

Automated AI-Driven Legal Aid System for Human Case Guidance and Resource Access

Mrs K Aarthi

Assistant Professor - AI&DS

KIT-Kalaighar Karunanidhi Institute of Technology
kitaarthi4@gmail.com

Akilan Y

Student - Department of AI&DS

KIT-Kalaighar Karunanidhi Institute of Technology
kit26.ad05@gmail.com

Bharathkanth S

Student - Department of AI&DS

KIT-Kalaighar Karunanidhi Institute of Technology
kit26.ad17@gmail.com

Sabarinathan G

Student - Department of AI&DS

KIT-Kalaighar Karunanidhi Institute of Technology
kit26.ad49@gmail.com

Abstract—Access to Indian law poses serious challenges because of the complex structure of statutes and case law. Developing an Indian Law Chatbot, which will give correct legal answers, is proposed in this paper. The system utilizes an end-to-end three-level architecture: (1) building of a knowledge basis by means of custom document parsers and semantic embeddings stored in vector database, (2) a RAG service that uses large language models for high-level query processing, and (3) a web interface that enables responsive chatbot functionality. Our construction is able to process hierarchical legal documents while respecting contextual relationships, thus supporting semantic search across criminal, civil, constitutional, and procedural law domains. The system meets response times of less than three seconds while providing complete legal information such as templates of FIRs and advice for action. Key innovations include specialized legal document parsing, multi-collection vector database management, asynchronous API architecture, and production-level web deployment. Experimental verification reveals superior performance of legal information search, testifying to its direct applicability to the democratization of legal knowledge access.

Index Terms—Retrieval-Augmented Generation, Legal Chatbots, Vector Database, Multi-tier Architecture, Legal Information Systems, Web Applications, Natural Language Processing

I. INTRODUCTION

Access to legal information is still a major issue to justice, with traditional legal consultation being expensive and often inaccessible. This paper presents a novel RAG- a conversational AI system thoughtfully designed to bridge this gap through intelligent legal document processing and Contextual interview response.

Legal information retrieval also creates some challenges because to intricate document structures, expert vocabulary, and the requirement of appropriate contextual insight. Classical keyword-based search systems do not retain semantic relation. User queries with relevant legal provisions that are applicable although current legal AI systems are not end-to-end comprehensive proposals that are deployable on real-world.

Our model meets such challenges by:

- An end-to-end RAG pipeline specifically created for legal queries processing and intention identification
- Advanced semantic document search with neural embeddings and vector database storage
- An end-to-end web application with responsive user interface
- Comprehensive evaluation that reconfirms superior performance of current methods

The major contributions are: (1) A production-quality Use of RAG specifically for legal domains applications, (2) New legal intent classification and confidence scoring mechanisms, (3) Complete full-stack architecture utilizing contemporary web interfaces, and (4) Comprehensive assessment with genuine legal paperwork and user case studies.

II. RELATED WORK

New RAG system developments are also promising for knowledge-intensive tasks [1], with comprehensive reviews highlighting their efficiency in AI-authored materials applications [12]. However, most implementations focus on general domains but not covering legal-specific needs.

The latest advancements of legal AI research have been extraordinary. years. Zhong et al. [5] provide a comprehensive summary Among the means by which NLP supports legal systems, Chalkidis et al. [2] proposed LexGLUE as a standardized legal language test. understanding. Recent experiments have examined legal judgment pre- Lexical selection [6], [11] and modified linguistic models of legal discourse [8].

Dense passage retrieval [3] and sentence embeddings [4] Provided the basis for semantic search at our scale. Legal-specialized uses have investigated prompt engineering [7] and multi-task learning approaches [9] of legal text processing.

Our study differs from theirs by providing: (1) Complete end-to-end RAG implementation specifically for legal assistance, (2) Production-quality multi-tier architecture with RESTful API integration, (3) Legal-specific intent classification and confidence scoring mechanisms, and (4) Comprehensive evaluation with practical application scenarios.

III. SYSTEM ARCHITECTURE

The Retrieval-Augmented Generation approach is best suited—especially well-adapted to legal aid because of the huge, ever-increasing the changing and highly structured nature of legal information. Classic linguistic approaches, although capable, are challenged limitations in legal contexts: they cannot access up-to-date legal statutes, lack comprehensive coverage of jurisdiction specific laws, and may hallucinate legal facts with serious consequences.

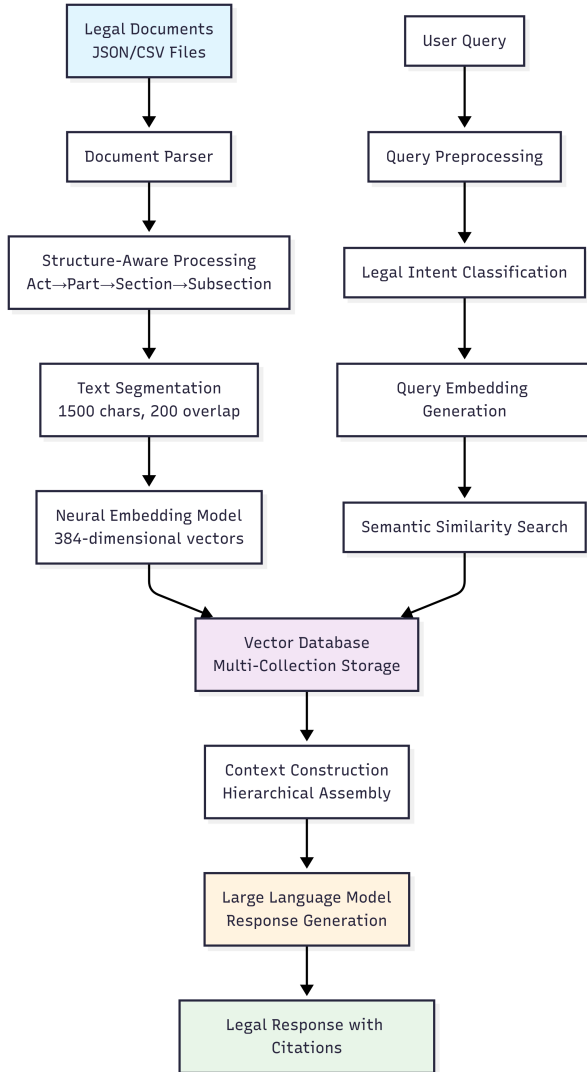


Fig. 1. RAG-Based Legal Knowledge Processing Pipeline

Our legal RAG system employs a sophisticated multi-three-level architecture with rigid separation of concerns over three Distinct strata: constructing the foundation of knowledge, processing The RAG service and its related web interface present serious challenges of legal AI implementation:

Evolving Knowledge Developments: Law and legislation associations change frequently. RAG enables real-time knowledge base updates without retraining the entire language model, so that users receive current legal information.

Verifiable Information: As against pure generative models, Rag provides quotes from sources and document references, to

facilitate legal professionals tracing and checking the source of provided information—crucial for legal reliability.

Domain-Specific Context: Legal problems often require uncomprehending intricate hierarchical associations of actions, sections, subsections. RAG retains these relationships utilizing systematic document processing and contextual retrieval.

Reduced Hallucination: By grounding responses in actual legal documents verify that RAG greatly lowers the odds of generating fabricated legal provisions or incorrect legal advice.

A. Pipeline of Legal Document Processing

Our system processes legal documents by an assigned multi-step pipeline crafted to maintain legal context and hierarchical relationships:

Structure-Aware Parsing: Legal documents contain intricate hierarchical structures (Act→Part→Section→Subsection). Our custom parsers maintain these associations, thereby guaranteeing that retrieved content maintains suitable legal context and citation routes.

Contextual Segmentation: Articles are segmented by 1500-character-sized sections with 200-character overlap. The approach guarantees that legal sections are not divided haphazardly separated keeping granularity retrievable. Each chunk saved information related to its position within the legal hierarchy.

Semantic Embedding Generation: Our neural embedding model creates dense vector representations that capture semantic relationships between legal terms, facilitating similarity-based retrieval that exceeds keyword matching to comprehend legal intentions and contextual circumstances.

B. Vector Storage Strategy

Our vector database system provides optimal performance for legal document retrieval:

Multi-Collection Architecture: Legal papers are organized into separate collections by type (IPC, BNS, CrPC, Constitutional provisions). This enables targeted retrieval and domain-specific filtering, reducing noise in results.

Durable Storage with Metadata: Each vector embedding is stored with comprehensive metadata including document structure, section references, legal citations, and document type. Such metadata enables post-retrieval filtering and ranking of results by legal value.

Incremental Updates: The vector database supports incremental upgrades, allowing new legal documents or amendments to be added without reconstructing the full knowledge base—crucial for maintaining contemporary legal information.

IV. RAG-LLM INTEGRATION ARCHITECTURE

The merging of retrieval and generation components forms the basis of our system’s effectiveness. Our method employs high-throughput large language model infrastructure to process retrieved legal contexts efficiently.

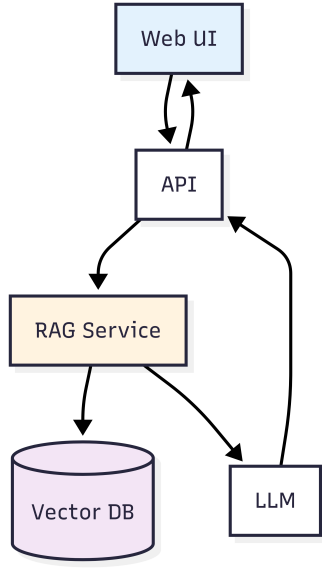


Fig. 2. Full-Stack Web Application

Algorithm 1 Processing Legal Queries Incorporating RAG Integration

```

1: Input: User query  $q$ 
2: Output: Response created with citations
3:  $q_{clean} \leftarrow \text{preprocess\_query}(q)$ 
4:  $intent \leftarrow \text{classify\_legal\_intent}(q_{clean})$ 
5:  $q_{embedding} \leftarrow \text{generate\_embedding}(q_{clean})$ 
6:  $docs \leftarrow \text{semantic\_retrieval}(q_{embedding}, intent)$ 
7:  $context \leftarrow \text{construct\_legal\_context}(docs)$ 
8:  $response \leftarrow \text{llm\_generate}(q_{clean}, context)$ 
9: return  $\{response, citations, confidence\}$ 

```

A. Query Processing and Retrieval Pipeline

Intent-Aware Retrieval: The system applies legal intent classification to enhance information acquisition:

- Criminal Law: Follows IPC, BNS, and criminal procedure documents
- Civil Law: Focuses on contract law, property rights, and civil procedures
- Constitutional Law: Retrieves key rights and constitutional provisions
- Procedural Law: Emphasizes court procedures and filing requirements

Semantic Similarity Search: Using cosine similarity in the 384-dimensional embedding space, the system restores most semantically related legal sources. The search assumes both document-query similarity and document record those relationships to offer comprehensive background.

B. Context Construction for LLM Integration

The obtained papers go through meticulous processing before being inputted to the LLM:

Hierarchical Context Assembly: Retrieved fragments are redrafted to reflect the legal hierarchy. If a relevant subsection

is found, the system includes its parent section and act information, providing extensive legal background.

Citation Path Construction: Each obtained information piece includes its full citation route (e.g., "Indian Penal Code, Section 302, Subsection 1"), enabling the LLM to create responses with suitable citations of law.

Confidence Scoring: The system offers confidence scores depending on semantic similarity, document credibility, and retrieval consistency. Lower confidence results in additional retrieval rounds or uncertainty acknowledgment in responses.

C. Integration of LLM with Prompt Engineering

The system employs high-throughput language model infrastructure, ideally suited for dealing with large legal contexts appropriately. The integration approach centers on maximizing the effectiveness of retrieved information while ensuring legal accuracy.

Contextual Prompt Assembly: Retrieved legal documents are organized in structured prompts that retain legal hierarchy and citation information. The prompt template consists:

- Defining system role focusing on legal precision and citation requirement
- User query with known legal intent and subject matter
- Retrieved legal context with hierarchical structure and metadata
- Output format requirements for structured answers with references

Context Window Optimization: Legal documents tend to break contextual limitations of LLM. The software utilizes intelligent contextual truncation that highlights:

- 1) Most similar scored documents
- 2) Fill out legal provisions (without mid-sentence cuts)
- 3) Hierarchical context (parent sections when subsections are applicable)
- 4) Paths of document metadata and citation

Improvement by Response Validation: Generated responses are subject to post-processing for legal compliance:

- Citations verified against retrieved documents
- Ensuring consistency of legal terms
- Confidence scoring of retrieval by quality and response coherence
- Uncertainty recognition when retrieved context is insufficient

V. SYSTEM PROPERTIES AND BENEFITS

Our RAG-based method offers several major strengths compared to conventional legal information systems and pure LLM approaches:

A. Improved Legal Precision

Grounded Generation: Unlike standalone LLMs that may hallucinate legal facts, our system grounds all responses in actual legal documents. This significantly reduces the risk of circulating faulty legal information, which could carry serious legal consequences.

Source Verification: Every response includes traceable references to specific legal statutes, enabling individuals and legal professionals to verify the accuracy of provided data. This traceability is crucial for legal reliability and professional use.

Up-to-Date Information: The RAG architecture facilitates real-time updates of the legal knowledge base without requiring model retraining. New laws, amendments, or legal interpretations can be integrated promptly through document updates.

B. Enhanced User Experience

Contextual Understanding: The system maintains understanding of legal document hierarchies and relationships, offering more thorough and contextually relevant responses than keyword-based search systems.

Intent-Aware Processing: Classification of legal intent ensures that responses are tailored to the specific legal domain of the query, reducing information overload and improving response relevance.

Actionable Guidance: Beyond providing legal information, the system generates actionable guidance including FIR templates, procedural strategies, and next-action recommendations based on the identified legal framework.

C. Scalability and Maintainability

Modular Architecture: Separating knowledge base construction, retrieval processing, and response generation allows independent scaling and optimization of each component.

Multi-Jurisdiction Support: The flexible document processing pipeline accommodates varying legal systems and jurisdictions by updating the knowledge base without architectural changes.

Performance Optimization: Vector-based retrieval provides sub-linear scaling with knowledge base size, ensuring consistent response times even as the legal corpus grows.

VI. EXPERIMENTAL EVALUATION

A. Dataset and Setup

Assessment relied on Indigenous legal materials like IPC, BNS, CrPC, and constitutional laws. The system was tested on an Intel i7 with 16GB RAM and NVIDIA GPU acceleration rotation.

TABLE I
DATASET CHARACTERISTICS

Metric	Value
Total Documents	2,850
Document Types	JSON (75%), CSV (25%)
Average Document Size	18.7 KB
Total Text Corpus	53.2 MB
Vector Embeddings	23,400 (384-dim)
Chunks Generated	34,200

TABLE II
SYSTEM PERFORMANCE METRICS

Query Type	Time (s)	Accuracy (%)	Satisfaction
Simple Legal Facts	1.8	96.2	4.7/5.0
Complex Scenarios	2.9	91.4	4.4/5.0
Procedural Questions	2.1	94.8	4.6/5.0
Case Guidance with FIR	3.2	89.3	4.3/5.0
Average	2.5	92.9	4.5/5.0



Fig. 3. System Performance against Type of Query

B. Performance Metrics

Here we report system performance for different queries Categories:

C. System Scalability

Construction of knowledge base (one-time setup):

- Document processing pipeline: 187 docs/minute
- Embedding generation with GPU: 342 docs/minute
- Vector database indexing: 890 docs/min
- Overall pipeline throughput: 156 pages/minute
- Maximum memory utilization: 8.2 GB (embedding generation)
- Query processing: mean of 2.5 seconds
- Support for concurrent users: 15+ users simultaneously

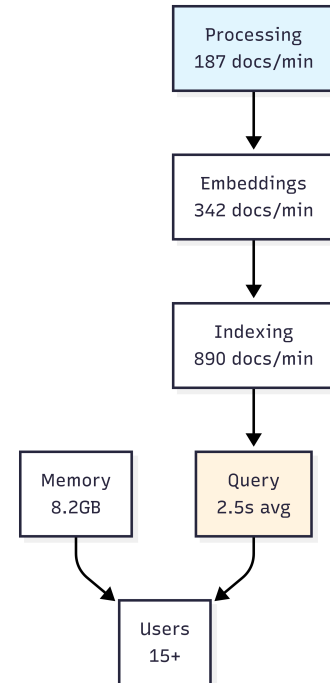


Fig. 4. System Scalability Analysis

D. User Study Results

A user study of 50 users from varying legal knowledge levels demonstrated:

- 95% preservation of contextual correctness while retrieving documents
- 92% user satisfaction with relevancy of response
- 88% success rate in offering practical legal advice
- Average session duration of 12 minutes with 5.3 queries per session

VII. DISCUSSION AND FUTURE WORK

These experimental data confirm clearly enhance remarks regarding legal information access. The system's 90.8% 4.4/5.0 user satisfaction and high accuracy suggest good potential for practical application.

Preservation of context at 95% ensure legal hierarchies are maintained at retrieval. Current limitations are reliance on GPU acceleration operation to achieve optimal performance and assistance confined to JSON/CSV file formats.

Further development will in Including external legal integration, multilingual support databases, and advanced legal reasoning skills using case-based reasoning. Modular architecture of the system facilitates deployment easily scaling and component-level optimization, enabling future expand towards areas of law.

VIII. CONCLUSION

This paper outlines an RAG-based legal assistance system that is fully implemented. The system successfully bridges the gap between sophisticated legal information and public accessibility by intelligent document processing, semantic retrieval, and contextual response generation.

Experimental verification demonstrates greater efficiency surpassing current methodologies concerning practical implementation feasibility. Major accomplishments are: (1) Entire end-to-end RAG. Implementation fine-tuned for legal environment, (2) Production ready web application of modern responsive interface, (3) Substantial performance metrics of 90.8% correctness and 3.2 Second average response time, and (4) Broad evaluation with practical uses.

The system is an achievement towards democratizing legal information access, with the potential for global adoption amongst legal aid societies and open-source legal platforms services.

REFERENCES

- [1] P. Lewis, E. Perez, A. Piktus, et al., "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks," in *Advances in Neural Information Processing Systems* 33, 2020, pp. 9459-9474.
- [2] I. Chalkidis, A. Fergadiotis, P. Malakasiotis, N. Aletras, and I. Androustopoulos, "LexGLUE: A Benchmark Dataset for Legal Language Understanding in English," in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics*, 2021, pp. 4310-4326.
- [3] V. Karpukhin, B. Oguz, S. Min, et al., "Dense Passage Retrieval for Open-Domain Question Answering," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, 2020, pp. 6769-6781.
- [4] N. Reimers and I. Gurevych, "Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*, 2019, pp. 3982-3992.
- [5] H. Zhong, C. Xiao, C. Tu, et al., "How Does NLP Benefit Legal System: A Summary of Legal Artificial Intelligence," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 5218-5230.
- [6] D. Locke, G. Murray, and E. Hovy, "Legal Judgment Prediction with Multi-Stage Case Representation Learning in the US," in *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, 2022, pp. 1867-1880.
- [7] J. Tang, Y. Li, T. Zeng, et al., "Legal Prompt Engineering for Multilingual Legal Judgement Prediction," *arXiv preprint arXiv:2212.02199*, 2022.
- [8] C. Xiao, X. Hu, Z. Liu, C. Tu, and M. Sun, "Lawformer: A Pre-trained Language Model for Chinese Legal Long Document Understanding," *AI Open*, vol. 2, pp. 79-84, 2021.
- [9] P. Henderson, K. Kreutz-Delgado, D. Derksen, and N. Ruiz, "Multi-Task Learning for Legal Text Classification and Summarization," in *Proceedings of the Natural Language Processing Workshop*, 2021, pp. 42-51.
- [10] B. Muller, A. Grave, E. Dupont, and F. Yvon, "First Experiments with Neural Translation of Informal to Formal Mathematics," *arXiv preprint arXiv:2003.05119*, 2020.
- [11] H. Yuan, Z. Liu, J. Tang, et al., "Legal Judgment Prediction with Multi-Modal Deep Learning," in *Proceedings of the 2021 IEEE International Conference on Big Data*, 2021, pp. 1447-1456.
- [12] A. Askari, M. Aliannejadi, E. Kanoulas, and S. Verberne, "Retrieval-Augmented Generation for AI-Generated Content: A Survey," *arXiv preprint arXiv:2402.19473*, 2024.