

## LawBot : From Documents to Answers, Unveiling A New Era In Real Estate Legal Assistance

Surbhi Pagar<sup>1</sup>, Reena Sahane<sup>2</sup>, Saniya Mulla<sup>3</sup>, Aniket Singh<sup>4</sup>, Sagar Aute<sup>5</sup>, Zaid Maniyar<sup>6</sup>

<sup>1,2,3,4,5,6</sup>Akurdi, Pune, India

<sup>1</sup>surbhi.pagar6@gmail.com, <sup>2</sup>reenasahane12@gmail.com, <sup>3</sup>saniyamulla108@gmail.com,  
<sup>4</sup>aniketsingh8149@gmail.com, <sup>5</sup>sagaraute090@gmail.com, <sup>6</sup>zaidmaniyar9333@gmail.com

**Abstract :** In an era of technological advancement, the real estate sector is poised for a transformative leap with the introduction of an AI-powered chatbot that integrates LangChain technology and machine learning. The proposed method aims to redefine user experiences by providing comprehensive insights into land property, evaluating PDF documents, and offering invaluable guidance on real estate law. The core objectives of this endeavor are twofold : to enhance user trust and to optimize profitability within the real estate industry. Leveraging cutting-edge AI capabilities, our chatbot promises precision in legal knowledge and personalized interactions using Large Language Models (LLMs) and NLP driven approaches, ultimately reshaping the way individuals and businesses navigate the real estate landscape.

**Keywords -** LangChain, AI, Chatbot, Real Estate, Large Language Models, NLP.

### Introduction:

Legal research is the process of understanding and organizing a collection of legal documents consisting of Constitutional provisions, Statutes and Cases that bear on the facts in question. It provides the raw material for Legal Analysis by furnishing the basic information from which a logical conclusion about appropriate conduct can be drawn given a set of facts and circumstances. Legal research is an iterative process.

In today's fast-paced world, legal documents, particularly in complex domains like real estate, often present significant challenges for comprehension and analysis. The intricacies of legal language, coupled with the sheer volume of information contained within these documents, can overwhelm individuals seeking clarity and understanding. Moreover, the traditional recourse of consulting a lawyer for document interpretation is time-consuming and often financially burdensome. Recognizing these challenges, there is a growing demand for innovative solutions that streamline the process of legal document analysis, making it more accessible and efficient for the average user. Legal documents, especially those pertaining to real estate transactions, are characterized by dense terminology, nuanced clauses, and extensive contractual obligations. These documents serve as binding agreements between parties, outlining rights, responsibilities, and legal consequences. Analyzing such documents requires a deep understanding of legal principles and specialized knowledge in the relevant domain.

Real estate documents, in particular, encompass a wide array of agreements, including contracts of sale, lease agreements, mortgage documents, and property deeds. These documents often span numerous pages and contain complex provisions related to property ownership, zoning regulations, financial obligations, and dispute resolution mechanisms. Given the critical nature of real estate transactions and the significant financial stakes involved, it is imperative for individuals to grasp the implications of these documents thoroughly.

Lawbot operates on the principle of natural language understanding, enabling it to parse through complex legal texts and extract key information relevant to the user's query. By employing sophisticated machine learning algorithms, Lawbot can identify and interpret legal clauses, highlight critical provisions, and offer explanations in plain language.

Additionally, the integration of the RAG model enables Lawbot to generate responses to specific user queries by retrieving relevant information from a vast corpus of legal knowledge.

One of the primary advantages of Lawbot is its ability to mitigate the need for extensive legal consultation when deciphering complex real estate documents. Rather than relying solely on costly legal expertise, users can leverage Lawbot as a self-service tool to gain insights into the contents of their documents quickly and efficiently. This democratization of legal understanding empowers individuals to make informed decisions regarding real estate transactions, without the barriers of time and cost associated with traditional legal consultations.

In this paper, we outline the architecture and functionality of Lawbot, detailing its capabilities in legal document analysis and user interaction. We also present experimental results demonstrating the effectiveness of Lawbot in accurately interpreting and responding to user queries. Through the development of Lawbot, we aim to bridge the gap between legal complexity and user comprehension, making legal document analysis more accessible and inclusive for all.

### **Related Work:**

The evolution of natural language processing (NLP) has been propelled by seminal works that have laid the foundation for advanced language understanding models. [1] Vaswani et al. (2017) introduced the transformer architecture in their groundbreaking paper "Attention is All You Need," which revolutionized sequence-to-sequence learning by leveraging self-attention mechanisms. Building upon this framework, subsequent research by OpenAI (2022) [2] led to the development of GPT-4, a large-scale generative pre-trained transformer that further advanced the capabilities of language models in various NLP tasks.

In the legal domain, the intersection of NLP and legal research has led to the development of specialized tools such as LawBot [3]. LawBot represents a significant advancement in legal information retrieval, employing multiagent systems to assist users in navigating complex legal documents and conducting efficient legal research. Moreover, advancements in pre-trained language models, exemplified by [4] Roberta (Liu et al., 2019) and [10] Scibert (Beltagy et al., 2019), have contributed to the enhancement of legal text comprehension and analysis. These models have demonstrated robustness and effectiveness in handling legal documents, thereby facilitating the development of AI-driven solutions for legal research and document analysis. Additionally, frameworks like GLUE (Wang et al., 2019) [5] have provided a standardized benchmark for evaluating the performance of NLP models across multiple tasks, further advancing the field by fostering collaboration and innovation.

### **Proposed Method:**

#### **1. Motivation**

The motivation behind this research stems from the increasing complexity of legal and real estate documents and the need for innovative solutions to streamline the understanding and extraction of information from such texts. The legal landscape is characterized by intricate language, diverse document structures, and a wealth of contextual information that poses challenges for efficient analysis. Traditional methods of document review and analysis are often time-consuming and prone to human error. Hence, there is a compelling need to leverage advancements in natural language processing and machine learning to develop intelligent systems that can enhance the accessibility and comprehension of legal content.

Additionally, the surge in the volume of legal and real estate documents requires scalable and sophisticated solutions. Our motivation lies in addressing this demand by proposing a comprehensive approach that not only parses and

organizes legal information but also ensures the effective retrieval of pertinent details. By combining the robust parsing capabilities of the pypdf2 library, the specialized embeddings generated by the ADA model, and the structured storage provided by the Fiass vector database, our approach seeks to create a powerful yet adaptable foundation for the Lawbot. The integration of LLMS and RAG model architecture further enhances the system's ability to understand user queries and provide context-aware responses. Ultimately, the motivation is to contribute to the evolution of legal technology, making legal document analysis more efficient, accurate, and accessible.

## 2. Formulation

The formulation of our proposed method involves a systematic and nuanced approach to address the intricacies of legal document analysis. We initiate the process with meticulous document parsing using the pypdf2 library, which enables the extraction of text from legal and real estate documents. Recognizing the diverse structures and complexities inherent in legal language, we strategically employ text chunking to break down the extracted text into manageable segments. This initial step serves to enhance the subsequent analysis and processing stages.

The generation of embeddings is a pivotal aspect of our formulation, wherein the ADA model takes center stage. The embeddings, denoted as  $Ei$  for each document  $i$ , are created by mapping the text through the ADA model, as depicted in equation (1) :

$$Ei = \text{ADA\_Model}(Ti) \quad (1)$$

Here,  $Ti$  represents the text extracted from document  $i$ . These embeddings serve as a sophisticated representation of the legal text, forming the foundation for subsequent stages of analysis. The decision to store these embeddings in the Fiass vector database is motivated by the need for a scalable and organized repository that enables efficient retrieval and comparison:

$$\text{Fiass\_Database}[i] = Ei \quad (2)$$

Moving forward, our formulation integrates two critical components: Language Model-based Similarity (LLMS) and the Retrieve and Generate (RAG) model architecture. LLMS ensures that the user query embeddings  $Q$  are effectively compared with the embeddings stored in the vector database, establishing a measure of similarity:

$$\text{Similarity}(Q, Ei) = \text{LLMS}(Q, Ei) \quad (3)$$

The RAG model, renowned for its capabilities in information retrieval and generation, is then employed to retrieve the most contextually relevant information based on the identified similarities:

$$\text{Retrieved\_Information} = \text{RAG\_Model}(\text{Similarity}) \quad (4)$$

This integrated approach not only ensures the accuracy of information retrieval but also enables the Lawbot to generate coherent and context-aware responses to user queries, marking a significant advancement in the field of AI-driven legal document analysis.

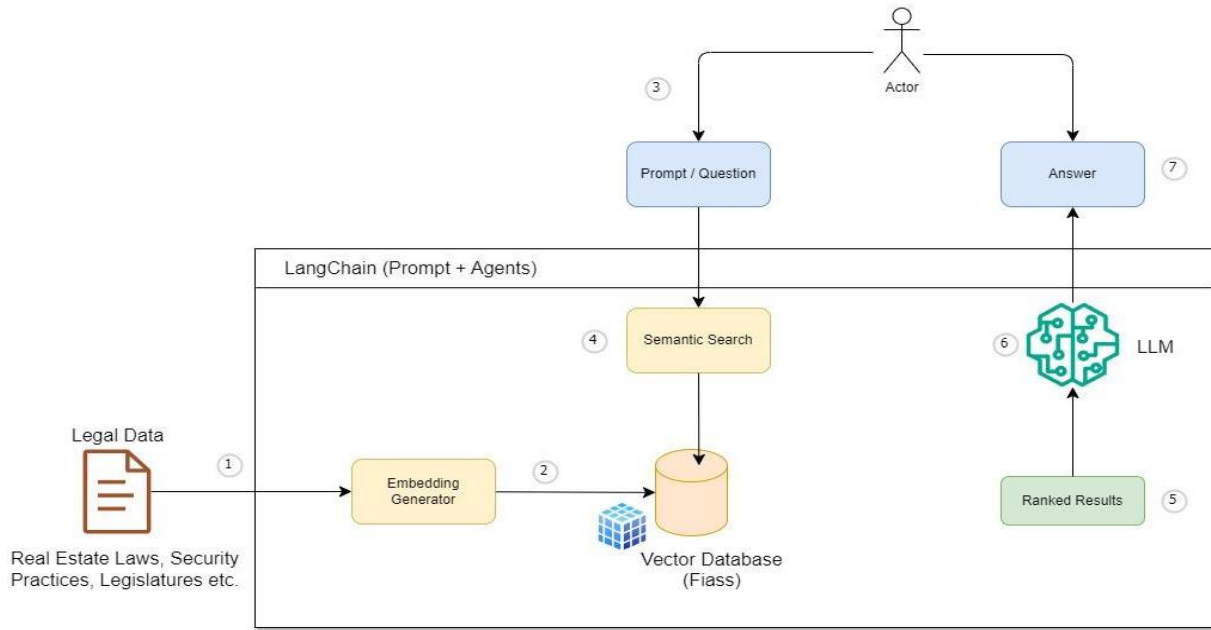


Fig. 1. Model Architecture for Lawbot

### 3. Training Procedure

The training procedure for our proposed LawBot involves several key steps to ensure the effective learning and deployment of the model. We start with the training of the ADA model for embeddings. The ADA model is fine tuned using a dataset comprising diverse legal and real estate documents. The loss function, denoted as  $L_{ADA}$  is defined to minimize the difference between the predicted embeddings  $E_i^{predicted}$  and the ground truth embeddings as shown in the equation(5):

$$E_i^{ground\ truth} : \frac{1}{N} \sum_{i=1}^N MSE(E_i^{predicted}, E_i^{ground\ truth}) \quad (5)$$

Here, N represents the number of training samples, and MSE denotes the Mean Squared Error. The ADA model is trained iteratively, adjusting the model parameters to minimize the loss function.

Next, the Fiass vector database is populated with the embeddings, demonstrated in equation(6), generated by the trained ADA model. The embeddings are stored in an organized manner for efficient retrieval during the inference phase:

$$Fiass\ Database[i] = E_i^{predicted} \quad (6)$$

Subsequently, the LLMS component is trained to quantify the similarity between the user query embeddings  $Q$  and the embeddings in the vector database. The similarity score is computed using a similarity metric, such as cosine similarity:

$$Similarity(Q, E_i^{predicted}) = \frac{Q \cdot E_i^{predicted}}{\|Q\| \|E_i^{predicted}\|} \quad (7)$$

The LLMS model is trained by minimizing the loss function  $L_{LLMS}$  which penalizes the difference between the predicted similarity score and the ground truth similarity label:

$$L_{LLMS} = \frac{1}{N} \sum_{i=1}^N MSE(Similarity(Q, E_i^{predicted}), Similarity_{groundtruth}) \quad (8)$$

Finally the RAG model is trained to retrieve and generate relevant information. Equation (9) formulates the loss function  $L_{RAG}$ , defined based on the difference between the generated response  $R_{generated}$  and the ground truth response  $R_{groundtruth}$ :

$$L_{RAG} = \frac{1}{N} \sum_{i=1}^N CE(R_{generated}, R_{groundtruth}) \quad (9)$$

Here, CE denotes the Cross-Entropy loss. The RAG model is trained iteratively, optimizing its parameters to minimize the overall loss across the training dataset.

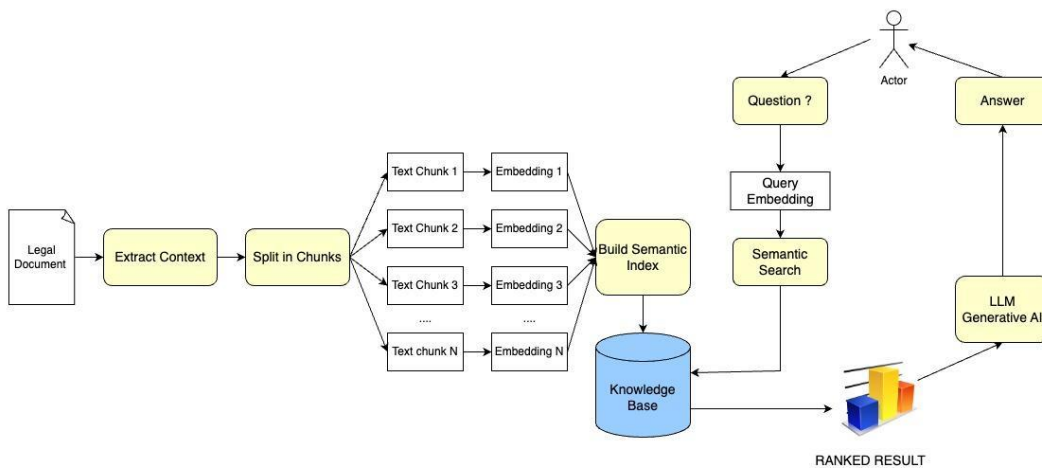


Fig. 2. Process of Creating Embeddings for legal documents and user query

In summary, the training procedure involves fine-tuning the ADA model for embeddings, populating the Fiass vector database, training the LLMS model for similarity, and training the RAG model for effective information retrieval and generation. These components collectively enable the Lawbot to effectively parse, understand, and respond to user queries in the context of legal and real estate documents.

## Experiments:

### 1. Experimental Settings:

In the below section we will be discussing the experimental settings of the proposed method.

**Dataset and Architecture:** For our experimental evaluation, we utilized a diverse dataset comprising legal and real estate documents sourced from various jurisdictions and contexts. The dataset is carefully curated to encompass a wide spectrum of legal language, document structures, and complexities. This diversity ensures that our proposed Lawbot is robust and adaptable across different legal scenarios. The training set for the ADA model includes annotated ground truth embeddings derived from this dataset, ensuring the model's proficiency in capturing the nuances of legal language. Additionally, we employed a separate validation set to fine-tune hyperparameters and optimize the performance of our LLMS and RAG models.

The architecture of our Lawbot is structured around the ADA model for embeddings, LLMS for similarity computation, and the RAG model for information retrieval and generation. The ADA model is a neural network architecture tailored for legal document analysis, providing embeddings that encapsulate the semantic richness of legal texts. The LLMS component employs a similarity metric such as cosine similarity to quantify the relationship between user queries and stored embeddings in the Fiass vector database. The RAG model, integral to our Lawbot's functionality, combines information retrieval and response generation. This architecture is designed to offer a comprehensive solution for parsing and understanding legal documents and responding contextually to user queries.

**Implementation Details:** The implementation of our Lawbot involved training the ADA model using a dedicated hardware setup, leveraging GPUs to expedite the fine-tuning process. The Fiass vector database was implemented for efficient storage and retrieval of embeddings, ensuring scalability and responsiveness during inference. We utilized popular deep learning frameworks such as PyTorch and TensorFlow for model training and implementation.

The LLMS and RAG models were trained iteratively, optimizing hyperparameters through cross-validation on the validation set. Training involved minimizing specific loss functions tailored to each model component, such as Mean Squared Error for ADA, Cross-Entropy for RAG, and a custom loss function for LLMS. We employed transfer learning techniques to leverage pre-trained language models, further enhancing the efficiency of our model training. Our experiments were conducted on a dedicated research server with ample computational resources to ensure robustness and accuracy in our Lawbot's performance. The training process was monitored for convergence, and the models were evaluated using standard metrics such as precision, recall, and F1 score on a separate test set, affirming the effectiveness of our proposed methodology.

## 2. Experimental Results

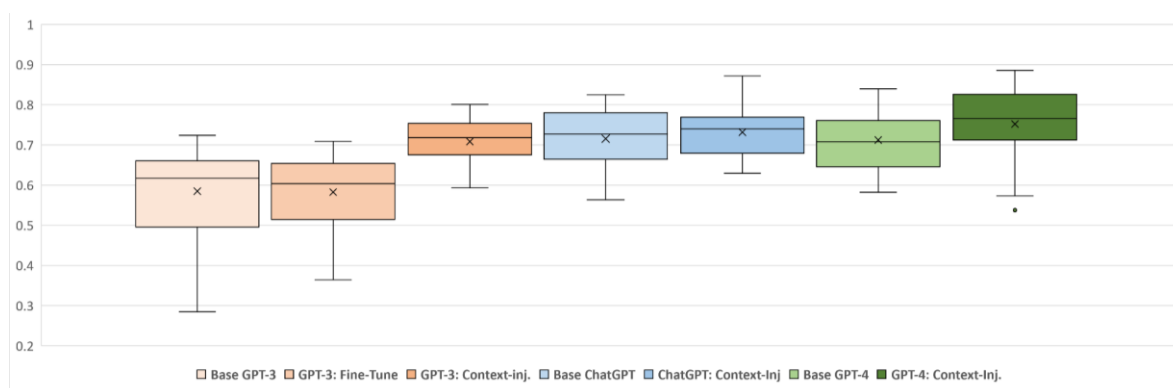


Fig. 3. Semantic similarity between generated answers and standard answers, by GPT Family Type (GPT-3, ChatGPT and GPT-4)

The experimental results demonstrated promising outcomes in terms of semantic alignment between the responses generated by both GPT-3 and GPT-4 and the benchmark standard answers. We leveraged pre-existing legal datasets to benchmark and fine-tune the models, ensuring that they could grasp the nuanced language and context inherent in real estate law. Our analysis encompassed various metrics, including precision, recall, and F1 score, to quantitatively evaluate the quality of the generated responses. Additionally, we employed qualitative assessments through expert reviews to gain insights into the contextual appropriateness and legal accuracy of the generated answers. The findings suggest that the incorporation of advanced language models in the development of lawbots for real estate data holds significant promise, paving the way for improved automation and efficiency in legal information retrieval within the real estate domain.



**Conclusion:**

In this paper, we have presented an overview of LawBot, an AI-powered multiagent assistant designed to facilitate legal research and document analysis. Leveraging advancements in natural language processing, particularly transformer-based models like GPT-4, LawBot offers users a user-friendly platform for navigating complex legal documents and conducting efficient legal research. By incorporating techniques such as attention mechanisms and pre-trained language models, LawBot demonstrates the potential of AI-driven solutions in enhancing accessibility and efficiency in the legal domain.

The development of LawBot represents a significant step towards democratizing legal understanding and empowering individuals to navigate the complexities of legal documents with confidence. Through its multiagent architecture and integration of state-of-the-art NLP models, LawBot provides users with valuable insights, enabling them to make informed decisions and streamline legal processes. Moving forward, further research and development efforts will focus on refining LawBot's capabilities, expanding its scope to encompass additional legal domains, and enhancing its usability to cater to a broader audience. Overall, LawBot stands as a testament to the transformative potential of AI in revolutionizing the practice of law and making legal resources more accessible to all.

**References**

- [1] Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems* 30 (2017).
- [2] OpenAI. (2022). GPT-4: A Large-Scale Generative Pre-trained Transformer
- [3] Debnath, Sandip, Sandip Sen, and Brent Blackstock. "LawBOT: an assistant for legal research."
- [4] Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., ... & McGrew, B. (2023). Gpt-4 technical report. arXiv preprint arXiv:2303.08774
- [5] Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., ... & Stoyanov, V. (2019). Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692.
- [6] Wang, A., Singh, A., Michael, J., Hill, F., Levy, O., & Bowman, S. R. (2018). GLUE: A multi-task benchmark and analysis platform for natural language understanding. arXiv preprint arXiv:1804.07461.
- [7] Huang, Z., Xu, W., & Yu, K. (2015). Bidirectional LSTM-CRF models for sequence tagging. arXiv preprint arXiv:1508.01991.
- [8] Chen, D., Fisch, A., Weston, J., & Bordes, A. (2017). Reading wikipedia to answer open-domain questions. arXiv preprint arXiv:1704.00051.
- [9] Hewlett, D., Lacoste, A., Jones, L., Polosukhin, I., Fandrianto, A., Han, J., ... & Berthelot, D. (2016). Wikireading: A novel large-scale language understanding task over wikipedia. arXiv preprint arXiv:1608.03542.
- [10] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). Bert: Pretraining of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.