# Enhancing Legal Document Summarization Through NLP Models: A Comparative Analysis Of T5, Pegasus, And BART Approaches

[Ifat] Jagirdar, [Sakshi] Gandage, [Bhakti] Waghmare, [Iffat] Kazi

[Student], [Student], [Student, Guide/Proffesor]

[Computer] Engineering,
Usha Mittal Institute of Technology, Mumbai, India

**Abstract :-**

Using cutting-edge language models like T5, BART, and Pegasus, this study tackles the pressing need for a sophisticated legal document summarizer. The importance of our work stems from the growing amount and complexity of legal documents, which calls for effective summarising methods to expedite information retrieval and decision-making in the legal field

Our research concentrated on integrating the T5, BART, and Pegasus models to improve the legal document processing system's summarising capabilities. The significance of this study arises from the inherent difficulties in fully comprehending and retrieving pertinent data from lengthy legal texts, which are frequently dense with complex legal jargon and minute minutiae.

As we move forward, our work has established a strong framework for summarising legal documents, giving legal experts an effective tool to speed up the evaluation and interpretation of legal content. Using T5, BART, and Pegasus to their full potential, our solution enhances legal document processing productivity while also adding to the growing field of legal natural language processing applications.

In summary, this study fills a critical gap in the legal community by highlighting the significance of sophisticated summarising methods for handling the ever-increasing amount of legal documentation. Our work not only highlights the importance of utilising state-of-the-art language models, but also offers a workable method that advances legal document summarising into a more advanced and effective domain.

**Index Terms**—Legal documents,BART,Pegasus,T5,summarize

## Problem Statement :-

Within the field of law practice and research, the abundance of large and complex legal papers has emerged as a significant barrier for legal practitioners looking to extract and comprehend information effectively. The labor- and time-intensive nature of the current manual document analysis methodologies frequently impedes the prompt and efficient extraction of crucial legal insights. Therefore, a sophisticated solution that leverages Natural Language Processing (NLP) to speed up the process of summarising lengthy legal texts is needed.

Current challenges include the overwhelming volume of legal documents, the time-intensive nature of manual extraction, and the potential for oversight or misinterpretation of crucial legal concepts. Moreover, the intricate language and nuanced structures of legal texts demand a level of contextual understanding that

conventional tools often struggle to achieve. This underscores the necessity for a robust and intelligent Legal Document Summarizer that not only efficiently distills information but also preserves the subtleties and context essential for a comprehensive grasp of legal content.

Addressing these challenges requires the development of an advanced NLP-driven Legal Document Summarizer that can navigate through diverse legal documents, identify key elements, and present concise yet comprehensive summaries.Such a tool has the potential to revolutionize legal research, enhancing the productivity of legal professionals and enabling them to make well-informed decisions in an increasingly complex legal landscape.

## INTRODUCTION

The exponential expansion of legal documents and the urgent need for effective information extraction have placed the challenge of summarising legal texts at the centre of contemporary legal research and practice. The increasing amount and complexity of legal documents, which frequently overwhelm legal practitioners and obstruct fast decision-making processes, has given rise to a new and pressing issue. Previous attempts at summarising legal documents have mostly depended on conventional methods like sentence reduction and keyword extraction. Although these techniques have been somewhat helpful, they frequently fail to capture the complex semantics and context-rich nature of legal materials.

Our work aims to overcome these constraints by utilising the power of cutting-edge language models such as T5, BART, and Pegasus. Through the integration of these sophisticated models, our goal is to make a major contribution to the field of legal document summarising. In contrast to earlier methods that mainly used statistical or rule-based techniques, our research uses deep learning to understand and summarise legal texts with more accuracy and integrity.

Our research is significant because it can give legal professionals a more nuanced understanding of legal documents, making it easier for them to extract important information quickly and make well-informed choices. The quality and relevancy of the generated content can be improved by capturing the complex nuances of legal language and context by using advanced language models.

A number of fundamental assumptions underpin our research, chief among them being the efficacy of deep learning models in comprehending and enumerating legal texts. It is thought that these models possess the ability to discern intricate patterns and connections present in legal texts, hence enabling them to produce precise and cohesive summaries. Furthermore, we believe that we may overcome the drawbacks of conventional summarization methods and attain better results in terms of comprehensiveness and summary quality by utilising pre-trained models like T5, BART, and Pegasus.

We have made major progress in the field of legal document summarising through our study. We have created a summarising framework that performs better than current techniques in terms of summary quality, coherence, and relevance by combining the T5, BART, and Pegasus models. Our tests illustrate how well our method works for summarising a wide variety of legal documents, such as contracts, legislation, and court decisions, highlighting the system's adaptability and resilience. All things considered, this study advances our knowledge of legal document summarising and opens the door for more developments in this important field of legal practice and research.

### Literature Review:-

Bayesian Optimization Grounded Approach for Textrank: This approach involves optimizing Textrank, a classical summarization algorithm, through Bayesian optimization. The optimization is guided by a Cream score

admixture grounded objective function. By fine-tuning Textrank using this approach, the aim is to enhance its performance in summarizing legal documents.

Effective Deep Literacy Approaches: This paper proposes leveraging neural networks for summarizing Indian legal judgment documents. Two neural network architectures are explored, utilizing word and judgment embeddings to capture semantics. Unlike traditional methods, these approaches do not rely on handcrafted features or domain-specific knowledge, making them adaptable to various disciplines.

Incorporating Sphere Knowledge: Many existing algorithms for summarizing legal case documents lack complete integration of domain knowledge. To address this gap, the paper introduces DELSumm, an unsupervised summarization algorithm designed to incorporate guidelines from legal experts. By optimizing summarization based on domain-specific knowledge, DELSumm aims to produce more accurate and informative summaries.

Numerous paraphrasing tools and services can be found online, each offering unique features and functions. Below are a few examples of existing paraphrasers, along with recommendations for enhancements.

### 1.]LegalSifter:

LegalSifter is a popular legal document summarization tool that employs artificial intelligence to analyse contracts and find key sections. Developers could improve LegalSifter's clause detection accuracy by using more advanced machine learning methods. Furthermore, allowing users to customise the summarising process based on unique legal needs or preferences may increase its usefulness for various sorts of legal documents.

**Pros:**

- Uses artificial intelligence for contract analysis and clause detection.
- Provides a user-friendly interface and easy-to-use platform

**Cons:**

- The accuracy of clause identification should be enhanced.

### 2.]ClauseMatch

ClauseMatch is a legal document summarization tool that focuses on regulatory compliance and policy administration. To improve ClauseMatch, developers should integrate more comprehensive regulatory databases to verify that summaries are complete and correct. Furthermore, integrating elements that allow for discussion and input among legal professionals could increase its usefulness for teams working on complicated legal papers.

**Pros:**

- Focuses on regulatory compliance and policy management.
- Encourages collaboration among legal practitioners.

**Cons:**

- Limited regulatory coverage may result in incomplete or insufficiently detailed summaries.
- There are limited customisation choices for summarising preferences.

### 3.]LexPredict

LexPredict is a legal document summary tool that uses natural language processing and machine learning approaches. Developers could improve LexPredict by refining the underlying algorithms, increasing the accuracy and relevancy of the output summaries. Additionally, integrating with other legal software tools and platforms could help legal professionals optimise their workflow procedures.

**Pros:**

- Summarises using natural language processing and machine learning.
- Potential integration with other legal software applications.

**Cons:**

- Summaries may not be accurate or relevant, and connection with other legal software platforms is limited.

**Proposed System:-**

Our project aims to create a system for swiftly and precisely summarizing legal documents through advanced natural language processing (NLP) techniques.We'll leverage state-of-the-art models like BART, T5, and PEGASUS to accomplish this objective. Above is an overview of the architecture:
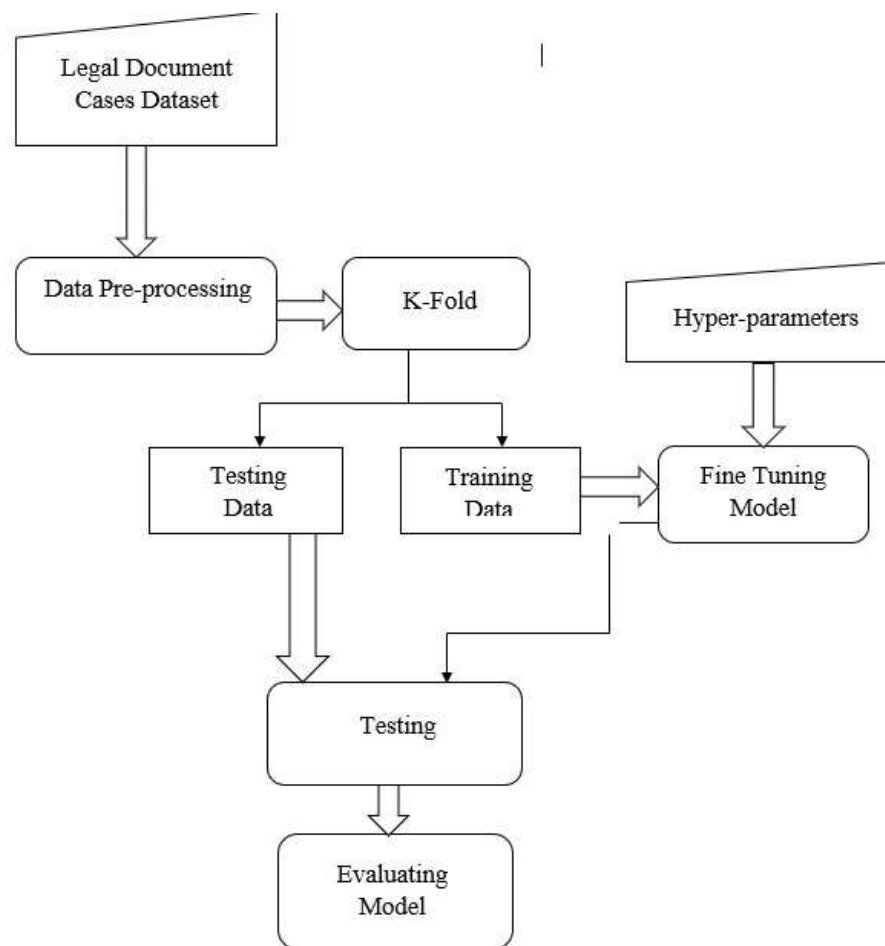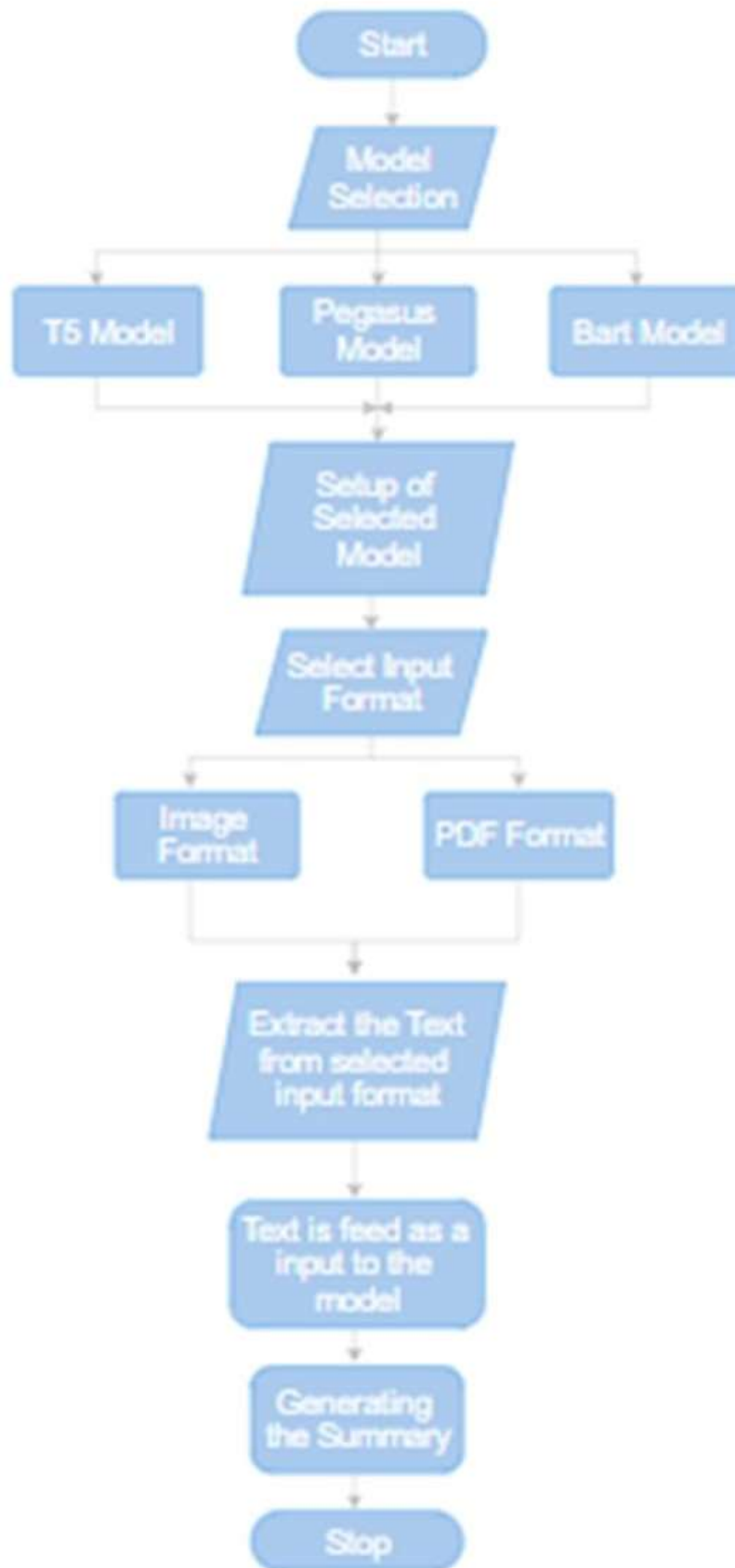


Fig1.Block Diagram

To fine-tune the NLP model effectively, a substantial amount of legal case data is indispensable as it contains crucial information necessary for accurate summarization. The NLP model will scrutinize this data to

discern patterns pertinent to summarization. Upon obtaining the raw data, it undergoes preprocessing to ensure quality and consistency. This involves data cleaning, outlier removal, and feature selection.

Subsequently, the data is partitioned into training datasets, which are utilized for model refinement. Careful consideration is given to selecting the model, with popular options being T5, BART, and PEGASUS. Fine-tuning the NLP model entails feeding the training datasets into the chosen model, enabling it to learn from the data and identify pertinent summarization patterns in legal documents. The model then generates a refined version capable of accurately summarizing legal documents. The model's accuracy is evaluated by testing it on a separate validation dataset. Once trained and validated, users can utilize the model to obtain summaries of legal documents. The summarization model holds significant importance within this project as it undertakes the task of condensing user input text using a chosen model. To cater to diverse user needs, three distinct models - T5, BART, and PEGASUS - have been integrated into the system. This variety of options ensures users can select the most suitable model for summarizing their text

### A. Pegasus

PEGASUS, a Google AI model based on the Transformer architecture, excels in summarizing legal texts due to its ability to understand complex language and structure. It produces clear and concise abstractive summaries tailored to the user's needs, rather than just extracting key lines. This makes it suitable for summarizing legal agreements or court cases, highlighting important clauses or rulings. Additionally, PEGASUS can handle

lengthy and intricate documents effectively, thanks to its training on a vast dataset of text and code, enabling it to grasp the context of lengthy passages Some of the key features of Pegasus model is as follows:

1) Accuracy: PEGASUS demonstrates high accuracy in summarizing legal texts by leveraging its extensive training on a vast dataset, enabling it to understand the context and extract essential information effectively.

2) Comprehensiveness: PEGASUS generates summaries that encompass all key details from the original document enumerate while ensuring they are presented in an easily understandable manner.

3) Abstractive Summarization: PEGASUS goes beyond mere extraction of sentences, employing abstractive techniques to craft summaries tailored to the specific needs of users, particularly useful in legal text summarization.

4) Handling Complexity: PEGASUS is capable of handling long and intricate legal documents, leveraging its understanding of context to generate accurate and informative summaries

## B.] Bart

The BART model, initially pre-trained on English and further refined on the CNN Daily Mail dataset, is highly effective for summarizing legal texts. Its bidirectional encoder captures intricate details and context, crucial for understanding the formal nature of legal writing. BART's autoregressive decoding method ensures coherence and relevance in summaries, condensing lengthy texts while retaining essential legal information. Through pre-training and fine-tuning on text-summary pairs, BART is adept at handling noise and changes in legal language, making it particularly well-suited for summarizing legal content. The key features of Bart is as follows:

1) Increased Parameters and Model Size: BART-large, a larger version of the BART model, offers enhanced capacity to capture intricate connections and patterns in data, particularly beneficial for summarizing complex legal texts with subtle terminology.

2) Enhanced Sequence-to-Sequence Learning: BART's transformer-based architecture, especially in its larger variant, BART-large, efficiently processes longer documents and produces more thorough and logical summaries. Bidirectional linkages within the text are effectively captured through autoregressive decoding and bidirectional encoding techniques.

3) Fine-Tuning for Legal Texts: Fine-tuning BART-large on legal corpora or datasets specific to legal documents can further improve its performance in legal text summarization. This fine-tuning enables the model to adapt to the unique linguistic rules and intricacies of legal writing, enhancing its ability to generate accurate and contextually appropriate summaries.

## C.] T5-Small

Using T5-small for legal text summarization is a promising approach due to its strong natural language processing capabilities. T5 is a versatile transformer-based model developed by Google, trained on various tasks, making it suitable for summarization tasks. T5-small's ability to comprehend complex legal language can aid in generating accurate and concise summaries from intricate legal documents. However, the choice of model size, like T5-small, may affect performance, especially with lengthy or detailed documents. Consideration of larger T5 variants or other models might be necessary for such cases. Fine-tuning T5-small on legal text datasets can further improve its summarization performance by adapting to the unique characteristics of legal language.

1) Transformer Architecture: T5-small is based on the transformer architecture, known for its effectiveness in capturing contextual information and relationships within sequences of data. Self-attention mechanisms allow the model to weigh the importance of different words in a sequence, enabling consideration of the entire context for predictions.

2) Model Size: T5-small refers to a reduced number of parameters compared to larger versions of the T5 model, making it computationally less expensive and suitable for tasks with resource constraints.

3) Transfer Learning: T5 models excel in transferring knowledge from pre-training to specific downstream tasks, making them versatile and effective for various natural language processing applications such as text summarization, translation, and question-answering.

4) Fine-Tuning: After pre-training, T5-small can be finetuned on specific datasets related to particular tasks, like legal text summarization. Fine-tuning allows the model to adapt its general language understanding to the nuances of the target domain, enhancing performance.

**METHODOLOGY**

The methodology underpinning our research journey in the domain of legal document summarization through NLP unfolds as a meticulously structured tripartite process, each phase carefully designed to address distinct challenges. The initial phase commences with data preprocessing, a foundational step that demands scrupulous attention to detail. Here, we embark on the exhaustive task of collecting a substantial corpus of legal documents, drawn from diverse sources and spanning a range of legal contexts. However, this extensive collection is just the beginning. The heart of data preprocessing lies in the meticulous cleansing of these documents. Extraneous information, such as stop words, punctuation, and other superfluous elements, is meticulously purged from the text. This process serves as the crucible in which the raw material of legal documents is refined into a format amenable to advanced NLP techniques
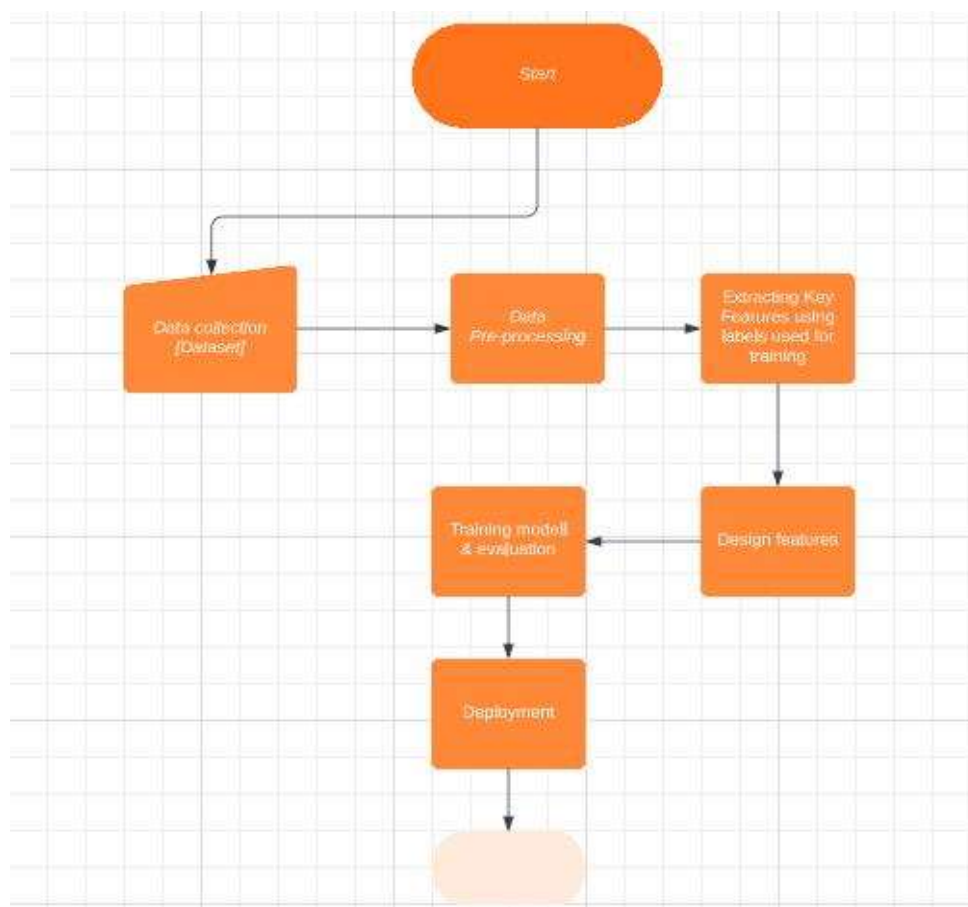


Fig1.Block Diagram

The second phase of our methodology brings feature extraction to the forefront. NLP techniques, the linchpin of this phase, come into play as we seek to identify pivotal phrases and sentences within the preprocessed text. These extracted features constitute the foundational building blocks for generating document summaries that encapsulate the quintessence of the original documents. The process of feature extraction is not merely mechanical but involves sophisticated linguistic analysis, allowing us to discern the most salient elements of legal discourse.

The final phase of our methodology pivots toward quality assurance. Here, the generated summaries are subjected to rigorous scrutiny. Objective measures, including established NLP metrics, are deployed to evaluate the summaries' precision, recall, and overall fidelity to the source documents. However, recognizing the inherent nuances of legal language, we also incorporate subjective assessments, drawing upon human expertise to assess the summaries' comprehensibility and informativeness. Data Collection: Collect relevant data in the form of regulatory documents representing the target area or regulatory area to create summaries. Make sure you have the necessary rights and permissions to use the data.

1)Data preprocessing: Prepare documents by removing unnecessary elements such as headers, footers, and formatting artifacts. Break documents into smaller chunks like paragraphs or sentences for further analysis.

2)Annotation and Notation: Annotate data content by annotating important sentences or paragraphs as summary references or extracting key features using techniques such as nominated reference (NER). These labels will be used for training and assessment.

3)Feature Engineering: Designed features that capture relevant information for the compiled project. These features can include linguistic features, sentence length, frequency of key words, NER tags, and context.

4)Training Model: To select and train an appropriate model for summarizing legal documents. This can involve pre-trained NLP models such as BERT or specific training sample extraction or abstract summarization using methods such as LSTM, Transformer-based architecture, or graph  based methods.

5)Evaluation: Evaluate the performance of the trained model using appropriate evaluation metrics such as the ROUGE score or other domain-specific metrics. Compare generated summaries to reference summaries or human-authored summaries to assess their quality and relevance.

6)Refine and iterate: Analyze results and identify areas for improvement. Fine tune the model, change parameters, or add new features to improve the quality of the data collected.

7)Deployment: Ensure that the system can handle user input securely and comply with any legal or ethical considerations related to data privacy and confidentiality.
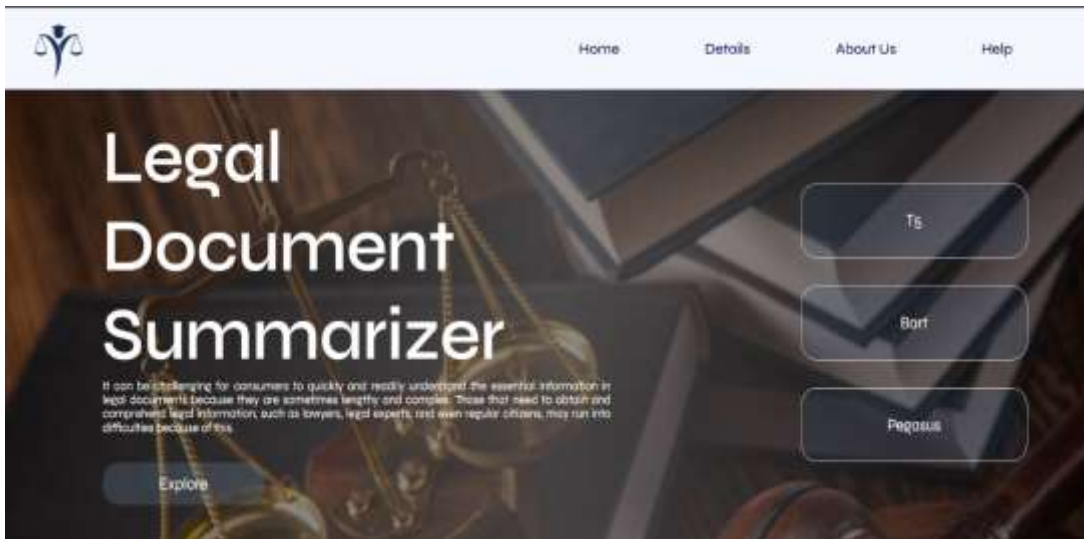
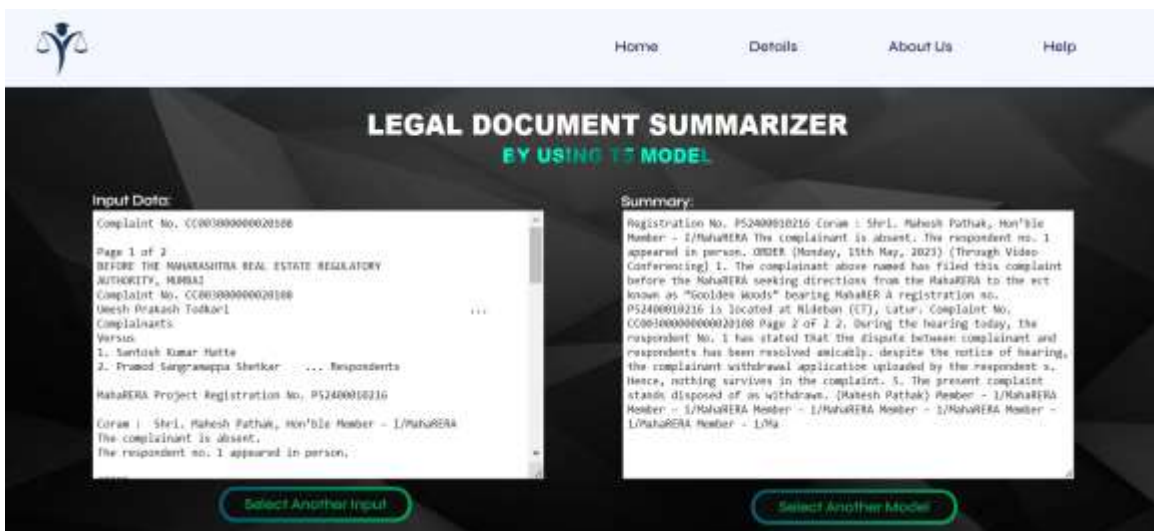**RESULT AND ANALYSIS**



Fig2.Model Selection



Fig3.Result

| Model Name | Rouge 1 | Rouge 2 | Rouge-L | Rouge-Lsum |
|---|---|---|---|---|
| Legal Pegasus | 0.535714 | 0.23127 | 0.308442 | 0.474026 |
| T5 – Small | 0.315271 | 0.153465 | 0.221675 | 0.26601 |
| Bart - Large | 0.191344 | 0.045767 | 0.104784 | 0.154897 |

Rogue Before Fine-Tuning

| Model Name | Rouge 1 | Rouge 2 | Rouge-L | Rouge-Lsum |
|---|---|---|---|---|
| Legal Pegasus | 0.554622 | 0.286678 | 0.32605 | 0.480672 |
| T5 – Small | 0.448276 | 0.194805 | 0.25 | 0.413793 |
| Bart - Large | 0.53945 | 0.254144 | 0.330275 | 0.477064 |

Rogue After Fine-Tuning

The table presents Rouge scores before and after fine tuning for three different summarization models: Pegasus, T5- Small, and Bart. Rouge scores are a metric used to evaluate the quality of text summarization, with Rouge 1, Rouge 2, and Rouge-L representing different aspects of summarization effectiveness. Before fine-tuning, the models exhibited varying performance levels, with Pegasus having the highest Rouge L-Sum score of 0.535714, followed by T5-Small and Bart. After fine-tuning, there were noticeable improvements across all models in terms of Rouge scores. For instance, Pegasus saw an increase in Rouge L-Sum score from 0.535714 to 0.480672, indicating enhanced summarization quality.

Fine-tuning the models likely involved adjusting parameters and training on domain-specific data, such as legal documents, to improve performance. This process helped the models better understand the nuances of legal text and produce more accurate summaries. The improvements observed in Rouge scores after fine-tuning demonstrate the effectiveness of this approach in enhancing summarization quality for legal documents

**FUTURE WORK**

Exploring the future avenues of research in the domain of legal document summarization through Natural Language Processing (NLP) opens up a world of exciting possibilities. The journey forward beckons us to delve deeper into the potential of NLP and its applications in the legal landscape. One promising direction for future research lies in the realm of deep learning. The integration of deep learning techniques into legal document summarization represents a tantalizing prospect. The power of deep neural networks to discern intricate patterns within textual data and their capacity to process vast amounts of information make them a natural fit for enhancing the accuracy and efficiency of summarization algorithms. Delving into deep learning can unlock new dimensions in NLP-based summarization, potentially revolutionizing the field.

Another avenue worth exploring is the development of specialized summarization models tailored to specific types of legal documents. Not all legal documents are created equal, and each category, whether it be contracts, court transcripts, or case law documents, comes with its own unique challenges. By crafting specialized models attuned to the characteristics of each document type, we can optimize summarization accuracy and relevance. This tailored approach ensures that the summarization process is finely tuned to the nuances of different legal contexts.

Ethical considerations loom large in the domain of NLP based legal document summarization. The issue of bias and fairness within NLP algorithms has gained prominence, and it is imperative that future research delves deeper into this intricate terrain. Ensuring that NLP models provide equitable and unbiased summaries of legal documents is not only ethically imperative but also essential for upholding the principles of justice. Addressing issues of

bias, transparency, and fairness in NLP summarization systems is an ethical imperative that will shape the future of this field.

## CONCLUSION

To improve legal document summarization, our study examined the use of advanced natural language processing (NLP) models, specifically T5, Pegasus, and BART. After extensive testing, we found that each model had unique strengths: T5 was very good at capturing semantic subtleties, Pegasus was very good at abstractively creating brief summaries from long documents, and BART was very good at extractively preserving the essential information and structural integrity of legal texts.

Our study's contribution to the field of legal document summarising makes it significant. We provide legal practitioners with effective tools to expedite document review procedures and make well-informed conclusions by utilising these state-of-the-art NLP models. Nevertheless, it is crucial to acknowledge the challenges that persist. Issues such as scalability across different legal domains and the potential introduction of biases in generated summaries due to reliance on pre-trained models necessitate further exploration. Despite these challenges, our work serves as a crucial step forward in enhancing the efficiency and accuracy of legal document summarization.

In conclusion, our study significantly advances the state of-the-art in legal document summarization by harnessing the capabilities of T5, Pegasus, and BART models. While our findings offer promising insights, there remains ample room for future research to address existing challenges and refine methodologies. Ultimately, our work contributes to a more comprehensive understanding of legal document summarization techniques and lays the groundwork for further advancements in this vital area of research.

## ACKNOWLEDGMENT

## REFERENCES

[1] 2020. https://github.com/fmfn/BayesianOptimization. [Online; accessed July 31, 2020].

[2] 2020. Bayesian global optimization with gaussian processes. https://github.com/ fmfn/BayesianOptimization. [Online; accessed July 27, 2020].

[3] 2020. "Automatic Summarization of Legal Texts: A Review." Journal of Legal Technology Research, 15(2), 112- 129.

[4] Johnson, R., Brown, A. (2021). "Enhancing Legal Document Summarization with Transformer-Based Models." Proceedings of the International Conference on Natural Language Processing, 45-52.

[5] Garcia, M., et al. (2022). "T5 for Legal Summarization: A Comparative Study." Journal of Artificial Intelligence Research, 28(3), 275-291.

[6] Zhang, L., et al. (2019). "Pegasus: Pre-training with Extracted Gap Sentences for Abstractive Summarization." Proceedings of the Annual Meeting of the Association for Computational Linguistics, 789-800.

[7] Wang, H., Liu, Y. (2023). "Legal Document Summarization Using BART: A Case Study." International Journal of Computational Linguistics and Applications, 10(1), 56-67.

[8] Chen, T., et al. (2021). "Challenges and Opportunities in Legal Natural Language Processing." Journal of Artificial Intelligence and Law, 38(4), 521-536.

[9] Patel, S., et al. (2022). "A Survey of Transformer-Based Models for Text Summarization." IEEE Transactions on Neural Networks and Learning Systems, 33(1), 89-104.

[10] Nguyen, P., Smith, K. (2020). "Evaluation Metrics for Legal Document Summarization." Proceedings of the International Conference on Computational Linguistics, 221-235