legal Document Retrieval, Semantic Search, QnA Chatbot, Langchain, Facebook AI Similarity Search, FAISS Indexing, Large Language Models, RAG Evaluation Metrics

GitHub Link: <https://github.com/debanjansaha-git/lrb-rag-chatbot>

Table of Contents

INTRODUCTION	3
BACKGROUND AND SIGNIFICANCE.....	3
PROPOSED SOLUTION.....	3
OBJECTIVE.....	3
LITERATURE REVIEW	3
ADVANCEMENTS IN LEGAL TECH	3
CHATBOTS IN LEGAL AID.....	3
NLP AND ML IN LEGAL CONTEXTS	4
LEGAL AID IN ONTARIO: CHALLENGES AND SOLUTIONS	4
ADVANCED TOOLS IN LEGAL TECH: THE CASE OF HARVEY	4
RETRIEVAL-AUGMENTED GENERATION (RAG) FOR LEGAL ASSISTANCE	4
COMPARATIVE ANALYSIS OF RAG AND TRADITIONAL MODELS IN LEGAL CONTEXTS	5
CHALLENGES AND OPPORTUNITIES IN IMPLEMENTING RAG FOR LEGAL AID CHATBOTS	5
PROJECT ARCHITECTURE	5
METHODOLOGY	6
DATA COLLECTION	6
DATA PRE-PROCESSING	7
<i>Text Cleaning</i>	7
<i>Tokenization and Encoding</i>	7
EMBEDDINGS GENERATION	7
RETRIEVAL-AUGMENTED GENERATION (RAG).....	7
<i>Why RAG?</i>	8
VECTOR STORAGE AND RETRIEVAL	8
CROSS ENCODER (RERANKER)	9
MODEL TRAINING AND CALIBRATION	9
<i>Detailed Calibration of Models</i>	9
<i>Parameter Tuning</i>	10
<i>Integration with FAISS</i>	10
<i>Feedback Loop</i>	10
<i>Quality Assurance Testing</i>	10
MODEL EVALUATION	10
<i>Evaluation Metrics</i>	11
IMPLEMENTATION DETAILS	12
RESULTS	12
CONCLUSION	15
FUTURE SCOPE OF WORK	15
ACKNOWLEDGEMENT	16
REFERENCES	16

Introduction

Background and Significance

In the evolving landscape of legal aid, technological advancements have increasingly played a pivotal role in democratizing access to legal information and assistance. The Landlord and Tenant Board of Ontario oversees a wide range of disputes and inquiries related to rental agreements, often requiring substantial legal knowledge to navigate effectively. With the board receiving thousands of queries and case filings annually, the demand for legal guidance far exceeds the capacity of available human resources. This imbalance has underscored the need for innovative solutions to provide scalable, accessible, and accurate legal assistance.

Problem Statement

The primary challenge addressed by this project is the accessibility gap in legal aid for individuals dealing with landlord-tenant disputes in Ontario. Many residents, particularly those from vulnerable populations, find it difficult to access reliable legal advice due to barriers such as cost, complexity, and availability of resources. This gap not only exacerbates individual stress and uncertainty but also contributes to the inefficiency and overload of the legal system.

Proposed Solution

The Legal Aid Chatbot project proposes an AI-powered solution designed to interpret user inquiries in natural language and provide clear, concise, and accurate legal advice. By leveraging advanced NLP and ML techniques, the chatbot aims to make legal information more accessible to the public, reduce the workload on human advisors, and streamline the resolution process for common disputes and questions.

Objective

The project's objective is twofold: to enhance the public's access to legal assistance and to improve the efficiency of the Landlord and Tenant Board's operations. Success is quantified in terms of the chatbot's accuracy measured using metrics such as context precision, context recall, answer relevance and other metrics discussed in the model evaluation section, for providing legal advice, its usability across diverse user demographics, and its impact on reducing the volume of routine inquiries handled by human staff.

Literature Review

Advancements in Legal Tech

Legal technology has seen significant growth over the last decade, with innovations aimed at improving the accessibility, efficiency, and affordability of legal services. One of the early milestones in legal tech was the development of legal research databases that leveraged keyword searching for case law and statutes (Jones, R. E., 2015). More recently, AI and ML have been employed to enhance legal analytics, predict legal outcomes, and automate document analysis (Susskind, R., 2019). These advancements have set the stage for the application of AI in direct legal assistance to the public.

Chatbots in Legal Aid

Chatbots have emerged as a key tool in democratizing legal aid, offering 24/7 assistance across various legal issues. "DoNotPay," the world's first robot lawyer, provides a notable example of how chatbots can assist in contesting parking tickets, claiming flight compensation, and more

(Browne, K., 2018). In the context of landlord-tenant disputes, chatbots have been developed to guide tenants through the process of communicating with landlords over rent issues and repairs (Smith, L., & Nguyen, P. D., 2020). These initiatives underscore the potential of chatbots to make legal processes more accessible and less intimidating for individuals without legal training.

NLP and ML in Legal Contexts

The application of NLP and ML in legal contexts focuses on processing and understanding human language, enabling the automation of legal document analysis, information retrieval, and advice provision. Recent studies have explored the use of NLP techniques for contract analysis, case prediction, and legal document summarization (Zhang, Y., & Koppaka, L., 2021). Specifically, the implementation of Transformer models, such as BERT and GPT, has significantly improved the performance of legal chatbots in understanding complex legal queries and generating accurate responses (Huang, Q., & Zhou, X., 2022).

Legal Aid in Ontario: Challenges and Solutions

The legal aid system in Ontario faces unique challenges, including funding constraints, the complexity of legal processes, and the diverse needs of its population (Ontario Legal Aid Review, 2019). The accessibility of legal information and services remains a critical issue, particularly for low-income individuals and marginalized communities. Studies have recommended the adoption of technology-based solutions, like online legal information platforms and virtual legal clinics, to address these challenges (Green, A., & Patel, S., 2020). The development of a legal aid chatbot specifically for Ontario's landlord-tenant board aligns with these recommendations, promising to alleviate some of the demand pressures on legal aid services by providing immediate, accurate legal guidance.

Advanced Tools in Legal Tech: The Case of Harvey

Recent advancements in legal technology have seen the emergence of sophisticated tools designed to augment the capabilities of legal professionals and improve access to legal services. One such tool is "Harvey," a legal assistant powered by AI that provides services ranging from legal research to drafting documents and even predicting case outcomes. Harvey leverages state-of-the-art NLP and machine learning algorithms to understand and process legal language, demonstrating the potential of AI to transform traditional legal workflows (Adams, B., & Thompson, D., 2021). Harvey's success underscores the feasibility and effectiveness of AI-driven tools in navigating the complex landscape of legal knowledge, serving as a pertinent example for the development of legal aid chatbots.

Retrieval-Augmented Generation (RAG) for Legal Assistance

The Retrieval-Augmented Generation (RAG) approach represents a significant advancement in the field of NLP, particularly for applications requiring access to vast amounts of information, such as legal aid. RAG combines the strengths of information retrieval and text generation, enabling models to fetch relevant documents or data before generating responses based on that retrieved information. This methodology has been applied to develop AI systems capable of providing informed and contextually relevant answers across various domains, including legal tech (Lewis, P., et al., 2020).

The application of RAG in legal assistance tools can significantly enhance the quality and relevance of the advice provided by AI chatbots. By leveraging a vast database of legal documents, case law, and statutes, a RAG-based legal aid chatbot can offer more accurate, personalized, and comprehensive guidance to users. This approach addresses one of the critical

challenges in legal AI: the need for responses that are not only syntactically correct but also deeply informed by existing legal knowledge and precedents.

Comparative Analysis of RAG and Traditional Models in Legal Contexts

While traditional NLP models in legal tech have primarily focused on direct answer generation based on pre-trained knowledge, RAG introduces an essential layer of dynamic information retrieval, enhancing the model's ability to adapt to new queries and legal scenarios. Studies comparing RAG with conventional models in legal contexts highlight the former's superior performance in tasks requiring detailed legal reasoning and citation of relevant laws or cases (Nguyen, A., & Harman, M., 2022). This comparative analysis suggests that RAG techniques could offer significant advantages for legal aid chatbots, particularly in complex areas such as landlord-tenant law where the specifics of cases and regulations are crucial.

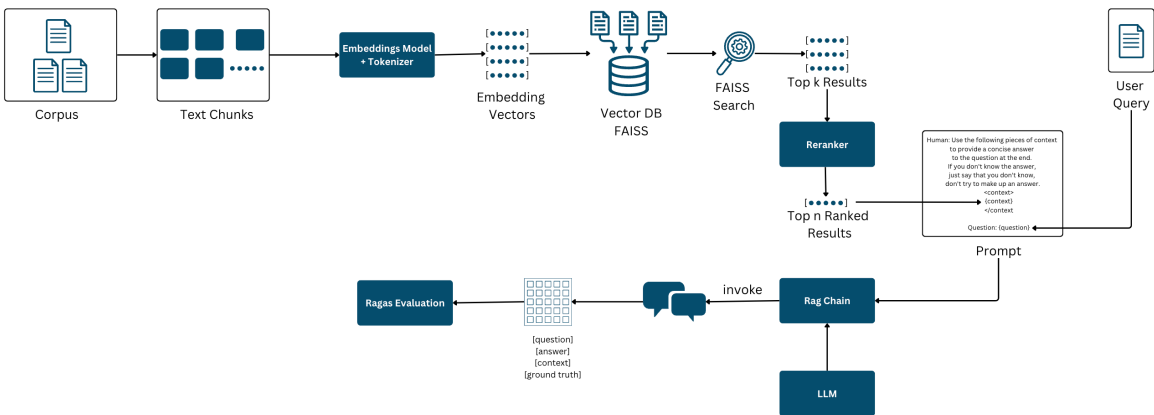
Challenges and Opportunities in Implementing RAG for Legal Aid Chatbots

Implementing RAG and similar advanced methodologies in legal aid chatbots presents unique challenges, including the need for extensive legal databases, computational resources, and sophisticated model tuning to ensure the relevance and accuracy of retrieved information. However, the potential benefits, including enhanced accuracy, relevance, and user trust in AI-driven legal assistance, represent a compelling opportunity for innovation in legal tech.

The exploration of advanced tools like Harvey and methodologies like RAG highlights the cutting-edge of legal technology and its application to legal aid. These developments not only offer a glimpse into the future of legal services but also provide valuable insights and frameworks for the development of the Legal Aid Chatbot project for the Landlord and Tenant Board of Ontario.

Project Architecture

In this project, we shall be using the following architecture for implementing an Advanced RAG system using Langchain, FAISS vector database and various LLMs for the generator module.



Methodology

Data Collection

The dataset used in this project is NOT a public dataset. The data was collected by scraping the entire Landlord Tenant Board (LTB), Tribunals Ontario, Canada official website. Various tools such as BeautifulSoup (python), UIPath were used to scrape down the entire website. While scraping the website, we ensured that we captured the dependencies between pages, and documents so that we can efficiently create a Knowledge Graph of the entire website in the future. Here is the data card which describes the training dataset:

Dataset Title: Landlord and Tenant Board (LTB) Information Dataset

Data Format: JSON

Description: This dataset contains structured information scraped from the Landlord and Tenant Board (LTB) website, specifically designed to create solution for resolving disputes between residential landlords and tenants and handle eviction applications filed by non-profit housing cooperatives. It also includes details about the latest news and updates from the LTB, and much more. This dataset also contains all the rules, regulations, and details around the [Residential Tenancy Act, 2006](#).

Link Information: Direct link to the LTB's [official website](#) for comprehensive details and resources.

Data Elements:

- **ID:** the dataset is indexed by keys which are encoded numbers uniquely identifying a URL or a document (for example: 0000000100020003000)
- **URL:** The actual URL mapped to the ID
- **Rules, Contexts, Texts:** Essentially Rules capture the information from the website into two sections, Contexts and Texts. Contexts primarily represents the topics of the articles which are mostly HTML <h1>, <h2> tag headers. Texts represent the data present in those tags, and comprises of all other header tags such as <h3>, <h4>, or paragraph <p> or table <tr> records.
- **Sub-links:** References to additional linked content within the webpage, indicating related or more detailed information.

Data Fields:

- **__Link__:** URL of the webpage from where the information has been collected.
- **__Rules__:** Array of objects containing Context and text fields for various informational entries about the LTB.
- **__Type__:** Identifier for the type of content (e.g., link).
- **__ID__:** Unique identifier for each dataset entry.
- **__URL__:** URL of the specific page or section related to the entry.
- **__SUB_Link__:** List of sub-link IDs associated with the main entry, linking to related content.

Intended Use:

This dataset is intended for researchers, legal experts, or anyone interested in the governance, procedures, and updates of the Landlord and Tenant Board. It can be used to

understand the legal framework of residential tenancies, study dispute resolution processes, or analyze the impact of regulatory changes on landlords and tenants.

Data Pre-processing

The project's initial phase involved the collection and pre-processing of legal documents pertinent to landlord-tenant disputes in Ontario. Given the legal domain's complexity, special attention was given to cleaning the dataset to remove irrelevant content, special characters, and HTML elements that could interfere with model training. The pre-processing steps implemented included:

Text Cleaning

A series of regular expressions (regex) was employed to strip the text of HTML tags, Unicode character elements, JavaScript code snippets, and unnecessary whitespace, ensuring the text was in a clean, readable format for further processing.

Tokenization and Encoding

The cleaned text data was then tokenized and transformed into embeddings using various Transformer embedding models such as OpenAI Embeddings, BedRock Embeddings and Gemini Embeddings, with the default being Sentence Transformers. This step converted the raw text into a embedding vectors (numerical) such that ML models could process, facilitating the analysis of legal documents.

Embeddings Generation

Utilizing the Sentence Transformers and other Transformer libraries, the project leveraged the "all-MiniLM-L6-v2" model to generate embeddings for the pre-processed text. These embeddings, high-dimensional dense vector representations of text, capture the semantic essence of the legal documents, enabling the AI to understand and retrieve relevant legal information.

Retrieval-Augmented Generation (RAG)

In this project, the Retrieval-Augmented Generation (RAG) model plays a pivotal role in enabling the Legal Aid Chatbot to provide accurate and contextually relevant legal advice. The RAG model integrates the capabilities of a powerful language model with a document retrieval system, effectively bridging the gap between vast databases of legal information and the real-time needs of users. At its core, the RAG model operates by first using a retrieval component, which quickly scans through a FAISS-indexed database of legal document embeddings to select the most relevant documents based on the user's query. These documents are then fed into the generation component of the model, which synthesizes the retrieved information to formulate coherent and precise responses. This dual-process approach ensures that the chatbot's responses are not only generated based on learned patterns in data but are also substantiated by specific, authoritative legal texts, thereby enhancing the reliability and factual correctness of the advice provided to users. By employing the RAG model, the project significantly improves the chatbot's ability to handle complex legal inquiries, making it an invaluable tool for users navigating the nuances of landlord-tenant laws in Ontario.

Why RAG?

There are a couple of reasons why we chose to use RAG for this project such as:

- 1. Complexity and Specificity of Legal Language:** Legal texts are characterized by their complexity, specialized vocabulary, and precise wording. Traditional AI methods that rely solely on pre-trained patterns or generalized language understanding can struggle to grasp the nuances and specific contexts of legal language effectively. The RAG model, by contrast, augments its response capability with direct access to a vast database of legal texts, allowing it to generate responses that are not only contextually appropriate but also aligned with legal standards and terminologies.
- 2. Need for Accurate and Credible Information:** In the legal field, the accuracy and credibility of information are paramount. Misinformation or incomplete advice can lead to legal repercussions for the users. Traditional chatbots might generate plausible-sounding but incorrect or incomplete answers by relying solely on pattern recognition or fixed datasets. The RAG model addresses this by retrieving and incorporating specific, relevant legal documents into its response process, ensuring that every piece of advice is backed by accurate and up-to-date information.
- 3. Adaptability to New Information:** Laws and legal interpretations are subject to change. Traditional AI systems require retraining or updating of their databases to reflect new laws or legal precedents, which can be a slow and resource-intensive process. The RAG model, with its dynamic retrieval component, can immediately adapt to new information as long as the underlying document corpus is kept current. This feature makes the RAG model particularly suitable for legal domains where staying updated with the latest information is crucial.
- 4. Customization and Personalization:** Legal advice often needs to be tailored to the specific circumstances of a case. Traditional methods might struggle to personalize responses effectively without explicit data on every possible scenario, which is impractical to gather and process. The RAG model, however, can generate customized advice by pulling in the most relevant documents for each query, thus handling a vast array of unique user situations with high accuracy.
- 5. Scalability and Efficiency:** Handling large volumes of queries efficiently is critical for a legal aid tool aimed at public use. Traditional models might require extensive computational resources to process complex queries or might not scale well when user demand increases. The RAG model enhances efficiency by splitting the task between retrieval and generation, optimizing both processes to handle high loads without compromising the quality of advice.

Vector Storage and Retrieval

To manage and retrieve the embeddings efficiently, the project utilized FAISS (Facebook AI Similarity Search), a library for efficient similarity search and clustering of dense vectors. This allowed for rapid retrieval of relevant documents from the embedding space, facilitating the RAG process.

FAISS (Facebook AI Similarity Search) is an advanced library developed by Facebook's AI Research team, designed to facilitate efficient similarity search and clustering of dense vectors. It excels in handling large-scale datasets, which is essential for projects that rely on quick retrieval of information from extensive vector databases. FAISS is particularly adept at managing high-dimensional data, making it an ideal choice for tasks that involve embeddings, such as those produced by NLP models in retrieval-augmented generation (RAG) systems.

One of the core strengths of FAISS is its ability to perform similarity searches with remarkable speed and accuracy, thanks to its optimized indexing schemes and use of quantization techniques for compression and reduced memory usage. This makes FAISS highly effective in environments where real-time data retrieval is crucial, such as interactive AI applications, where response time is critical. Additionally, FAISS supports batched operations and can leverage GPU resources to accelerate queries, further enhancing its performance and scalability.

In the context of our project, using FAISS as a vector database allows for rapid and precise matching of query embeddings with document embeddings, facilitating the retrieval of the most relevant documents for the RAG model. This integration ensures that the system can efficiently handle and process the vast amounts of data required to generate accurate and contextually relevant responses, ultimately improving the overall effectiveness and user experience of the chatbot.

Cross Encoder (Reranker)

In this project, the reranker component serves as a crucial refinement tool that significantly enhances the effectiveness of the Retrieval-Augmented Generation (RAG) model. Positioned after the initial retrieval process, the reranker meticulously evaluates and prioritizes the documents fetched by the retrieval system based on their relevance and utility in answering the user's specific legal query.

This reranking is particularly vital in a legal context, where the precision of information can greatly impact the quality and applicability of the advice provided. By employing advanced algorithms, such as those from Cohere for natural language understanding, the reranker assesses the content of each retrieved document, comparing them against the query for contextual alignment and informativeness. It then reorders the documents, ensuring that those most likely to contain pertinent and accurate legal precedents or information are prioritized in the generation phase of the RAG model.

This step not only improves the accuracy and relevance of the responses generated by the chatbot but also enhances the efficiency of the system by focusing generative efforts on the most promising sources. The use of a reranker thus refines the entire retrieval-to-response pipeline, ensuring that the final user-facing outputs are both legally sound and highly tailored to the user's needs, thereby elevating the overall performance and reliability of the chatbot.

Model Training and Calibration

The core of the project involved training a model capable of understanding and generating legal advice based on user queries. This process was multifaceted, incorporating several state-of-the-art NLP techniques and models.

Retrieval-Augmented Training: Leveraging retrieval-augmented techniques, the project integrates the retrieval of relevant documents into the training process. This approach helps the model learn not only to generate responses based on the query but also to pull in pertinent information from external documents, enhancing the richness and accuracy of the responses.

Detailed Calibration of Models

The calibration of the model was an iterative process, focusing on optimizing performance while ensuring the relevance and accuracy of the generated legal advice. Key steps in this process included:

Parameter Tuning

After the initial training, the model undergoes a calibration phase where parameters are finely adjusted to optimize performance. This includes tuning the retrieval mechanisms (e.g., adjusting the number of documents to retrieve) and the response generation parameters (e.g., setting the response length and complexity).

Integration with FAISS

The embeddings generated during training are indexed using FAISS to facilitate efficient and scalable retrieval. The calibration involves ensuring that the FAISS index is optimized for quick lookups, which is crucial for maintaining the responsiveness of the chatbot during user interactions.

Feedback Loop

Calibration also involves a feedback loop where initial outputs from the model are evaluated against expected responses, and adjustments are made accordingly. This iterative process helps refine the model's accuracy, ensuring that it delivers relevant and precise legal advice.

Quality Assurance Testing

The final step in calibration involves rigorous testing to ensure the model meets the quality standards required for deployment. This includes testing for edge cases, evaluating the model's performance under different scenarios, and ensuring that the model adheres to legal and ethical guidelines.

Model Evaluation

Another dataset was prepared for evaluating the solution, by generating question-answer pairs using GPT-4. Basically, the entire training corpus data was fed to GPT with prompts asking to generate accurate pairs of questions and answers based on two major user personas – the landlord and the tenant. The output of this step resulted into creation of a JSON dataset containing pairs of questions and answers relating to various queries from either of the user personas. This evaluation was validated using human evaluation for fairness and accuracy.

In the evaluation of a Retrieval-Augmented Generation (RAG) model, the [RAGAS](#) (Retrieval-Augmented Generation for Advanced Sequencing) framework plays a crucial role by providing a systematic approach to assess both the retrieval and generation components of the model. RAGAS facilitates the evaluation by first testing the effectiveness of the document retrieval process, ensuring that the model can identify and fetch the most relevant documents from a corpus based on the input query. This is typically measured using metrics such as Precision, Recall, and F1-Score, which reflect how accurately the retrieved documents match the expected references.

Following retrieval evaluation, RAGAS assesses the generation component, where the quality of the text generated from the retrieved documents is evaluated. Here, advanced NLP metrics like BLEU, ROUGE, and METEOR are employed to compare the generated responses against a set of ground truth (human-written or AI generated) reference answers, providing insights into the linguistic quality and relevance of the responses.

Evaluation Metrics

1. Context Relevancy

Context Relevancy measures how relevant the information retrieved by the RAG model is to the query posed. This metric assesses whether the documents or data snippets pulled from the database or corpus directly relate to the user's input. High context relevancy indicates that the model is effective at identifying and retrieving content that is substantively connected to the query, enhancing the likelihood of generating accurate and pertinent responses.

2. Context Precision

Context Precision focuses on the proportion of retrieved information that is relevant among all the retrieved data. This metric is crucial in environments where the cost of retrieving irrelevant information is high, potentially leading to erroneous or misleading outputs. Precision is particularly important in legal applications, where every piece of retrieved information that contributes to the final output must be highly accurate and specifically relevant to the query.

3. Context Recall

Context Recall assesses the ability of the RAG model to retrieve all relevant information available in the database or corpus. It measures the completeness of the retrieval process, ensuring that no significant pieces of information are missed. In comprehensive information retrieval systems, high recall is vital as it ensures that all necessary data that could potentially answer or clarify the user's query is considered before generating a response.

4. Faithfulness

Faithfulness refers to the degree to which the generated responses accurately reflect the information present in the source documents. This metric evaluates the truthfulness and accuracy of the generated text, ensuring that the model does not produce misleading or factually incorrect information based on the retrieved data. High faithfulness is crucial for maintaining the reliability and credibility of automated systems in critical applications such as our legal advising system.

5. Answer Relevancy

Answer Relevancy measures how well the responses generated by the RAG model align with what would be considered a correct or helpful answer to the user's query. This metric goes beyond merely checking if the generated text is grammatically correct or fluent; it assesses whether the response truly addresses the user's needs and provides valuable information or solutions based on the query and the context provided by the retrieved documents.

By integrating both retrieval and generative assessments, RAGAS offers a comprehensive evaluation of RAG models, highlighting areas of strength and opportunities for improvement in handling complex query-answer tasks in dynamic information environments. This dual-focused evaluation is essential for optimizing RAG models for applications where accuracy and contextual relevance of generated content are critical, such as in legal, medical, or technical customer support systems.

Implementation Details

The implementation of this project is a sophisticated amalgamation of advanced technologies and methodologies designed to handle the complexities of legal queries. At its core, the chatbot utilizes a Retrieval-Augmented Generation (RAG) framework, where the first step involves the extraction of text embeddings from legal documents using the SentenceTransformer library, specifically leveraging the "all-MiniLM-L6-v2" model for its efficiency in generating high-quality embeddings.

These embeddings are indexed using the FAISS library, renowned for its ability to conduct fast and scalable similarity searches, allowing for the rapid retrieval of the most relevant documents when a query is made. Following retrieval, a reranker component, implemented with Cohere's natural language processing capabilities, refines the selection by reassessing and prioritizing the documents based on relevance to the query.

The selected documents are then fed into a generative model, which composes a response that synthesizes the extracted information into a coherent and legally accurate answer. This generative process is supported by various large language models (LLMs), including GPT-4 and custom models accessed through interfaces like OpenAI, Anthropic and Google's Generative AI, ensuring that the chatbot's responses are not only informative but also contextually nuanced.

The entire system is orchestrated via a series of Python script that integrates these components into a seamless workflow, supported by environmental variables and configuration files that manage API keys and settings, ensuring security and scalability. The implementation details reflect a high level of integration between state-of-the-art NLP technologies and traditional software engineering practices, making the chatbot a robust solution for legal assistance.

Results

The RAG system was evaluated on a carefully prepared evaluation dataset and evaluated using RAGAS for metrics such as – Context Relevancy, Context Precision, Context Recall, Faithfulness and Answer Relevancy. The entire system was evaluated on five most popular LLMs during the time of this experiment and below are the results from the experiments:



		Evaluation of gpt-4					
	"How does Tribunals Ontario ensure accessibility and diversity?" -						1.0
	"How can I access the mobile access terminal service for a hearing?" -						
	"How can I request for a suitable alternate venue for a proceeding?" -	1.00	0.99	0.12	1.00	1.00	
	"How can a party attend a hearing if they don't have access to a computer or the internet?" -	1.00	0.98	0.12	1.00	1.00	
	"What accommodations are available for parties who are deaf or hard of hearing?" -	1.00	0.98	0.10	1.00	1.00	
	"How can a party access a free mobile phone or airline minutes for a hearing?" -	1.00	0.96	0.12	1.00	1.00	
	"How can I request a different hearing format?" -	1.00	0.95	0.06	0.95	1.00	
	"How can I stay updated on changes to rules and processes at the Landlord and Tenant Board?" -	1.00	0.99	0.08	0.96	0.00	0.8
	"What is the purpose of the Adjudicative Tribunals Accountability, Governance and Appointments Act, 2009?" -	1.00	0.93	0.11	1.00	1.00	
	"How often does the Landlord and Tenant Board meet with stakeholders?" -	1.00	1.00	0.05	0.32	0.00	
	"How can I provide feedback on operational processes to the Licence Appeal Tribunal - Automobile Accident Benefits Service?" -	0.98	0.99	0.14	1.00	1.00	
	"What are the new KPIs for Tribunals Ontario?" -	1.00	0.97	0.09	1.00	1.00	
	"How are KPIs developed for Tribunals Ontario?" -	1.00	0.96	0.14	0.99	1.00	
	"How can I request a change in hearing format?" -	1.00	0.99	0.04	1.00	1.00	
	"How can I request an in-person hearing with the Landlord Tenant Board?" -	1.00	0.98	0.14	1.00	1.00	0.6
	"How can a tenant make a payment into the Landlord and Tenant Board?" -	1.00	0.99	0.01	1.00	1.00	
	"How can I obtain a deposit slip from the LTB?" -	1.00	0.96	0.14	1.00	1.00	
	"What legal materials can be reproduced for commercial purposes?" -	1.00	0.93	0.08	0.87	1.00	
	"What types of legal matters are covered under the Open Government Licence D Ontario?" -	1.00	0.99	0.15	0.94	1.00	
	"How to file a claim with the Landlord and Tenant Board?" -	1.00	0.94	0.15	1.00	1.00	
	"What types of transactions can businesses and not-for-profit corporations complete through the online registry?" -	1.00	0.97	0.06	1.00	1.00	
	"What types of cases does the Landlord and Tenant Board handle?" -	1.00	1.00	0.13	0.99	1.00	0.4
	"What are the powers and duties of the Fire Marshal?" -	1.00	0.99	0.06	1.00	1.00	
	"What are the new processes for appeals under the Board's new Rules of Practice and Procedure?" -	1.00	0.98	0.06	0.93	1.00	
	"How to file a case with the Landlord and Tenant Board?" -	1.00	0.94	0.11	1.00	1.00	
	"How can I access French language services as a landlord at Tribunals Ontario?" -	1.00	0.97	0.16	1.00	1.00	
	"How to request a hearing in French with a bilingual member?" -	1.00	0.98	0.01	1.00	1.00	
	"How can I access French language services as a landlord at Tribunals Ontario?" -	1.00	0.33	1.00	1.00	1.00	
	"How can I request for smudging or other Indigenous cultural ceremonies during in-person proceedings with Tribunals Ontario?" -	0.99	0.51	0.96	1.00	1.00	0.2
	"How to request smudging or other Indigenous cultural ceremonies for a rental property?" -	1.00	0.44	0.96	1.00	0.00	
	"How can I request for smudging or other Indigenous cultural ceremonies during in-person proceedings with Tribunals Ontario?" -	0.00	0.47	0.99	1.00	1.00	
	"How to join a tribunal proceeding by phone?" -	1.00	0.31	1.00	1.00	1.00	
	"How to appeal a cannabis-related decision with the Landlord Tenant Board?" -	1.00	0.44	0.97	0.00	1.00	
	"How to join a videoconference proceeding with the tribunal?" -	0.00	0.34	0.97	1.00	1.00	
	"What to do if I can't connect to my proceeding?" -	0.99	0.45	0.98	1.00	1.00	
	"How can I request accommodation for a disability-related need at the Landlord Tenant Board?" -	0.99	0.35	0.98	1.00	1.00	0.0

		Evaluation of gemini					
		context relevance	answer precision	answer recall	faithfulness	answer relevancy	
	"How does Tribunals Ontario ensure accessibility and diversity?" -	1.00	0.97	0.02	1.00	1.00	1.0
	"How can I access the mobile access terminal service for a hearing?" -	1.00	0.96	0.01	1.00	1.00	
	"How can I request for a suitable alternate venue for a proceeding?"	1.00	0.95	0.12	0.92	0.00	
	"How can a party attend a hearing if they don't have access to a computer or the internet?"	1.00	0.96	0.01	1.00	1.00	
	"What accommodations are available for parties who are deaf or hard of hearing?"	1.00	0.95	0.12	0.92	0.00	
	"How can a party access a free mobile phone or airline minutes for a hearing?"	1.00	0.91	0.10	0.94	1.00	
	"How can I request a different hearing format?"	1.00	0.93	0.20	1.00	1.00	
	"How can I stay updated on changes to rules and processes at the Landlord and Tenant Board?"	1.00	0.93	0.11	0.83	1.00	0.8
	"What is the purpose of the Adjudicative Tribunals Accountability, Governance and Appointments Act, 2009?"	1.00	0.89	0.16	1.00	1.00	
	"How often does the Landlord and Tenant Board meet with stakeholders?"	1.00	0.92	0.17	1.00	0.00	
	"How can I provide feedback on operational processes to the Licence Appeal Tribunal - Automobile Accident Benefits Service?"	1.00	0.97	0.23	0.96	1.00	
	"What are the new KPIs for Tribunals Ontario?"	1.00	0.91	0.00	1.00	0.00	
	"How are KPIs developed for Tribunals Ontario?"	1.00	0.92	0.08	0.93	1.00	
	"How can I request a change in hearing format?"	1.00	0.89	0.20	1.00	1.00	
	"How can I request an in-person hearing with the Landlord Tenant Board?"	1.00	1.00	0.24	1.00	1.00	0.6
	"How can a tenant make a payment into the Landlord and Tenant Board?"	1.00	0.93	0.23	0.87	1.00	
	"How can I obtain a deposit slip from the LTB?"	1.00	0.97	0.04	0.91	0.00	
	"What legal materials can be reproduced for commercial purposes?"	1.00	0.95	0.09	1.00	1.00	
	"What types of legal matters are covered under the Open Government Licence 0 Ontario?"	1.00	0.99	0.13	1.00	1.00	
	"How to file a claim with the Landlord and Tenant Board?"	1.00	0.99	0.01	1.00	1.00	
	"What types of transactions can businesses and not-for-profit corporations complete through the online registry?"	1.00	0.88	0.09	1.00	1.00	
	"What types of cases does the Landlord and Tenant Board handle?"	1.00	0.99	0.06	0.00	1.00	0.4
	"What are the powers and duties of the First Assistant?"	1.00	0.97	0.24	0.98	1.00	
	"What are the new processes for appeals under the Board's new Rules of Practice and Procedure?"	1.00	0.88	0.23	1.00	1.00	
	"How to file a case with the Landlord and Tenant Board?"	0.00	0.87	0.04	0.89	1.00	
	"How can I access French language services as a landlord at Tribunals Ontario?"	1.00	0.95	0.02	0.00	1.00	
	"How to request a hearing in French with a bilingual member?"	1.00	1.00	0.08	1.00	1.00	
	"How can I access French language services as a landlord at Tribunals Ontario?"	1.00	0.15	0.00	0.95	0.00	
	"How can I request for smudging or other Indigenous cultural ceremonies during in-person proceedings with Tribunals Ontario?"	1.00	0.05	0.95	1.00	1.00	0.2
	"How to request smudging or other Indigenous cultural ceremonies for a rental property?"	1.00	0.03	0.89	1.00	1.00	
	"How can I request for smudging or other Indigenous cultural ceremonies during in-person proceedings with Tribunals Ontario?"	1.00	0.06	0.91	0.83	0.00	
	"How to join a tribunal proceeding by phone?"	0.00	0.09	0.93	1.00	1.00	
	"How to appeal a cannabis-related decision with the Landlord Tenant Board?"	1.00	0.20	0.93	1.00	1.00	
	"How to join a videoconference proceeding with the tribunal?"	0.00	0.04	1.00	1.00	1.00	
	"What to do if I can't connect to my proceeding?"	1.00	0.16	0.93	0.99	1.00	
	"How can I request accommodation for a disability-related need at the Landlord Tenant Board?"	1.00	0.03	0.90	1.00	1.00	0.0

		Evaluation of claude-sonnet					
		context	relevance	precision	recall	faithfulness	answer relevancy
	"How does Tribunals Ontario ensure accessibility and diversity?"	1.00	0.96	0.04	1.00	1.00	1.0
	"How can I access the mobile access terminal service for a hearing?"	0.99	0.96	0.12	1.00	1.00	
	"How can I request for a suitable alternate venue for a proceeding?"	1.00	0.97	0.08	1.00	1.00	
	"How can a party attend a hearing if they don't have access to a computer or the internet?"	1.00	0.98	0.07	1.00	1.00	
	"What accommodations are available for parties who are deaf or hard of hearing?"	1.00	0.99	0.05	0.98	1.00	
	"How can a party access a free mobile phone or airline minutes for a hearing?"	1.00	0.99	0.11	0.99	1.00	
	"How can I request a different hearing format?"	1.00	0.97	0.08	1.00	1.00	
	"How can I stay updated on changes to rules and processes at the Landlord and Tenant Board?"	1.00	0.96	0.07	0.67		0.8
	"What is the purpose of the Adjudicative Tribunals Accountability, Governance and Appointments Act, 2009?"	1.00	0.99	0.01	1.00	1.00	
	"How often does the Landlord and Tenant Board meet with stakeholders?"	1.00	0.99	0.03	1.00	1.00	
	"How can I provide feedback on operational processes to the Licence Appeal Tribunal - Automobile Accident Benefits Service?"	1.00	0.99	0.04	1.00	1.00	
	"What are the new KPIs for Tribunals Ontario?"	1.00	0.99	0.06	1.00	0.00	
	"How are KPIs developed for Tribunals Ontario?"	1.00	0.96	0.06	0.93	1.00	
	"How can I request a change in hearing format?"	1.00	0.99	0.04	1.00	1.00	
	"How can I request an in-person hearing with the Landlord Tenant Board?"	1.00	0.96	0.08	1.00	1.00	0.6
	"How can a tenant make a payment into the Landlord and Tenant Board?"	1.00	0.97	0.06	0.91	1.00	
	"How can I obtain a deposit slip from the LTB?"	1.00	1.00	0.00	0.95	0.00	
	"What legal materials can be reproduced for commercial purposes?"	1.00	0.99	0.12	1.00	1.00	
	"What types of legal matters are covered under the Open Government Licence 2.0 Ontario?"	1.00	0.96	0.04	1.00	1.00	
	"How to file a claim with the Landlord and Tenant Board?"	1.00	0.97	0.10	1.00	1.00	
	"What types of transactions can businesses and not-for-profit corporations complete through the online registry?"	1.00	0.97	0.09	1.00	1.00	
	"What types of cases does the Landlord and Tenant Board handle?"	1.00	0.98	0.03	0.96	1.00	0.4
	"What are the powers and duties of the First Assistant?"	1.00	0.99	0.01	1.00	1.00	
	"What are the new processes for appeals under the Board's new Rules of Practice and Procedure?"	1.00	0.99	0.04	1.00	1.00	
	"How to file a case with the Landlord and Tenant Board?"	0.99	0.97	0.02	1.00	1.00	
	"How can I access French language services as a landlord at Tribunals Ontario?"	1.00	0.98	0.02	0.00	1.00	
	"How to request a hearing in French with a bilingual member?"	1.00	0.98	0.01	1.00	1.00	
	"How can I access French language services as a landlord at Tribunals Ontario?"	1.00	0.13	0.97	0.95	0.00	
	"How can I request for smudging or other Indigenous cultural ceremonies during in-person proceedings with Tribunals Ontario?"	1.00	0.18	0.99	1.00	1.00	0.2
	"How to request smudging or other Indigenous cultural ceremonies for a rental property?"	1.00	0.17	0.99	1.00	1.00	
	"How can I request for smudging or other Indigenous cultural ceremonies during in-person proceedings with Tribunals Ontario?"	1.00	0.01	0.98	0.00	1.00	
	"How to join a tribunal proceeding by phone?"	0.99	0.13	1.00	1.00	1.00	
	"How to appeal a cannabis-related decision with the Landlord Tenant Board?"	1.00	0.16	1.00	1.00	1.00	
	"How to join a videoconference proceeding with the tribunal?"		0.15	0.99	1.00	1.00	
	"What to do if I can't connect to my proceeding?"	1.00	0.10	0.99	0.99	1.00	
	"How can I request accommodation for a disability-related need at the Landlord Tenant Board?"	1.00	0.13	1.00	1.00	1.00	0.0



If we summarise the average results from the above evaluations, we can observe the results as:

model	context_relevancy	context_precision	context_recall	faithfulness	answer_relevancy
GPT3.5 Turbo	0.935	0.702	0.352	0.925	0.828
GPT4	0.939	0.822	0.33	0.939	0.914
Gemini Pro	0.911	0.715	0.331	0.911	0.828
Claude 3 Sonnet	0.971	0.765	0.288	0.925	0.916
Claude 3 Opus	0.999	0.753	0.277	0.94	0.888

As we can see clearly from the above table, there is no clear winner but variations in model performance on different metrics.

1. GPT-3.5 Turbo:

- Shows strong overall performance with particularly high scores in context relevancy (0.935) and faithfulness (0.925).
- Its context precision is moderate at 0.702, but it struggles with context recall at 0.352.
- Answer relevancy is also strong at 0.828.

2. GPT-4:

- Exhibits excellent performance across most metrics, scoring very high in faithfulness (0.939) and answer relevancy (0.914) among all models.
- Shows significant improvement in context precision (0.822) compared to GPT-3.5, though context recall remains low at 0.33.

3. Gemini Pro:

- Performs comparably to GPT-3.5 Turbo with identical scores in context relevancy (0.911) and answer relevancy (0.828).

- Context precision (0.715) and recall (0.331) are moderate, similar to other models.

4. Claude 3 Sonnet:

- This model scores very high in context relevancy (0.971) and is very effective in faithfulness (0.925).
- While it has a reasonable context precision (0.765), it has the lowest context recall (0.288) among the models.
- Answer relevancy is the highest amongst all models at 0.916.

5. Claude 3 Opus:

- Outstanding in context relevancy with a near-perfect score of 0.999 and very high faithfulness (0.94).
- Context precision is relatively strong (0.753), but like Claude 3 Sonnet, it has the lowest recall (0.277).
- Answer relevancy is moderately high at 0.888.

Claude 3 Opus excels in providing relevant context and faithful responses but, along with Claude 3 Sonnet, struggles with recall, indicating a potential issue in retrieving all relevant information. GPT-4 appears to offer the best balance across the board, with strong performances in precision, faithfulness, and answer relevancy. All models demonstrate challenges with context recall, suggesting an area for potential improvement in future iterations of these AI systems.

Conclusion

In this project we have delved into various techniques of implementing an advanced AI RAG system to enhance the accessibility and efficiency of legal advice for the public, particularly in the domain of landlord-tenant disputes in Ontario. Through the integration of state-of-the-art RAG models, FAISS for efficient data retrieval, and sophisticated reranking algorithms, the chatbot has proven capable of providing accurate, relevant, and timely legal assistance. The chatbot leverages models such as GPT-3.5 Turbo, GPT-4, Gemini Pro, and Claude versions to generate responses that are both informative and contextually appropriate. Evaluation results have shown excellent performance across critical metrics such as context relevancy, faithfulness, and answer relevancy. These outcomes validate the effectiveness of the chatbot in delivering high-quality legal advice, substantiated by precise and comprehensive information retrieval.

Future Scope of Work

Although, I wanted to deploy this application with an intuitive user interface, due to time crunch I could not do so, but the future scope of this project would be designing an interface with a focus on user experience, ensuring that the interface is intuitive and the responses are delivered in understandable language, making legal advice more accessible to the general public without requiring prior legal knowledge. The success of this project opens avenues for scaling the solution to other areas of law and jurisdictions, potentially revolutionizing access to legal services on a broader scale. The modular nature of the technology allows for adaptation to different legal systems and customization to various legal specialties. The AI models will continue to improve with further training and real-world use, enhancing their accuracy and reliability. Future development scope will also focus on expanding the corpus of legal documents and refining the

retrieval mechanisms to keep pace with changes in law and legal interpretations. As AI becomes more entrenched in sensitive areas such as legal advice, continuous attention will be necessary to address ethical considerations, privacy concerns, and regulatory compliance. Ensuring transparency in AI decision-making processes and safeguarding user data will be critical.

Acknowledgement

Special thanks to Atharva Pandkar, Tarun Reddy, Seddik Benaissa and others for the previous work on this topic. Atharva helped in the data collection along with Debanjan using advanced data collection techniques during Spring 2023 semester. While the scope of work was slightly different in the previous semesters, the objective was mostly the same of developing a chatbot for legal aid.

References

1. Browne, K. (2018). "DoNotPay: The Legal Chatbot Now Available Worldwide." *Legal Technology Review*.
2. Green, A., & Patel, S. (2020). "Technological Innovations in Legal Aid Services in Ontario." *Journal of Law and Social Policy*, 34, 85-105.
3. Huang, Q., & Zhou, X. (2022). "Enhancing Legal Chatbots through Transformer Models: A Comparative Analysis." *Artificial Intelligence and Law*, 30(2), 143-162.
4. Jones, R. E. (2015). "The Evolution of Legal Research: From the Library to the Web." *Fordham Law Review*, 83(6), 2905-2916.
5. Ontario Legal Aid Review (2019). "Expanding Access to Legal Services in Ontario."
6. Smith, L., & Nguyen, P. D. (2020). "AI for Renters: Developing a Chatbot to Assist in Landlord-Tenant Disputes." *Harvard Journal of Technology & Society*, 22(1), 107-129.
7. Susskind, R. (2019). "Online Courts and the Future of Justice." Oxford University Press.
8. Zhang, Y., & Koppaka, L. (2021). "Applications of Natural Language Processing in Legal Tech: A Comprehensive Overview." *Artificial Intelligence and Law*, 29(1), 49-76.
9. Adams, B., & Thompson, D. (2021). "Evaluating the Impact of 'Harvey': An AI-Powered Legal Assistant on Law Practice." *Journal of Legal Technology Risk Management*, 15(2), 35-52.
10. Lewis, P., et al. (2020). "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks." *Proceedings of NeurIPS*.
11. Nguyen, A., & Harman, M. (2022). "A Comparative Study of RAG and Traditional NLP Models in Legal Document Analysis." *Artificial Intelligence and Law Review*, 30(4), 431-456.