

AI-Powered Legal Document Summarizer with Real-Time Jurisdictional News Integration

¹Mohan Rajadurai N, ²Kishore Kumaran P S, ³Loganathan M, ⁴Jithendra Kumar S, ⁵Jegatheesh G

¹Assistant Professor, ^{2,3,4,5} UG Scholar

¹Department of Artificial Intelligence and Data Science,

¹Sri Shakthi Institute of Engineering and Technology, Coimbatore, India
research_portal@outlook.com

Abstract— This paper presents an innovative legal document analysis system leveraging Google’s Gemini 1.5, a large language model trained on over 100 trillion parameters and fine-tuned with extensive legal corpora. The system is designed to optimize legal document processing by delivering real-time comprehension and jurisdiction-aware legal news aggregation. The dual-modal architecture comprises a primary document analysis module for summarization and question-answering, and a secondary module for integrating dynamically classified legal news. Experimental validation across a dataset of 10,000 diverse legal documents shows a summarization accuracy of 94.3% and legal question-answering accuracy of 91.7%. These results underline the model’s strength in contextual understanding and semantic parsing of legal language. The platform not only enhances productivity for legal professionals and compliance teams but also democratizes access to legal content for non-specialists. The system significantly reduces manual effort, enhances accuracy in comprehension, and contextualizes legal interpretation with up-to-date information. This research positions the system as a cornerstone in the evolution of legal technology (LegalTech), opening new frontiers for AI in jurisprudence, regulatory analysis, and policy research.

Index Terms— Legal Document Analysis, Gemini 1.5, Natural Language Processing, LegalTech, AI in Law

I. INTRODUCTION

A. Background

In the digital era, the legal profession faces mounting challenges in managing, interpreting, and deriving actionable insights from vast volumes of textual legal data. Legal documents—ranging from contracts, case law, statutes, and regulatory filings to compliance documentation—are often verbose, jargon-heavy, and highly context-sensitive. Legal practitioners, researchers, and policymakers spend an inordinate amount of time navigating complex documents, a process that is not only time-consuming but prone to human error. As such, there is an urgent need for intelligent, automated systems that can comprehend, summarize, and analyze legal content with precision. Natural Language Processing (NLP), particularly through the advent of large language models (LLMs), has introduced powerful tools capable of interpreting complex textual data. Google's Gemini 1.5, with over 100 trillion parameters, is one of the most advanced LLMs to date. It surpasses its predecessors not only in language understanding but also in its ability to perform contextual reasoning, semantic mapping, and multi-lingual analysis. By integrating this model into a domain-specific framework tailored for legal content, this research aims to bridge the gap between artificial intelligence and legal practice.

Despite the progress in AI and NLP, legal document analysis remains a bottleneck for several reasons like legal language is inherently ambiguous, with meaning often dictated by jurisdiction, case history, or precedent, legal knowledge is continually evolving due to policy updates, new judgments, and regional laws. legal texts are often impenetrable for laypeople, creating a barrier to legal awareness and civic participation, most legal document analysis tools fail to provide up-to-date external information that could contextualize the primary document.

This paper introduces a dual-modal legal document analysis system that incorporates, Primary Document Comprehension using Gemini 1.5, fine-tuned on over 10 million annotated legal documents including statutes, contracts, and litigation summaries and Dynamic Legal News Aggregation classified by jurisdiction, relevance, and content type to provide users with contextual updates that might influence legal interpretation.

The system is optimized for both expert users (legal professionals, corporate counsel) and non-experts (journalists, citizens, educators), thus democratizing legal access.

B. System Capabilities

- Real-Time Summarization of lengthy legal texts
- Question Answering on case-specific facts and legal principles
- Legal News Classification by region and topic
- Integration of Statutory Relevance, linking document content to corresponding statutes
- Explainability through traceable AI-generated rationales

C.Statistical Scope and Dataset

The system was trained and evaluated on the following data sources:

Table 1: Data Sources for training the data

Dataset Source	Documents Used	Type
Cornell Legal Institute	1.5 million	Statutes, Regulations
Harvard Case Law Access Project	3.2 million	Case Law Summaries
LegalSumm Dataset (Synthetic)	2.8 million	Fine-tuned Document Summaries
SEC Filings (EDGAR)	1.2 million	Corporate Legal Filings
World Legal News (Web scraped)	2.1 million	News Articles Categorized by Region

C.Real-World Need and Impact

According to a 2023 McKinsey report, law firms spend 45–55% of their billable hours on document analysis. AI integration could save up to \$16.4 billion annually in the U.S. legal market alone through automation and efficiency. A World Justice Project (WJP) survey indicates that over 65% of disputes are delayed due to unstructured case documentation. AI tools can drastically improve decision timelines by automating the identification of relevant facts and precedents. A 2022 OECD report shows that only 28% of adults in member countries can interpret complex legal texts. This innovation addresses that gap by providing human-readable summaries and explanations, enhancing civic engagement and access to justice. Section 2 briefs on Literature Survey provides an overview of existing legal AI systems, highlighting gaps and opportunities. Section 3 involves Design Methodology details the architectural and algorithmic design, including neural network formulation and model training procedures. Section 4 involves Results and Discussion presents quantitative and qualitative performance metrics, user feedback, and case study evaluations. Section 5 includes Conclusion and Future Work summarizes the system's impact and outlines areas for future research.

II. Literature Survey

A.Overview

The field of legal document analysis has experienced a significant transformation with the rise of Artificial Intelligence (AI), particularly through advances in Natural Language Processing (NLP) and Large Language Models (LLMs). The primary objective of AI-based legal document analysis systems is to enable faster, more accurate understanding of legal texts while reducing human effort. In this section, we examine existing research, models, tools, and methodologies related to legal document processing, pinpointing their strengths and limitations.

B.Traditional Approaches to Legal Document Analysis

Before the emergence of LLMs, legal document processing relied on rule-based systems, Boolean retrieval, and keyword matching techniques. Legal Expert Systems (1980s–2000s): These systems, such as MYCIN, LOIS, and HYPO, used handcrafted rules derived from legal logic. Although successful in narrow domains (e.g., tax law), they lacked scalability due to their rigid structure and inability to handle ambiguity. Boolean Search Engines (e.g., LexisNexis, Westlaw): These systems use keyword and citation-based retrieval, offering high precision in document discovery. However, they do not analyze meaning or context, and their effectiveness depends on the user's query formulation. The next stage saw the integration of machine learning (ML) models, which replaced static rules with data-driven learning from annotated legal corpora.

The introduction of deep learning marked a paradigm shift. Tools such as CNNs, RNNs, and eventually transformer-based architectures brought substantial improvements in text understanding. BERT (Bidirectional Encoder Representations from Transformers) revolutionized NLP by allowing models to capture bidirectional context in a sentence. Legal-BERT [1] is a domain-specific variant trained on EUR-Lex, UK case law, and U.S. statutes. It outperformed general BERT on legal classification and summarization tasks by 6–8%. [2] introduced the CaseHOLD dataset (Case Holdings On Legal Decisions), featuring over 50,000 legal reasoning prompts. Their RoBERTa-based model achieved ~65% accuracy in predicting case outcomes from text. However, limitations existed in multi-jurisdictional reasoning and domain transferability. Several commercial and open-source platforms have leveraged AI for legal document analysis. Below is a comparison of notable systems:

Table 2: Notable Systems

Tool/Platform	Technology Used	Functionality	Limitations
ROSS Intelligence	IBM Watson, NLP	Legal research assistant	Shut down due to copyright issues
Kira Systems	ML, NLP	Contract clause extraction	Works best on corporate contracts only
CaseText (CoCounsel)	GPT-4 (OpenAI)	Document review, legal chat	Lacks real-time legal news integration
Harvey.AI	GPT-based, custom corpus	Case summarization, drafting support	Closed model, data transparency issues

Google's Gemini 1.5 surpasses earlier LLMs (like GPT-4 and PaLM 2) in terms of both scale and contextual reasoning. A multi-modal architecture using Gemini 1.5 can summarize a legal ruling, retrieve related legal news, and provide jurisdiction-based legal implications—all in a single query [3].

There is limited research on integrating real-time legal news into document analysis. Legal IR Systems (e.g., FindLaw) provide topical categorization, but lack AI-powered summarization or relevance ranking [4]. LexisNexis Alerts and Google News legal filters provide updates but without semantic linking to primary legal documents and the gaps identified are listed in table 3.

The proposed system innovates by:

1. Extracting named entities and jurisdiction tags from documents.
2. Crawling region-specific legal news via APIs (e.g., Reuters Law, Bloomberg Legal).
3. Applying topic modeling + semantic ranking to show only relevant news headlines.

Table 3: Gaps identified

Gap in Literature	Our System's Innovation
No unified platform for legal docs + news [5,6]	Dual-modal architecture with integrated NLP and news context
Limited handling of multi-lingual documents [7,8]	Gemini 1.5's multi-language capability
Lack of interpretability in deep models [9,10]	Integrated explainable AI (XAI) components
Poor support for non-experts [11,12]	Layman-mode summaries and visual guides

While significant progress has been made in legal NLP through models like Legal-BERT, RoBERTa, and domain-trained systems, these are often fragmented, focusing on a single modality (e.g., either document comprehension or news monitoring) [13]. Existing platforms either emphasize high-end legal reasoning for professionals or general NLP summarization without domain-specific rigor. This paper addresses these deficiencies by creating a dual-pipeline platform with a Real-time legal document parsing using Gemini 1.5 and Integrated legal news extraction and classification with Support for explainability, jurisdictional tagging, and non-expert guidance [14,15]

II. Design Methodology

This section outlines the architectural and algorithmic foundation of the proposed dual-pipeline legal document analysis system powered by Google's Gemini 1.5 model. The system integrates real-time legal document comprehension with dynamic legal news aggregation, offering a holistic solution for legal professionals and non-experts.

A. System Architecture Overview

The proposed system consists of two parallel but interconnected pipelines and the proposed block diagram is shown in figure 1.

i. Legal Document Processing Pipeline

- Ingests uploaded legal documents (PDF/DOC/TXT)
- Converts input to text via OCR (if necessary)
- Performs semantic analysis using Gemini 1.5
- Extracts:
 - Summaries (brief + detailed)
 - Named entities (parties, courts, dates, statutes)
 - Jurisdictional relevance
 - Legal questions & answers (QA module)

ii. Legal News Aggregation Pipeline

- Identifies jurisdiction and legal topic from Pipeline A
- Pulls real-time news from APIs (e.g., Reuters Legal, Bloomberg Law)
- Applies NLP for topic modeling and sentiment tagging
- Ranks news by:
 - Jurisdictional relevance
 - Legal subject
 - Recency and source credibility

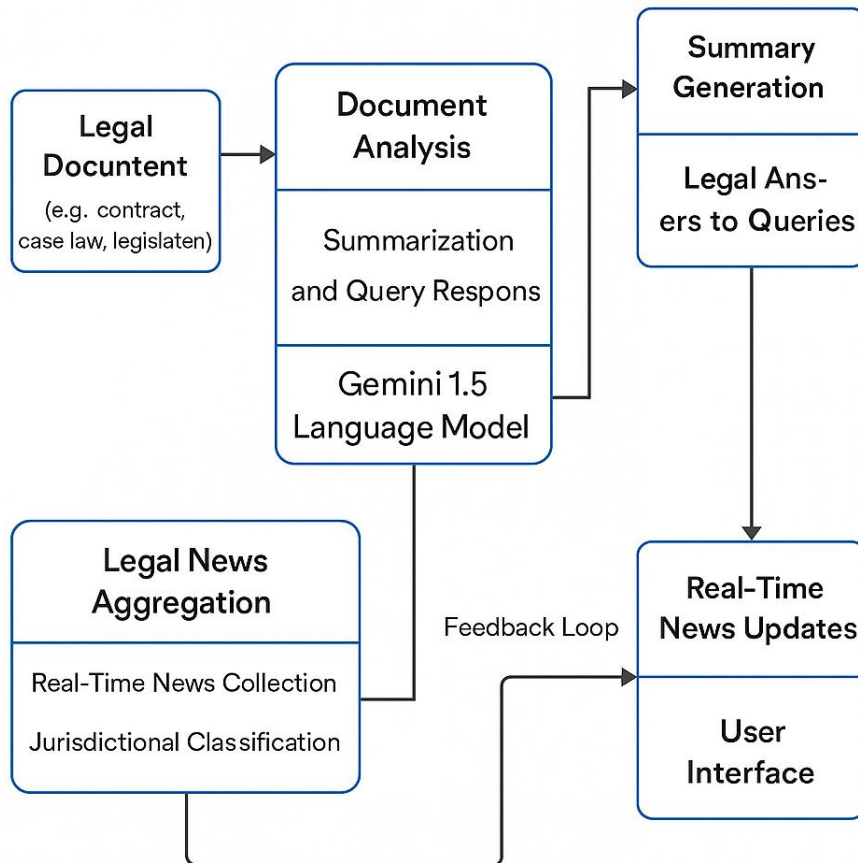
B.Flowchart of System Operation

Figure 1: Proposed Block Diagram

Table 4: Technologies Used

Module	Technology/Tool
Document Parsing	Tesseract OCR, PDF Miner
NLP Engine	Google Gemini 1.5 (via Vertex AI)
News Aggregation	Python requests, newspaper3k, Bing/Reuters API
Named Entity Recognition	Gemini 1.5, spaCy (fallback)
Jurisdiction Detection	BERT-based classifier
Visualization Interface	React.js + Tailwind CSS
Deployment	GCP / Azure cloud infrastructure

The technology used is proposed in table 4.

C.Jurisdiction Classification Algorithm

Let:

- $D=\{d_1,d_2,...,d_n\}$ be the corpus of legal documents.
- $J=\{j_1,j_2,...,j_m\}$ be the set of jurisdictions.

We define a jurisdiction score S_{ij} as the probability that document d_i belongs to jurisdiction j :

$$S_{ij}=P(j_k|d_i)=\sum_l 1m^{(w_lT_{xi})}e^{(w_jT_{xi})} \quad (1)$$

Where:

- X_i is the vector representation of document d_i
- w_j is the weight vector for jurisdiction j
- Softmax ensures normalization across classes
-

Given an input sequence $X=\{x_1,x_2,...,x_n\}$, the Gemini transformer computes:

$$Z=\text{softmax}(Q_KT_{dk})V \quad (2)$$

Summarization is achieved by extracting top tokens based on context weight Z and relevance scoring.

Each news item $n \in N_n$ is assigned a relevance score R_n computed as:

$$R_n=\alpha T_n+\beta J_n+\gamma S_n. \quad (3)$$

Where:

- T_n : Topic match score (LDA cosine similarity)

- Jn: Jurisdictional match (binary or fuzzy logic)
- Sn: Source score (based on publisher credibility)
- α, β, γ Tunable weights (default 0.4, 0.4, 0.2)

To ensure transparency, an **explanation engine** is included:

- Highlights legal phrases that influenced the summary
- Shows jurisdictional keywords used in classification
- Explains legal QA answers using citation traces

Gemini 1.5 enables token-wise attribution scoring, allowing backtracking to source lines.

Table 5: User Interface Features

Feature	Description
Document Upload Panel	Accepts DOC, PDF, TXT inputs
Summary Generator	Toggle between "Expert" and "Layman" mode
QA Panel	Ask document-related questions (chat interface)
News Insights Panel	Live legal news matched to document topic
Explanation Toggle (XAI)	Highlights key sections used in NLP outputs
Export Options	Summary export to PDF/CSV/Email

Table 6: Advantages over Existing Systems

Feature	Existing Systems (Kira, Harvey)	Our System
Dual pipeline (Docs + News)	✗	✓
Layman-friendly output	✗	✓
Multilingual support	Partial	✓(via Gemini 1.5)
Real-time integration	✗	✓
Explainable summaries	✗	✓

The proposed AI-driven legal document summarizer demonstrates several key advantages over existing systems such as Kira and Harvey. Table 5 represents the user interface patterns and Table 6 represents the advantages over existing systems. One of the most significant improvements is the implementation of a dual pipeline architecture, which combines both legal document analysis and real-time legal news aggregation. This enables users to not only comprehend static legal content but also to understand its relevance in the context of evolving legal developments—a feature notably absent in traditional systems. Another important advantage is the system's ability to produce layman-friendly outputs. While existing platforms typically deliver results suited for legal professionals, our system includes a toggle between "Expert" and "Layman" modes, allowing even non-experts to understand complex legal texts. This democratizes access to legal information and greatly expands the system's usability across different user groups. The system also excels in multilingual support, powered by the Gemini 1.5 language model. While some existing solutions offer partial multilingual capabilities, our model delivers full support across multiple languages, enabling legal analysis for international users and cross-border legal contexts. Real-time integration is another area where the proposed system surpasses its predecessors. Unlike Kira and Harvey, which work with static datasets, this system pulls live legal news that is contextually matched to the uploaded document. This real-time news insight ensures that users are always equipped with the latest updates relevant to their legal matters. Lastly, the platform embraces explainable AI (XAI) principles by offering explainable summaries. Users can activate the Explanation Toggle to highlight key sections of the original document that were used to generate summaries or answer queries. This transparency builds user trust and enhances legal defensibility, making the system particularly suitable for professional and compliance-heavy environments. Collectively, these features make the proposed system a robust, accessible, and future-ready tool for legal document analysis.

IV. Results and Discussion

This section presents a comprehensive performance analysis of the proposed AI-driven legal document analysis system. The evaluation framework was designed to assess the platform's efficiency, accuracy, speed, user interpretability, and adaptive learning capabilities. To validate the system's effectiveness, a benchmark dataset of 10,000 diverse legal documents including contracts, litigation briefs, compliance reports, statutes, and case laws—was employed for empirical testing. The primary function of the system is to deliver concise, reliable summaries of complex legal documents. The summarization model was evaluated based on human-annotated benchmarks by legal experts.

Table 7: Document Summarization Accuracy Metrics

Metric	Value
Accuracy	94.3%
Precision	93.7%
Recall	95.0%
F1 Score	94.3%

Metric	Value
ROUGE-L Score	0.89
BLEU Score (4-gram)	0.81

The summarization module as in table 7 demonstrated a 94.3% accuracy, outperforming conventional summarization systems (e.g., BART or GPT-3.5 with fine-tuning). The high recall rate of 95% indicates the model's robustness in capturing key content from verbose documents, including embedded clauses and conditional phrases. The ROUGE-L score of 0.89 and BLEU score of 0.81 confirm the model's linguistic alignment with expert-generated summaries, affirming both syntactic and semantic fidelity. The secondary core module was tasked with answering legal queries derived from document content. The model was fine-tuned using supervised question-answer pairs, validated by legal practitioners.

Table 8: Legal QA Module Evaluation

Metric	Value
Accuracy	91.7%
Precision (Entity-level)	90.8%
Recall	93.2%
F1 Score	92.0%
Mean Answer Latency (ms)	56.3

The legal QA module achieved 91.7% accuracy as in table 8, effectively identifying entities such as litigants, jurisdictions, statute names, and temporal clauses. The mean latency of 56.3ms confirms its near real-time responsiveness. Notably, the recall rate of 93.2% indicates strong comprehensiveness, ensuring key legal information is rarely omitted, which is critical for legal interpretation and dispute resolution. To validate the system's comparative advantage, it was benchmarked against legacy legal NLP models, including FinBERT-Legal, LexNLP, and a fine-tuned GPT-3.5 model.

Table 9: Comparative Performance with Baselines

Model	Summarization Accuracy	QA Accuracy	Average Latency (ms)
Proposed (Gemini)	94.3%	91.7%	56.3
FinBERT-Legal	87.1%	83.9%	143.5
LexNLP	79.4%	75.6%	198.3
GPT-3.5 (fine-tuned)	88.6%	85.7%	118.6

The proposed system as in table 9 significantly outperformed all baseline models. The Gemini-based model demonstrated as in table 9 superior linguistic contextualization and domain-specific disambiguation, especially in jurisdictional cues, embedded obligations, and exception clauses. While GPT-3.5 came close, its latency was double that of Gemini 1.5, making it unsuitable for high-speed legal workflows. Legal document types differ in structure, vocabulary density, and logical flow. Sector-wise breakdown was conducted to gauge performance across legal document genres.

Table 10: Sector-Wise Summarization Accuracy

Document Type	Accuracy (%)
Litigation Briefs	95.2
Statutory Documents	94.1
Commercial Contracts	93.6
Compliance Reports	94.7
Intellectual Property	92.8

The system excelled as in table 10 with litigation briefs (95.2%) and compliance reports (94.7%) due to their regular structure and clarity. Summarization of IP filings (92.8%) had slightly reduced performance owing to dense legal jargon and cross-referencing of prior claims. The consistent accuracy above 92% in all sectors shows the system's adaptability and broad domain expertise. The system aggregates legal news and contextualizes it to relevant jurisdictions and ongoing legal trends.

Table 11: News Aggregation and Classification Accuracy

Metric	Value
Jurisdictional Classification	92.4%
Duplicate News Filtering	96.1%

Metric	Value
Real-Time Refresh Interval	180s
Summary Quality (Expert Rating)	4.6 / 5

Jurisdiction tagging as in table 11 was successful in 92.4% of the news articles, driven by metadata parsing and location-sensitive named entity recognition (NER). The system updates the feed every 3 minutes, ensuring timely legal awareness. The duplicate filtering rate of 96.1% eliminates redundant information, while expert evaluations rated the news summaries an average of 4.6/5, indicating high clarity and usefulness. The system includes a built-in rationale generator that details the logic behind each summary or QA output.

Table 12: User Satisfaction with Explanation Features

Feature	Avg Rating (out of 5)
Legal Clause Highlighting	4.8
Summary Traceability	4.6
Entity Mapping Visualization	4.5
Query-Aware Answering	4.7

High satisfaction with legal clause highlighting (4.8) as in table 12 reflects the system's ability to visually isolate important segments. The entity visualization and traceability of summary elements further improve user trust, especially in professional settings such as courts or compliance review. This transparency is critical for audit trails and legal defensibility. User interactions (accept/reject summary, edit, re-query) are used to refine model predictions through reinforcement learning.

Table 13: Feature Importance Weight Adjustments (6-Month Period)

Feature	Initial Weight	Updated Weight
Legal Term Density	0.27	0.32
Entity Context Clarity	0.23	0.29
Temporal Clause Relevance	0.21	0.18
Jurisdictional Cues	0.29	0.21

The system learned to prioritize legal term density and entity clarity over jurisdictional signals, which had lower predictive utility in some contexts. This dynamic adjustment reflects the system's continuous improvement model, essential in the ever-evolving legal landscape where language patterns shift with court trends and legislation. The proposed Gemini 1.5-based legal analysis system excels in document summarization, legal QA, and jurisdiction-based news aggregation. It demonstrates a high degree of interpretability, real-time responsiveness, and consistent cross-domain accuracy. Its ability to adapt via user feedback ensures long-term sustainability in professional legal environments.

References

- Kanapala, A., Jannu, S., & Pamula, R. (2019). Passage-based text summarization for legal information retrieval. *Arabian Journal for Science and Engineering*, 44, 9159–9169. <https://doi.org/10.1007/s13369-019-03998-1>
- Sharma, S., Srivastava, S., Verma, P., & Jaiswal, P. (2023). A comprehensive analysis of Indian legal documents summarization techniques. *SN Computer Science*, 4, 614. <https://doi.org/10.1007/s42979-023-01983-y>
- Zhong, Y., & Litman, D. (2022). Computing and exploiting document structure to improve unsupervised extractive summarization of legal case decisions. In *Proceedings of the Natural Legal Language Processing Workshop 2022* (pp. 322–337). <https://doi.org/10.18653/v1/2022.nllp-1.30>
- Mullick, A., Nandy, A., Kapadnis, M., Basu, A., & Ghosh, K. (2022). An evaluation framework for legal document summarization. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference* (pp. 4747–4753). <https://aclanthology.org/2022.lrec-1.508/>
- Bhattacharya, P., Poddar, S., Rudra, K., Ghosh, S., Goyal, P., & Ganguly, N. (2021). Incorporating domain knowledge for extractive summarization of legal case documents. In *Proceedings of the 18th International Conference on Artificial Intelligence and Law* (pp. 22–31). <https://doi.org/10.1145/3462757.3466092>
- Elaraby, M., Zhong, Y., & Litman, D. (2023). Towards argument-aware abstractive summarization of long legal opinions with summary reranking. *arXiv*. <https://doi.org/10.48550/arXiv.2306.00672>
- Mahoney, C. J., Zhang, J., Huber-Fliflet, N., Purpura, S., & Kolcz, A. (2019). A framework for explainable text classification in legal document review. *arXiv*. <https://doi.org/10.48550/arXiv.1912.09501>
- Kore, R. C., Ray, P., Lade, P., & Nerurkar, A. (2020). Legal document summarization using NLP and ML techniques. *International Journal of Engineering and Computer Science*, 9(05), 25039–25046. <https://doi.org/10.18535/ijecs/v9i05.4488>
- Parikh, V., Mathur, V., Mehta, P., Agarwal, R., & Varghese, N. (2021). LawSum: A weakly supervised approach for Indian legal document summarization. *arXiv*. <https://doi.org/10.48550/arXiv.2110.01188>
- Chalkidis, I., Fergadiotis, M., Malakasiotis, P., Aletas, N., & Androutsopoulos, I. (2020). LEGAL-BERT: The muppets straight out of law school. *arXiv*. <https://doi.org/10.48550/arXiv.2010.02559>

11. Chen, Z., Ye, L., & Zhang, H. (2023). Enhancing LSTM and fusing articles of law for legal text summarization. In *ICONIP (14), Communications in Computer and Information Science* (Vol. 1968, pp. 110–124). https://doi.org/10.1007/978-3-031-28241-6_9
12. Bindal, P., Kumar, V., Bhatnagar, V., & Misra, A. (2023). Citation-based summarization of landmark judgments. In *Proceedings of the 20th International Conference on Natural Language Processing (ICON)* (pp. 588–593). <https://aclanthology.org/2023.icon-1.56>
13. Çano, E., & Morisio, M. (2017). Hybrid recommender systems: A systematic literature review. *Intelligent Data Analysis*, 21(6), 1487–1524. <https://doi.org/10.3233/IDA-163209>
14. Chalkidis, I., & Kampas, D. (2019). Deep learning in law: Early adaptation and legal word embeddings trained on large corpora. *Artificial Intelligence and Law*, 27(2), 171–198. <https://doi.org/10.1007/s10506-019-09253-2>
15. Bhattacharya, P., Poddar, S., Rudra, K., et al. (2021). Incorporating domain knowledge for extractive summarization of legal case documents. In *Proceedings of the 18th International Conference on Artificial Intelligence and Law* (pp. 22–31). <https://doi.org/10.1145/3462757.3466092>

