



A Comprehensive Tool for Legal Document Interpretation and Summarization using Large Language Models

1stVeena Gode Swamy Rao 1*, 2ndSuhas Katrahalli 2, 3rd Dhruthi Bhat 3, 4th Tanya Arora

¹Department of Computer Science and Engineering, Ramaiah Institute of Technology, Bangalore, India

²Department of Computer Science and Engineering, Ramaiah Institute of Technology, Bangalore, India

³Department of Computer Science and Engineering, Ramaiah Institute of Technology, Bangalore, India

⁴ Department of Computer Science and Engineering, Ramaiah Institute of Technology, Bangalore, India

*Corresponding Author: veenags@msrit.edu

Citation: Veena Gode Swamy Rao and Suhas Katrahalli, Dhruthi Bhat, Tanya Arora "A Comprehensive Tool for Legal Document Interpretation and Summarization using Large Language Models," *International Journal of Communication Networks and Information Security (IJCNIS)*, vol. 16, no.4, 2024, pp.818-824.

ARTICLE INFO

Received: 17 Aug 2024

Accepted: 24 Sep 2024

ABSTRACT

The proposed system in this paper introduces a user-friendly software solution leveraging cutting-edge AI technology called Large Language Models (LLMs) to simplify the understanding of legal documents and ensure fairness within the legal system. With LLMs at its core, the system offers two primary functions. Firstly, users can upload various legal documents, such as contracts or statutes, and ask questions related to their content. Using sophisticated natural language processing techniques, the system analyses these documents and provides accurate answers, aiding both legal professionals and individuals without legal expertise in navigating complex legal texts effortlessly.

By harnessing the power of LLMs, this software revolutionises how we interact with legal documents. Its advanced capabilities enable users to better understand legal papers and ensure they're fair and transparent. With its user-friendly interface and focus on leveraging LLM technology, the system aims to empower users to make informed decisions and promote fairness and accountability within the legal domain.

Keywords: – LLMs, NLP, Legal Documents, Prompt Chaining, Collaborative Benchmarking

INTRODUCTION

In the proposed work, we address a prevalent issue under scored by compelling data: many individuals enter into legal agreements without a thorough comprehension of their terms. Research indicates that this trends pans across various contexts, encompassing scenarios such as business negotiations and tenancy agreements.

The consequences of this lack of understanding manifold, ranging from potential disputes to perpetuating disparities in legal literacy. Individuals without a legal background often find themselves disadvantaged in contractual negotiations, further exacerbating existing societal inequalities.

Our work aims to remedy this situation by leveraging sophisticated technology and linguistic strategies to simplify legal documents. By enhancing accessibility and promoting comprehension, we aspire to empower individuals to make informed decisions and advocate for their rights, thereby fostering fairness and equity in legal interactions. Ultimately, our objective is to bridge the gap between legal expertise and everyday understanding, facilitating a more equitable distribution of legal knowledge and resources.

question-answer based system which can solve the queries of users

Fig1:FlowchartofLegal Document Interpretation and Summarizationsystem

A. *ExtractingTextfromDocument*

The first and foremost thing in order to achieve. The text is to be extracted from the document which will be in pdf format. The extracted text will be the base for all the further operation to be done.

B. *SplitTextintoChunks*

The major advantage of this system is the ability to parse through huge documents. The extracted text will be divided into smaller chunks with a buffer. We are using a map-reduce method to make sure we can parse huge amounts of text through any LLM.

C. *EmbeddingText*

The text is converted to embeddings which is basically a vector representation of that particular text. It is done in order to store it in a vector database. There are two methods presented in this system, one of them is using HuggingFace Embeddings which is Open Source.

The other method is using Open AI Embeddings which is proprietary.

D. *StoringtheTextChunksasEmbeddings*

The system uses a Vector Database to store the text chunks for faster retrieval. Similarity search is carried out to find the chunks which have the same context as the question being asked.

E. *PassingtheprompttoLLM*

A chain is used for this to implement the conversation feature. There are two LLMs presented in this system, one of them is using Llama which is Open Source and one more is ChatGPT API which is proprietary.

F. *MemoryPersistence*

The chat history is stored in the memory to make sure it is taken into account for the upcoming prompts.

G. *CreatingaGUI*

Designing and implementing a GUI with pdf upload and chat input options so that the user can chat with the document.

H. *IntegratingLangSmith*

The project is integrated with LangSmith. LangSmith traces contain the full information of all the inputs and outputs of each step of the application. It helps quickly debug, test, and continuously improve the application.

IMPLEMENTATION

ToolsusedforComputation

Programming Language: The whole system was built and achieved using Python language. Python brings an exceptional amount of power and versatility to machine learning environments. The language's simple syntax simplifies data validation and streamlines the scraping, processing, refining, cleaning, arranging and analysing processes, thereby making collaboration with other programmers less of an obstacle.

A. *TechnologiesandToolsused*

HuggingFace is now widely associated with state-of-the-art developments in natural language processing (NLP). They were established in 2016 with the goal of democratising AI by giving developers, academics, and companies easy access to tools and information. The Hugging Face Transformers library, which offers cutting-edge pre-trained models for a variety of NLP applications, including sentiment analysis and language translation, is the foundation of their services. HuggingFace not only contributes to open-source projects but also provides enterprise solutions for businesses wishing to use natural language processing (NLP) in their goods and services. Hugging Face is a dynamic community of users and contributors that keeps pushing the limits of AI-driven language understanding.

Streamlit is a technology used for building web applications and data dashboards in Python. It is a popular open-source tool that simplifies the process of creating interactive data-driven applications. With Streamlit, developers can easily create interactive web-based data visualisations without the need for extensive web development knowledge. Streamlit allows developers to create responsive and interactive data applications with simple Python scripts. Streamlit also offers a range of features that make it easy to deploy and share your applications.

Llama was released in February 2023, Meta released a set of open-source artificial intelligence models known as **LLAMA**, which stands for "Large Language Model MetaAI". These models are intended for use in code authoring, translation, and text production operations. In contrast to certain rivals, **LLAMA** prioritises efficiency, attaining robust outcomes with a reduced quantity of parameters in comparison to analogous models. Originally available in four sizes (7, 13, 33, and 65 billion parameters), the 13B version even outperformed the OpenAI GPT-3, which was far larger, on a number of benchmarks. This emphasis on accessibility is maintained in the most recent version, **LLAMA 3**, which is freely available for both commercial and research use. Because of its open methodology, **LLAMA** is now a well-liked tool among developers and researchers that are working on cutting-edge AI.

Langchain is a framework designed to simplify building applications powered by large language models (LLMs) like OpenAI's GPT-4 or Meta's **LLAMA**. Launched in October 2022, **Langchain** provides a toolbox for developers, regardless of expertise level. Imagine Lego blocks for AI development – **Langchain** offers pre-built components like prompts, data retrieval functions, and decision-making algorithms. These can be snapped together to create complex applications. This modular approach allows developers to focus on the unique functionalities of their app without getting bogged down in the underlying LLM complexities. **Langchain** even offers pre-built "chains" for common tasks like chatbots and question answering, further accelerating development. With its open-source nature and focus on accessibility, **Langchain** is fostering a growing community of developers pushing the boundaries of what's possible with LLMs.

ChromaDB ChromaDB is not like other databases. It is an open-source vector storage created with large language models (LLMs) in mind. **Chroma DB**, which was introduced in 2023, is dedicated to storing and retrieving vector embeddings, a unique type of data. These numerical representations of language or code, known as embeddings, help LLMs comprehend the links and significance of data.

The document is uploaded by the user which is in pdf format. The text from the document is then extracted using **PyPDF2** library.

```
def get_pdf_text(legal_docs):
    text = ""
    for pdf in legal_docs:
        legal_doc_reader = PdfReader(pdf)
        for page in legal_doc_reader.pages:
            text += page.extract_text()
    return text
```

The text extracted from the document is then split into chunks using **CharacterTextSplitter** function from **Langchain** library. The chunk size and overlap is set according to requirement.

```
def get_text_chunks(raw_text):
    text_splitter = CharacterTextSplitter(
        chunk_size=1000, chunk_overlap=200, length_function=len
    )
    chunks = text_splitter.split_text(raw_text)
    return chunks
```

The text chunks are then converted to embeddings and stored to a vector database, **ChromaDB** is used for this system.

```
def get_vectorstore(text_chunks):
    embedding_function = HuggingFaceEmbeddings()
    db = Chroma.from_texts(text_chunks, embedding_function)
    return db
```

The question from the user is then taken as input and then passed to the LLM as prompt with the context of the document. The chat history is stored in the memory using the function **ConversationBufferMemory**. A chain is used for the conversation in which the LLM, retriever and memory is passed as parameters. The system works with different types of documents which are uploaded as input. Here are a few snapshots of the system working with different types of documents (lease agreement, internship agreement). The users can interact with the document and ask for queries regarding the document. This can help the users' understanding of the legal document which will help them make better decisions.

RESULTS AND DISCUSSION

The system works with different types of documents which are uploaded as input. Here are a few snapshots of the system working with different types of documents (lease agreement, internship agreement). The users can interact with the document and ask for queries regarding the document. This can help the users' understanding of the legal document which will help them make better decisions. A snapshot of the couple of questions that are put forward and the response has been reflected in Figure 2, Figure 3, Figure 4, and Figure 5.

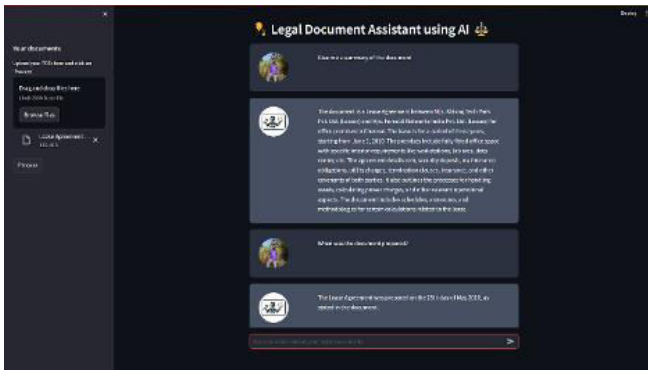


Fig2:LeaseAgreement1

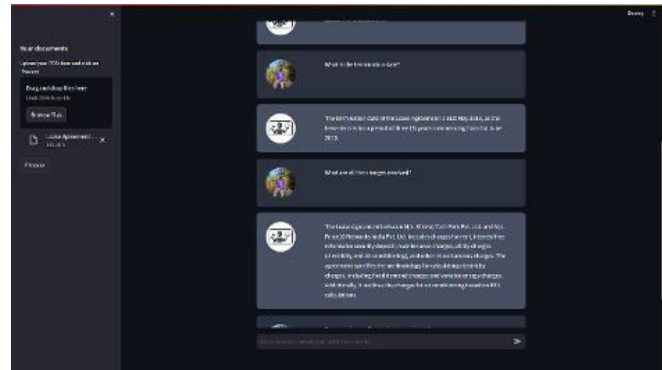


Fig3:LeaseAgreement2

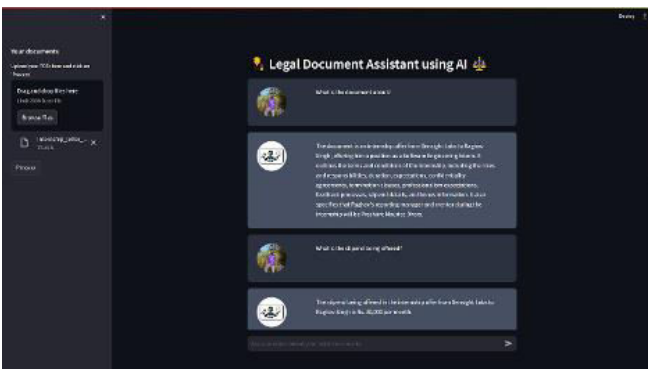


Fig4:InternshipAgreement1

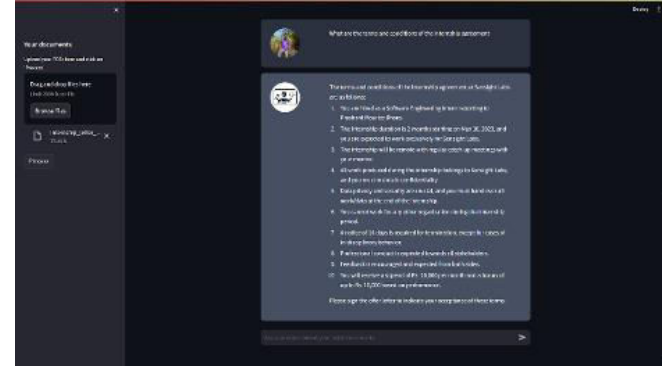


Fig5:InternshipAgreement2

CONCLUSION

This work in the development of the Legal Document Interpretation and Summarization serves a major purpose of helping a non-advocate person (with no background of law) understand a legal document and get the queries resolved. It is made sure that Open Source tools are used in the project so that there is no cost barrier for any user. Although the speed would not be the fastest, it will be enough for normal use. However, the journey doesn't end here; there exist avenues for future work and enhancement:

- Further refinement of the analysis algorithms can enhance the accuracy and efficiency of the system.
- Integration of machine learning techniques can enable the system to adapt and improve its analysis capabilities over time.
- Expanding the scope to include support for multiple languages and legal jurisdictions can broaden its applicability and reach.
- Continuous monitoring and updating of the system to stay abreast of evolving legal standards and document formats will be paramount.

Through these avenues of future work, the project can continue to evolve and advance, catering to the ever-changing needs of legal professionals and stakeholders.

ETHICAL DECLARATION

Conflict of interest: No declaration required. **Financing:** No reporting required. **Peer review:** Double anonymous peer review.

REFERENCES

- [1] Yi Luo, Xiaowei Xu, Comparative study of deep learning models for analysing online restaurant reviews in the era

- of the COVID-19 pandemic, *International Journal of Hospitality Management*, Volume 94, 2021, 102849, <https://doi.org/10.1016/j.ijhm.2020.102849>
- [2] M., M. ., ShivaKumar, S. ., J., T. ., & K. R., V. . (2023). Restaurant Based Emotion Detection Of Images From Social Media Sites Using Deep Learning Model. *International Journal of Intelligent Systems and Applications in Engineering*, 11(10s), 267–276. <https://ijisae.org/index.php/IJISAE/article/view/3250>
- [3] Hossain, Eftekhari & Sharif, Omar & Hoque, Moshiri & Sarker, Iqbal. (2021). SentiLSTM: A Deep Learning Approach for Sentiment Analysis of Restaurant Reviews. https://doi.org/10.1007/978-3-030-73050-5_19
- [4] Tanbin Siddique Eidul, Md.Alim Imran, & Amit Kumar Das. (2022). Restaurant Review Prediction using Machine Learning and Neural Network. *International Journal of Innovative Science and Research Technology*, 7(3), 1388–1392. <https://doi.org/10.5281/zenodo.6486696>
- [5] Bhoite, Sachin & Kulkarni, Atharva & Bhandari, Divya. (2019). Restaurants Rating Prediction using Machine Learning Algorithms. *International Journal of Computer Applications Technology and Research*. 8. 377-378. <https://doi.org/10.7753/IJCATRO809.1008>.
- [6] Shina, Sharma, S. & Singha ,A. (2018). A study of tree based machine learning Machine Learning Techniques for Restaurant review. 2018 4th International Conference on Computing on Computing Communication and Automation (ICCCA). <https://doi.org/10.1109/CCAA.2018.8777649>
- [7] Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, Ł.; and Polosukhin, I. 2017b. Attention is All You Need. *Advances in Neural Information Processing Systems*, 30.
- [8] Bernard J. Jansen, Soon-gyo Jung, Joni Salminen, Employing large language models in survey research, *Natural Language Processing Journal*, Volume 4, 2023, 100020, ISSN 2949-7191, <https://doi.org/10.1016/j.nlp.2023.100020>.
- [9] Zan, D., Chen, B., Zhang, F., Lu, D., Wu, B., Guan, B., Wang, Y., & Lou, J. (2023). Large Language Models Meet NL2Code: A Survey. <https://doi.org/10.18653/v1/2023.acl-long.411>
- [10] Wang, Y., Zhong, W., Li, L., Mi, F., Zeng, X., Huang, W., Shang, L., Jiang, X., & Liu, Q. (2023, July 24). Aligning Large Language Models with Human: A Survey. *ArXiv.org*. <https://doi.org/10.48550/arXiv.2307.12966>
- [11] Zhu, Y., Yuan, H., Wang, S., Liu, J., Liu, W., Deng, C., Dou, Z., & Wen, J.-R. (2023, August 15). Large Language Models for Information Retrieval: A Survey. *ArXiv.org*. <https://doi.org/10.48550/arXiv.2308.07107>
- [12] Dai, Y. S., Feng, D., Huang, J., Jia, H., Xie, Q., Han, M., Han, W., Tian, W., & Wang, H. (2023). LAiW: A Chinese Legal Large Language Models Benchmark (A Technical Report). *ArXiv (Cornell University)*. <https://doi.org/10.48550/arxiv.2310.05620>
- [13] Brynjolfsson, E., Li, D., & Raymond, L. R. (2023, April 1). Generative AI at Work. *National Bureau of Economic Research*. <https://www.nber.org/papers/w31161>
- [14] Cohen, M.C., Dahan, S., Khern-am-nuai, W. et al. The use of AI in legal systems: determining independent contractor vs. employee status.
- [15] Arbel, Y. A., & Becher, S. (2023). How Smart are Smart Readers? LLMs and the Future of the No-Reading Problem. *Social Science Research Network*. <https://doi.org/10.2139/ssrn.4491043>
- [16] Kabir, Md Shahin & Nazmul Alam, Mohammad. (2023). IoT, Big Data and AI Applications in the Law Enforcement and Legal System: A Review. 10. 1777-1789.
- [17] Myers, D., Mohawesh, R., Chellaboina, V.I. et al. Foundation and large language models: fundamentals, challenges, opportunities, and social impacts *Cluster Comput*(2023). <https://doi.org/10.1007/s10586-023-042037>
- [18] Demszky, D., Yang, D., Yeager, D.S. et al. Using large language models in psychology. *Nat Rev Psychol* 2, 688–701 (2023). <https://doi.org/10.1038/s44159-023-00241-5>
- [19] Agrawal, M., Hegselmann, S., Lang, H., Kim, Y., & Sontag, D. (2022, December 1). Large language models are few-shot clinical information extractors. *ACLWeb; Association for Computational Linguistics*. <https://doi.org/10.18653/v1/2022.emnlp-main.130>
- [20] X. Jing and S. Xiuhui, "Artificial Intelligence Legal Consulting Based on Categorization," 2020 International Conference on Artificial Intelligence and Education (ICAIE), Tianjin, China, 2020, pp. 148-151, doi: 10.1109/ICAIE50891.2020.00041.
- [21] T. McKeown, J. Mustafina, R. Magizov and C. Gataullina, "AI in Law Practices," 2020 13th International Conference on Developments in eSystems Engineering (DeSE), Liverpool, United Kingdom, 2020, pp. 27-32, doi: 10.1109/DeSE51703.2020.9450780.
- [22] C. Mahoney, P. Gronvall, N. Huber-Fliflet and J. Zhang, "Explainable Text Classification Techniques in Legal Document Review: Locating Rationales without Using Human Annotated Training Text Snippets," 2022 IEEE International Conference on Big Data (Big Data), Osaka, Japan, 2022, pp.2044-2051, doi:10.1109/BigData55660.2022.10020626.
- [23] R. Saha and S. Jyhne, "Interpretable Text Classification in Legal Contract Documents using Tsetlin Machines,"

2022 International Symposium on the Tsetlin Machine (ISTM), Grimstad, Norway, 2022, pp. 7-12, doi: 10.1109/ISTM54910.2022.00011.