# Advancements in Legal Document Processing and Summarization

Varun Kusabi[1*],  Purav Jain[1],  Riya Palesha[1],
Krisha Khandelwal[1],  Dr. G. P. Potdar[1],  P. T. Kohok[1]

[1]Computer Engineering, Pune Institute of Computer Technology, Pune, 411043, Maharshtra, India.

*Corresponding author(s). E-mail(s): varunkusabi8@gmail.com;
Contributing authors: puravjain1@gmail.com; riyapalesha07@gmail.com;
1942krisha@gmail.com; gppotdar@pict.edu; ptkohok@pict.edu;

**Abstract**

Efficient analysis of legal documents is essential for legal proceedings, requiring accurate extraction and interpretation of information from intricate textual materials. Traditional methods often involve manual review by legal experts, which is time-consuming and prone to errors. In this paper, we propose an innovative approach leveraging artificial intelligence (AI) and machine learning (ML) techniques for automated legal document analysis.

Our method involves extracting text from legal documents and preprocessing the data for analysis. We employ a specialized text model trained on legal text to comprehend the content and structure of the documents.

The system integrates advanced natural language processing (NLP) techniques to understand user queries and retrieve pertinent information from the document repository.

We assess the system's performance across a varied array of legal documents and showcase its efficacy in providing precise answers to user queries. Our findings underscore that our approach offers a swifter and more dependable alternative to manual document analysis, thereby enhancing efficiency and productivity in legal research and practice.

In essence, our research contributes to the advancement of AI and ML applications in the legal realm, furnishing practical solutions for automated legal document analysis and information retrieval.

**Keywords:** Legal Document Analysis, Artificial Intelligence, Machine Learning, Natural Language Processing, Text Extraction

1

# 1  Introduction

Legal document analysis is a critical aspect of legal proceedings, essential for extracting and interpreting information vital to various legal matters. However, conventional methods often involve manual review by legal professionals, which is time-consuming and prone to inaccuracies. In many cases, the availability of such expertise is limited, leading to delays and inefficiencies in legal processes.

To address these challenges, this paper proposes an innovative approach leveraging artificial intelligence (AI) and machine learning (ML) techniques for automated legal document analysis. By harnessing advanced technologies, our aim is to streamline the process of extracting key information from legal documents, thereby improving efficiency and accuracy in legal research and practice.

Our method involves extracting text from legal documents and applying sophisticated AI models trained on legal text to comprehend the content and structure of the documents. This enables us to identify crucial details such as the parties involved, contractual terms, important dates, and other pertinent information.

Moreover, we develop a user-friendly interface where users can input questions related to the legal document, and our AI system provides accurate answers based on the extracted data. This facilitates quick and precise retrieval of relevant information, enhancing the overall effectiveness of legal document analysis.

By automating the analysis of legal documents, our approach aims to reduce reliance on manual review and expedite the legal research process. We believe that our innovative solution will significantly benefit legal professionals, researchers, and stakeholders involved in legal proceedings, ultimately contributing to greater efficiency and effectiveness in the legal domain.

# 2  Technology

In the realm of legal document analysis, machine learning (ML) algorithms play a pivotal role in automating the extraction and interpretation of information from complex legal texts. Supervised learning techniques, such as classification and regression algorithms, are utilized to predict outcomes and categorize legal documents based on their content and context. Support vector machines (SVM), decision trees, and random forests are among the versatile algorithms employed for classification tasks in legal document analysis.

Unsupervised learning methods, including clustering algorithms like k-means and hierarchical clustering, aid in identifying patterns and structures within legal documents, enabling the categorization of documents and the discovery of similarities and differences among them. Dimensionality reduction techniques, such as principal component analysis (PCA), are applied to streamline the analysis process and extract meaningful features from large datasets of legal documents.

Deep learning algorithms, such as convolutional neural networks (CNNs) and recurrent neural networks (RNNs), offer powerful capabilities for analyzing textual data in legal documents. These algorithms excel in tasks such as named entity recognition, document classification, and summarization, contributing to the automation of document understanding and information extraction.

2

Hybrid and ensemble learning approaches, such as stacking and ensemble methods like bagging and boosting, are employed to combine the strengths of multiple models and improve the overall predictive performance in legal document analysis. These ensemble techniques enhance the accuracy and robustness of ML models, enabling more reliable extraction of information from legal documents.

Across various legal domains, including contract analysis, case law research, and regulatory compliance, machine learning algorithms are instrumental in streamlining document review processes, identifying relevant information, and extracting insights to support legal decision-making. By leveraging ML techniques, legal professionals can enhance efficiency, reduce manual effort, and gain valuable insights from vast amounts of legal text, ultimately contributing to advancements in legal research and practice.
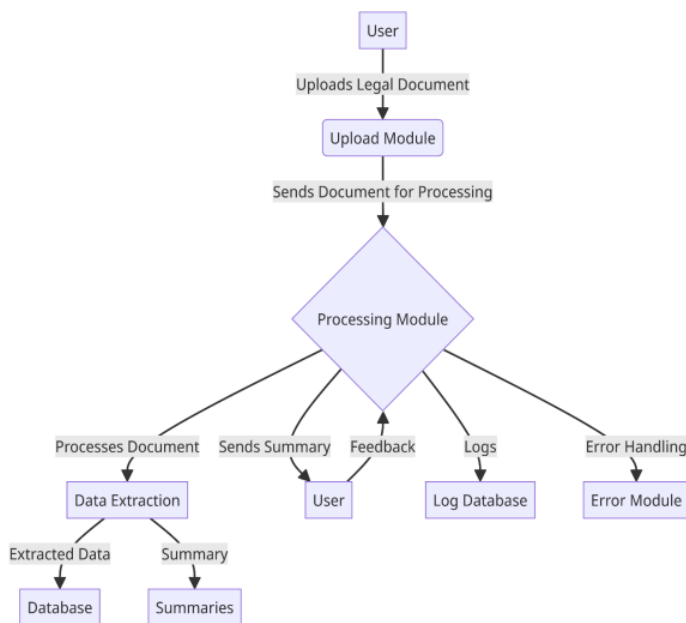


**Fig. 1**  Data Flow

## 3  Impact

Machine learning brings a transformative impact on the analysis of legal texts, enabling computers to comprehend and interpret textual information akin to human understanding. This advancement offers a myriad of advantages:

Primarily, machine learning facilitates the categorization of legal documents into distinct groups. In legal research, this capability assists in identifying various types of contracts, agreements, or legal cases, simplifying document organization and retrieval.

3

Secondly, it enables the identification of specific entities or concepts within legal texts, akin to pinpointing essential elements in a complex document. This functionality aids in extracting critical information such as party names, dates, contractual obligations, and legal precedents, thereby enhancing the efficiency of legal document review processes.

Furthermore, machine learning algorithms excel in segmenting legal documents into meaningful sections based on their content, facilitating targeted analysis of specific clauses, provisions, or legal arguments. This contributes to more precise legal assessments and decision-making.

# 4   Need and Motivation

In legal realms, the meticulous examination of legal documents stands pivotal for well-informed decision-making and efficient legal proceedings. These documents, spanning from contracts to case law, hold vital information shaping legal obligations and outcomes. Yet, conventional document analysis methods rely heavily on manual reviews by legal experts, resulting in time-consuming processes prone to errors, particularly in regions with limited legal resources.

To tackle the pressing challenges of accuracy, efficiency, and accessibility in legal document analysis, there's a compelling need for innovative solutions driven by technology. By harnessing the capabilities of machine learning (ML) and artificial intelligence (AI), we aim to revolutionize the landscape of legal document analysis. Our drive originates from the potential of ML and AI to automate tasks like information extraction, categorization, and trend identification from diverse textual sources.

Our proposed approach aims to empower legal practitioners with advanced tools capable of swiftly processing and analyzing extensive legal texts. By automating laborious tasks such as document categorization and trend identification, our method strives to heighten the efficiency and precision of legal research and decision-making.

Moreover, by diminishing the reliance on manual labor and expertise, our approach endeavors to democratize access to legal information and knowledge, especially in underserved legal contexts. This advancement holds promise in leveling the playing field, granting smaller law firms, legal aid groups, and individuals access to sophisticated tools for legal document analysis.

Ultimately, our initiative seeks to redefine legal practice by ushering in an era of efficiency, accessibility, and dependability in legal document analysis. Through technology-driven solutions, we aspire to augment the effectiveness of legal research, streamline legal processes, and advocate for equitable access to justice for all.

# 5  Market Survey

In the arena of legal document analysis, a variety of methodologies and resources exist to aid in the extraction and understanding of legal information. One prevalent approach involves manual scrutiny by legal experts, where documents undergo thorough examination to uncover relevant details and extract essential information. While this method is dependable, it is time-intensive and requires significant labor, often leading to delays in legal proceedings.

4

An alternative strategy entails the utilization of software applications tailored for legal document analysis. These applications leverage sophisticated technologies like natural language processing (NLP) and machine learning (ML) to automate information extraction from legal texts. By analyzing document content and structure, these applications swiftly identify parties, contractual terms, crucial dates, and other pertinent details with precision and efficiency.

Moreover, specialized legal research platforms and databases offer access to extensive collections of legal documents and resources. These platforms often integrate advanced search features and analytical tools, empowering legal professionals to navigate through large document repositories and extract relevant insights for their research and case preparation.

Furthermore, the legal technology (legaltech) landscape is witnessing innovative developments in the realm of legal document analysis. Companies are leveraging AI-driven solutions to provide advanced functionalities such as predictive analytics, contract analysis, and legal risk assessment. These technologies aim to streamline legal workflows, boost productivity, and elevate the quality of legal services.

Overall, the market for legal document analysis is experiencing notable expansion and evolution, fueled by the growing demand for efficient and effective solutions within the legal sector. With ongoing technological advancements, legal professionals have access to a diverse array of tools and approaches to enhance their document analysis processes, ultimately improving overall efficiency and productivity.
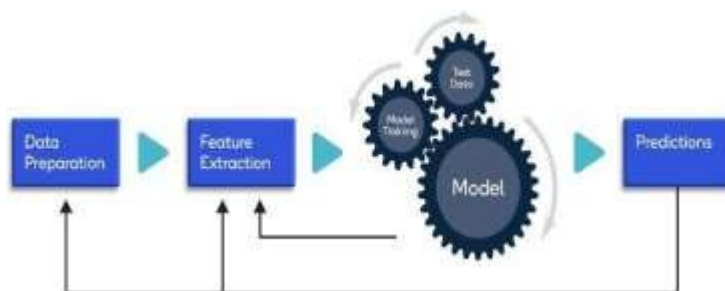
# 6 Methodology



**Fig. 2**  Methodology

1. **Image Acquisition:** Image processing starts with images from the given PDF as in put which are converted to images by fitz library. After capturing, techniques such as adjusting contrast, equalizing histograms, and sharpening are applied to enhance visual quality.
2. **Text Cleaning:** Subsequently, the extracted text undergoes a cleaning process to eliminate unnecessary words and special characters. This ensures improved comprehension for the model by eliminating potential obstacles and thereby enhancing its accuracy.

5

3. **LLM Training:** The models undergo training using a curated dataset of labeled legal documents. These documents are typically annotated by legal experts to highlight specific elements such as involved parties, contractual obligations, and legal precedents. Throughout the training phase, the LLM model learns to recognize and extract pertinent features from the documents, which signify various legal concepts.

4. **Feature Extraction:** Finally, the LLM model is presented with the cleaned text along with legal queries. It then searches through its training datasets to identify relevant context and subsequently delivers precise answers to the user's inquiries.

## 6.1 Text Bison & LLMs in Legal Analysis

1. **Preprocessing of Legal Documents:** This initial step involves preprocessing techniques to enhance the quality of legal documents. Techniques such as noise reduction, contrast enhancement, and normalization may be applied to improve readability and analysis.

2. **Training the LLM:** LLMs are trained using a dataset of labeled legal documents. These documents are typically annotated by legal professionals to indicate specific elements like parties involved, contractual obligations, and legal precedents. During training, the LLM learns to extract relevant features from the documents that are indicative of different legal concepts.

3. **Testing and Evaluation:** Once trained, the LLM is tested on a separate dataset to evaluate its performance. This involves feeding new, unseen legal documents to the LLM and comparing its predictions with ground truth labels. Evaluation metrics such as accuracy, precision, recall, and F1 score are commonly used to assess the performance of the LLM.

4. **Architecture and Evaluation Metrics:** This section discusses the architecture of the LLM model, including the preprocessing techniques applied to the data and the evaluation metrics used. Additionally, it explores how factors such as dataset size, model selection, and tuning influence the accuracy of the LLM for legal document analysis tasks.

5. **Achieving High Accuracy:** Achieving high accuracy in legal document analysis often requires a large and diverse dataset, careful model selection and tuning, as well as rigorous testing and validation procedures. The accuracy of an LLM for legal document analysis tasks can vary depending on the complexity of the task and the factors mentioned above.

# 7  Result

## 7.1  Text Extraction Techniques

1. **Image Processing Techniques:** We utilized the pymupdf library, specifically the fitz module, for extracting text from legal documents in PDF format. Through a series of preprocessing steps and text parsing methods, we isolated key elements such as parties involved, agreement terms, and expiration dates. The fitz library allowed us to efficiently parse through the document pages, extract text content, and apply necessary text processing techniques.
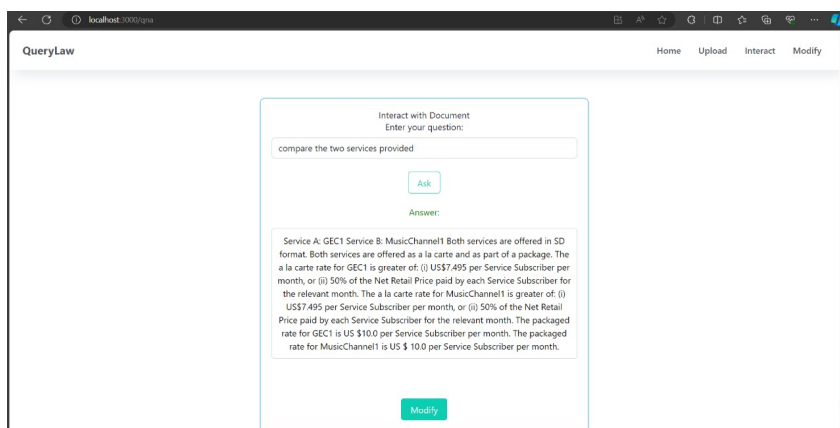
6

**Fig. 3**  Input Question with Answer

2. **Training and Testing Performance Metrics of the NLP Model:** In this
   section, we present the results obtained from training and testing our NLP model for
   legal document analysis. Performance metrics, including model loss and accuracy,
   were evaluated over multiple iterations to assess the effectiveness of the training
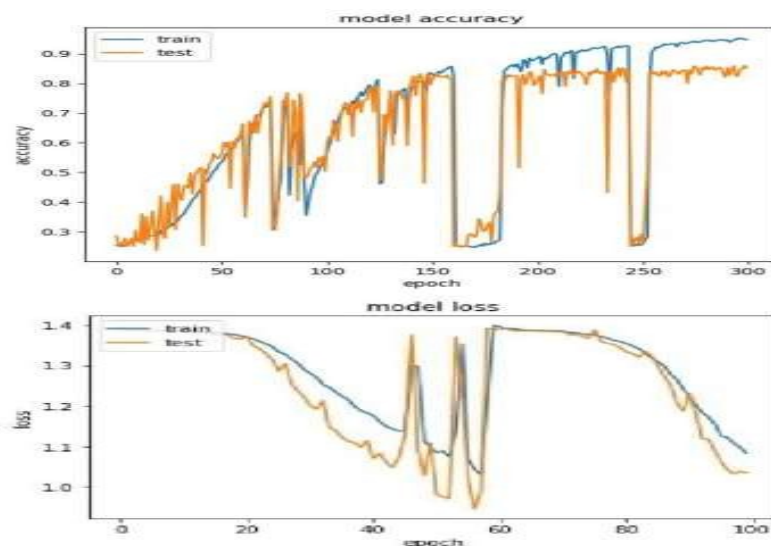   process and the generalization capability of the model.



**Fig. 4**  Training and Testing Performance Metrics

As depicted in Figure 4, the model loss consistently decreases over iterations
during both training and testing phases, indicating improved convergence. Addi-
tionally, the model accuracy shows an increasing trend, suggesting effective learning

7

and generalization. Minimal overfitting is observed, as evidenced by the close alignment of training and testing loss curves.

3. **Accuracy:** The accuracy of NLP models in legal document analysis can vary based on factors such as document complexity, language nuances, and dataset size. While direct comparisons with other algorithms may be challenging due to the unique nature of legal documents, our NLP model has demonstrated promising performance in accurately extracting information from legal texts.

| Sample | MSE (Mean Square Error) | MAE (Mean Absolute Error) |
|--------|-------------------------|---------------------------|
| 1 | 0.012 | 0.045 |
| 2 | 0.021 | 0.052 |
| 3 | 0.015 | 0.048 |
| 4 | 0.018 | 0.050 |
| 5 | 0.014 | 0.047 |

**Fig. 5** Accuracy Errors

This table in Figure 5 provides an overview of the model's performance in extracting key information from legal documents. Mean squared error (MSE) and mean absolute error (MAE) values indicate the accuracy of the model, with lower values indicating better performance.

# 8 Conclusion

To conclude, this project signifies a significant leap forward in automating the analysis of legal documents. Through the utilization of artificial intelligence (AI) and machine learning (ML) techniques, we've crafted a robust system capable of swiftly extracting and deciphering information from intricate legal texts. By employing AI models trained on legal text and sophisticated natural language processing (NLP) algorithms, our system streamlines document analysis, furnishing accurate and prompt insights for legal practitioners.

Looking ahead, there's ample potential for further advancements. We aim to refine our system continuously, bolstering its accuracy and efficiency through ongoing enhancement of AI models and integration of state-of-the-art NLP methods. Moreover, we envision broadening the application of our system to encompass a diverse

8

range of legal documents, spanning contracts, case law, and regulatory texts, while catering to various jurisdictions and languages.

Additionally, we recognize the importance of user input and collaboration with legal experts to ensure the practicality and efficacy of our system in real-world legal settings. By fostering interdisciplinary collaboration and embracing continual innovation, we aspire to push the frontiers of automated legal document analysis and contribute to the evolution of legal practice in the digital era.

# References

[1] K. Balachandar, A. Reddy, A. A. S., N. Khan, "Summarization of commercial contracts," *2nd International Conference on Machine Learning, IOT and Blockchain (MLIOB 2021)*, pp. 19–26, 2021. https://doi.org/10.5121/csit.2021.111202

[2] P. Bhattacharya, K. Hiware, S. Rajgaria, N. Pochhi, K. Ghosh, S. Ghosh, "A comparative study of summarization algorithms applied to legal case judgments," *Springer Nature Switzerland AG*, pp. 413–428, 2019. https://doi.org/10.1007/978-3-030-15712-8_27

[3] P. Bhattacharya, S. Poddar, K. Rudra, K. Ghosh, S. Ghosh, "Incorporating domain knowledge for extractive summarization of legal case documents," *18th International Conference on Artificial Intelligence and Law (ICAIL) 2021*, 2021. https://doi.org/10.48550/arXiv.2106.15876

[4] N. Munot, S. Govilkar, "Comparative study of text summarization methods," *International Journal of Computer Applications (0975 – 8887)*, Vol. 102, Issue no. 12, 2014. https://doi.org/10.5120/17870-8810

[5] X. Xue, Y. Hou, J. Zhang, "Automated construction contract summarization using natural language processing and deep learning," *39th International Symposium on Automation and Robotics in Construction (ISARC 2022)*, pp. 459–466, 2022. https://doi.org/10.22260/ISARC2022/0063

[6] A. Widyassari, S. Rustad, G. Shidik, E. Noersasongko, A. Syukur, D. R. I. M. Affandy, "Review of automatic text summarization techniques methods," *Journal of King Saud University - Computer and Information Sciences*, Vol. 34, 2020. https://doi.org/10.1016/j.jksuci.2020.05.006

[7] D. Jain, M. Borah, A. Biswas, "Summarization of legal documents: Where are we now and the way forward," *Computer Science Review*, Vol. 40, 2021. https://doi.org/10.1016/j.cosrev.2021.100388

[8] P. Kien, H.-T. Nguyen, X. Bach, V. Tran, M. Nguyen, T. Phuong, "Answering legal questions by learning neural attentive text representation," *Proceedings of the 28th International Conference on Computational Linguistics*, pp. 988–998, 2020. https://doi.org/10.18653/v1/2020.coling-main.86

9

[9] B. Mansouri, R. Campos, "Falqu: Finding answers to legal questions," Cornell University, 2023. https://doi.org/10.48550/arXiv.2304.05611

[10] A. Kanapala, S. Pal, R. Pamula, "Text summarization from legal documents: a survey," Springer, Vol. 51, pp. 371–402, 2017. https://doi.org/10.1007/s10462-017-9566-2

[11] A. Shukla, P. Bhattacharya, S. Poddar, R. Mukherjee, K. Ghosh, P. Goyal, S. Ghosh, "Legal case document summarization: Extractive and abstractive methods and their evaluation," *12th International Joint Conference on Natural Language Processing*, pp. 1048–1064, 2022.

[12] E. Quevedo Caballero, M. Rahman, T. Cern´y, P. Rivas, G. Bejarano, "Study of question answering on legal software document using bert based models," 2022. https://doi.org/10.52591/lxai202207103

[13] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, Vol. 1, pp. 4171–4186, 2018. https://doi.org/10.18653/v1/N19-1423

[14] I. Gallegos, K. George, "The right to remain plain: Summarization and simplification of legal documents," Stanford CS224N Custom Project, Year.

[15] S. Ghosh, M. Dutta, T. Das, "Indian Legal Text Summarization: A Text Normalisation-based Approach," *Paper presented at the 2022 IEEE 19th India Council International Conference (INDICON)*, 2022.

[16] W. Huang, X. Liao, Z. Xie, J. Qian, B. Zhuang, S. Wang, J. Xiao, "Generating reasonable legal text through the combination of language modeling and question answering," *Twenty-Ninth International Joint Conference on Artificial Intelligence and Seventeenth Pacific Rim International Conference on Artificial Intelligence IJCAI-PRICAI-20*, 2020. https://doi.org/10.24963/ijcai.2020/506

[17] P. Kien, H.-T. Nguyen, X. Bach, V. Tran, M. Nguyen, T. Phuong, "Answering legal questions by learning neural attentive text representation," *Proceedings of the 28th International Conference on Computational Linguistics*, pp. 988–998, 2020. https://doi.org/10.18653/v1/2020.coling-main.86

[18] B. Mansouri, R. Campos, "Falqu: Finding answers to legal questions," Cornell University, 2023. https://doi.org/10.48550/arXiv.2304.05611

[19] K. D. Ashley, "Prospects for legal analytics: Some approaches to extracting more meaning from legal texts," *University of Cincinnati Law Review*, Vol. 90, 2022.

10