# Bridging Legal Knowledge and AI: Retrieval-Augmented Generation with Vector Stores, Knowledge Graphs, and Hierarchical Non-negative Matrix Factorization

Ryan C. Barron
barron@lanl.gov
Theoretical Division,
Los Alamos National Laboratory
Los Alamos, New Mexico, USA
CSEE, UMBC
Baltimore, Maryland, USA

Maksim E. Eren
maksim@lanl.gov
Information Systems and Modeling,
Los Alamos National Laboratory
Los Alamos, New Mexico, USA

Olga M. Serafimova
OMSerafimova@hollandhart.com
Holland & Hart LLP
Santa Fe, New Mexico, USA

Cynthia Matuszek
cmat@umbc.edu
CSEE, UMBC
Baltimore, Maryland, USA

Boian S. Alexandrov
boian@lanl.gov
Theoretical Division,
Los Alamos National Laboratory
Los Alamos, New Mexico, USA

## Abstract

Agentic Generative AI, powered by Large Language Models (LLMs) and enhanced with Retrieval-Augmented Generation (RAG), Knowledge Graphs (KGs), and Vector Stores (VSs), represents a transformative technology applicable across specialized domains such as legal systems, research, recommender systems, cybersecurity, and global security, including proliferation research. This technology excels at inferring relationships within vast unstructured or semi-structured datasets. The legal domain we focus on here comprises inherently complex data characterized by extensive, interrelated, and semi-structured knowledge systems with complex relations. It comprises constitutions, statutes, regulations, and case law. Extracting insights and navigating the intricate networks of legal documents and their relations is crucial for effective legal research and decision-making. Here, we introduce a generative AI system, a jurisdiction-specific legal information retrieval that integrates RAG, VS, and KG, constructed via Hierarchical Non-Negative Matrix Factorization (HNMFk), to enhance information retrieval and AI reasoning and minimize hallucinations. In the legal system, these technologies empower AI agents to identify and analyze complex connections among cases, statutes, and legal precedents, uncovering hidden relationships and predicting legal trends—challenging tasks essential for ensuring justice and improving operational efficiency. Our system employs web scraping techniques to systematically collect legal texts, such as statutes, constitutional provisions, and case law, from publicly accessible platforms like Justia. It bridges the gap between traditional keyword-based searches and contextual understanding by leveraging advanced semantic representations, hierarchical relationships, and latent topic discovery. This approach is demonstrated in legal document clustering, summarization, and cross-referencing tasks. The framework marks a significant step toward augmenting legal research with scalable, interpretable, and accurate retrieval methods for semi-structured data, advancing the intersection of computational law and artificial intelligence.

## CCS Concepts

• **Computing methodologies → Natural language processing**; *Reasoning about belief and knowledge*; *Ontology engineering*; • **Applied computing → Law**; • **Information systems → Document topic models**; **Content analysis and feature selection**; *Ontologies*; *Document collection models*.

## Keywords

law, legal knowledge, nmf, topic labeling, llm, chain of thought, prompt tuning, information retrieval

## 1 Introduction

The legal domain is uniquely complex, encompassing constitutions, statutes, court rules, regulations, ordinances, and case law. Each source follows distinct structural logics, while constitutions and statutes typically feature hierarchical formatting, case law often appears as long-form unstructured prose. As legal documents are continuously produced, legal professionals and researchers increasingly require computational tools that can go beyond basic keyword-based searches and deliver contextually meaningful, explainable insights. However, traditional legal information retrieval

approaches, such as Boolean logic [52] and lexical matching techniques like TF-IDF [7], fall short by missing subtle semantic overlap and deeper contextual dependencies for legal reasoning.

RAG has recently emerged as a compelling paradigm integrating traditional information retrieval's strengths with LLMs' generative capabilities. In the legal domain, RAG systems retrieve relevant documents or knowledge fragments and synthesize them into coherent, grounded answers. This framework mitigates LLM-specific limitations such as hallucinations [26, 33] by anchoring outputs in authoritative source texts, enhancing performance on tasks like statute interpretation, legal summarization, and case retrieval. Yet to fully realize these benefits, such systems must be supported by infrastructures that can represent, retrieve, and reason over vast, heterogeneous legal corpora.

This work introduces a domain-specific RAG system that integrates three core components to meet these challenges:

- **VS**: Legal texts are embedded into dense semantic vector spaces using pre-trained models (e.g., BERT [17], GPT [8, 44]), enabling retrieval based on contextual similarity rather than surface keyword overlap.
- **KG**: Legal entities—such as cases, statutes, and doctrines—are interlinked through explicit relationships (citations, shared principles), forming a Neo4j-based graph that supports structured navigation and inference over legal concepts.
- **HNMFk** : To uncover latent semantic structure, we apply hierarchical NMF (HNMFk) to legal corpora, producing interpretable topic clusters that support document classification, trend discovery, and integration into the KG as abstract semantic anchors.

The resulting architecture fuses high-recall semantic search with logical reasoning via structured knowledge, enabling accurate, explainable responses across a range of legal tasks. Our contributions to this work consist of the following:

- **Jurisdiction-Specific legal information retrieval**:
  A pipeline tailored to New Mexico's legal documents: statutes, constitution, and case law, demonstrating state-level legal domain adaptation.
- **Legal RAG framework tailored to New Mexico's legal documents**: A modular system combining RAG, VS search, KGs, and LLMs for grounded, explainable legal responses.
- **The first knowledge graph of New Mexico's legal documents with latent topics**: Our HNMFk extracts and classifies the latent topics of New Mexico's legal documents to construct a Knowledge Graph of metadata and citation links for semantic reasoning.

## 2 Relevant Work

This section reviews contributions across RAG domains, semantic search using vector embeddings, knowledge graph construction, NMF, and legal information systems.

### 2.1 Non-Negative Matrix Factorization for Pattern Discovery

NMF is a dimensionality reduction technique used to uncover latent patterns in data. [32] analyzed NMF as an interpretable method for extracting features and topics from large datasets, explicitly highlighting its ability to identify meaningful and non-overlapping components. Building on this work, [23] introduced sparseness constraints for better interpretability, improving applications of NMF in real-world scenarios through more focused feature selection. In the legal domain, NMF has been valuable for analyzing complex textual data, such as case law and statutes, and assisting with topic discovery and clustering. For instance, [9] applied NMF to legal documents to extract latent topics and visualize relationships, demonstrating how NMF's interpretable structure aids researchers in identifying underlying topics not readily apparent in raw text. They even applied NMF hierarchically to find fine-grained topics. More recently, [34] proposed a guided semi-supervised NMF approach for topic discovery in legal documents, using domain knowledge to steer factorization and ensure the extracted topics remain highly relevant. This semi-supervised extension bridges the gap between fully automated unsupervised techniques and expert-driven analysis. NMF's use in legal contexts is significant, particularly for interpretable results on large textual datasets.

### 2.2 Retrieval-Augmented Generation

RAG has emerged as a foundational approach for improving AI systems across various domains, including law. Lewis et al. [33] introduced a framework that dynamically retrieves relevant documents to ground generative outputs, achieving notable gains in accuracy. Building on this idea, Guu et al. [22] proposed a retrieval-augmented pretraining method that integrates external knowledge for improved downstream task performance, while Izacard and Grave [25] demonstrated the effectiveness of retrieval in open-domain question answering. These advances lay the groundwork for applying RAG to the legal sector, where the method's ability to ground LLMs in authoritative texts reduces hallucinations and increases accuracy in tasks including law retrieval, statutory reasoning, and judgment prediction. Notable examples include CBR-RAG, which incorporates Case-Based Reasoning to structure retrieval for legal QA [54], and LegalBench-RAG. This benchmark suite tailors evaluation metrics to the demands of legal information synthesis [42]. Parallel work has demonstrated RAG's capabilities in other domains, such as malware data analysis, by combining embeddings, KGs, and NMF [4]. Other works show how LLMs can dynamically decide when and what to retrieve to improve legal reasoning [31] and how multi-step legal judgment prediction can benefit from iterative retrieval and generation [50], further demonstrating the effectiveness of combining retrieval strategies with LLMs.

### 2.3 Semantic Search with Vector Embeddings

Semantic search operates on dense vectors to find the deeper semantic relationships in texts, going beyond keyword matching to proper context retrieval [22, 28]. This search is especially valuable in the legal domain, where queries often demand conceptual rather than surface-level understanding. Domain-specific pretraining has become increasingly important, as in LEGAL-BERT [12], outperforming general-purpose models by better capturing legal language nuances. Benchmarks like LeCaRD [36] show the effectiveness of dense retrievers, such as SBERT, in legal case retrieval of Chinese law, and the abilities of hybrid approaches that integrate lexical

and dense retrieval methods. Hierarchical transformer architectures [24] and long-context models like Longformer [5] further address the challenges of lengthy legal documents. In contrast, citation-driven approaches such as SPECTER [15] leverage metadata to improve retrieval. Furthermore, work on neighborhood contrastive learning for scientific document representations (SciNCL) [41] demonstrates how controlled nearest neighbor sampling over citation graph embeddings can provide continuous similarity. This strategy could also inform citation-based retrieval improvements in legal domains. Challenges remain for scaling, explainability, and heterogeneous data sources in real-world legal workflows.

## 2.4 Knowledge Graphs in Legal Research

Although the legal knowledge is jurisdiction-specific, various KGs have been widely adopted across various legal systems in legal research, offering structured representations that support case retrieval, citation analysis, and question answering. In the United States, researchers already have used KGs for legal knowledge extraction [46], researchers also built various domain-specific graphs to model regulatory frameworks [19], and explored their potential in addressing concerns of trust, privacy, and transparency in legal AI systems [10].

In other jurisdictions, similar KG methodologies have proven valuable. In India, KGs have been applied to recommend comparable legal cases [18]. For example, in China, they have been used to link statutes and case law to improve retrieval and to enhance case law search in Chinese courts [6, 55]. In Germany, Milz et al. conducted a structural analysis of a legal citation network using KGs to map citation patterns in German court decisions [37].

Beyond case retrieval, KGs have also been incorporated into knowledge-aware machine reading systems designed for legal question answering [21]. These international efforts highlight the versatility of KGs in capturing the intricacies of legal systems. In this work, we build on these foundations within the context of New Mexico statutes and case law, demonstrating the applicability of KGs to state-level U.S. legal research and retrieval.

## 2.5 Legal Systems and Case Retrieval

Legal information systems have evolved rapidly with the advent of neural architectures and hybrid retrieval pipelines, enabling precise tasks such as precedent retrieval, statute matching, and judgment prediction [11, 14, 45]. Benchmarks like LeCaRD [36] and LexGLUE [13], along with LEGAL-BERT-based systems [12], have demonstrated the capability of these neural methods to improve accuracy in analyzing legal corpora. In particular, researchers have leveraged structured reasoning with transformers and graph representations to link statutes and precedents, as evidenced in the COLIEE competition [43]. Nonetheless, several limitations remain: data scarcity and jurisdictional bias continue to restrict the generalizability of such models. At the same time, resource-intensive retrievers like BERT-based cross-encoders [56] have challenges scaling to large-scale legal databases. Earlier works in juris-informatics have highlighted the potential of automating legal reasoning and document analysis [2], laying the groundwork for modern approaches that fuse knowledge graphs, transformers, vector stores, and agent-oriented RAG pipelines to deliver more explainable and efficient legal workflows.

## 3 Methods

Laws and their interpretations are limited in application to their respective jurisdictions. This examines the State of New Mexico, using the available Supreme Court and Court of Appeals case law, the state constitution, and state statutes as the primary data resource.

## 3.1 Overview and Interconnect

Our jurisdiction-specific AI system integrates RAG, VS, KG, HN-MFk, and LLMs in a modular, agentic architecture for New Mexico's legal information retrieval and reasoning, shown in Figure 1.
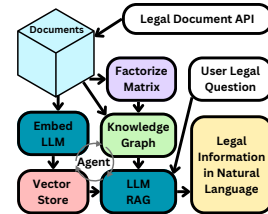


**Figure 1: System Overview**

Queries enter the system through a natural language interface. The input is embedded and used to perform two operations: (1) a semantic retrieval against a vector database containing constitutional provisions, statutes, and segmented case law, and (2) a traversal of the Neo4j-based legal knowledge graph containing latent HNMFk topics, citations, and structural metadata from the legal corpus.

An orchestration layer evaluates both results, aligning retrieved documents and KG nodes with the user's query context. The agent performs reasoning across case-to-case, statute-to-case, or doctrine-level connections if a relevant KG subgraph is identified. Otherwise, it defaults to the highest-confidence passages from vector search.

The selected materials are passed to an LLM agent in both cases, where a final response is generated, grounded in precise citations. If the user requests clarification or follow-up, the system maintains a short-term state and re-evaluates context via VS and KG pathways.

## 3.2 Data Collection

To compile a research-oriented corpus of legal documents, we used Justia's publicly accessible resources in full compliance with their terms of service and robots.txt guidelines. A restricted, responsible crawler was developed to access only permitted content such as statutes, constitutional provisions, and public case law, while respecting rate limits and avoiding disallowed sections. The scraper incorporates request delays, monitors HTTP status codes, and applies exponential backoff to ensure ethical data collection.

Starting from Justia's landing pages, the crawler identifies valid navigation links, resolves absolute URLs, and recursively follows permitted content. It extracts metadata and full text for statutes, constitutional provisions, citations, and judicial opinions for New Mexico's Supreme Court and Court of Appeals.

Compliance and transparency are maintained through:

- **Logging** to track URLs, prevent duplication, and resume interrupted crawls.
- **Rate limiting** to avoid overloading Justia's servers.
- **Respect for technical restrictions** in robots.txt / TOS.
- **Strictly non-commercial use**, limited to scholarly research.

These measures preserve ethical and legally compliant data collection that is aligned with Justia's content usage policies.

## 3.3 Dimension Reduction in Legal Texts

Legal documents- constitutions, statutes, and case law- are traditionally organized into chapters, articles, and sections, but these structures do not always match the latent patterns revealed through factorization. Each document type is analyzed separately using non-negative tensor and matrix factorization to explore these hidden relationships. A TF-IDF matrix, $\mathbf{X}$, is first constructed from the cleaned corpus. Constitutional provisions, statutory clauses, and case law paragraphs form the units of analysis for clustering.

In this study, we use **Tensor Extraction of Latent Features** (**T-ELF**)[1] [20], combined with automatic model selection, to decompose $\mathbf{X}$ into coherent topic H-clusters. **T-ELF** efficiently identifies latent topics, grouping constitutional provisions around themes such as "separation of powers" and clustering statutes and case law based on regulatory objectives and recurring legal doctrines, respectively. The optimal number of clusters, $k$, is determined using silhouette scores and is accelerated through a binary search [3].

*3.3.1 Application of Non-Negative Matrix Factorization to Legal Texts.* NMF approximates the matrix $\mathbf{X}$ with two non-negative matrices, $\mathbf{W}$ and $\mathbf{H}$ such that $\mathbf{X} \approx \mathbf{W} \cdot \mathbf{H}$, where $\mathbf{W}$ describes how terms distribute over topics, and $\mathbf{H}$ describes how these topics distribute across documents. Constitutional articles and sections reveal underlying governance or civil rights themes; statute clauses highlight regulatory objectives, and segmented judicial opinions expose common doctrines and legal reasoning patterns.

*3.3.2 Automatic Model Determination Using NMFk.* A central challenge in applying NMF is selecting the best number of latent features ($k$). We use *NMFk* [1], which combines clustering stability with reconstruction accuracy. Bootstrap resampling generates slightly perturbed versions of the original matrix, and repeated decompositions measure how consistently clusters form. Silhouette scores help ensure cohesive, well-separated topics, while reconstruction error verifies that the model effectively captures patterns in the original data.

By adapting this NMF approach to a hierarchical setting, legal texts can be organized into a tree-like structure. Constitutions may be segmented into articles and sections, statutes into chapters and clauses, and case law into layered precedents and sub-issues. This hierarchical perspective enhances the discovery of latent relationships at multiple levels of granularity, facilitating deeper analyses of large-scale legal corpora.

## 3.4 Knowledge Graph

Features derived from **T-ELF** and document metadata are transformed into a series of head, entity, and tail relations, forming directional triplets integrated into a Neo4j KG [38].

In the legal context, the KG incorporates metadata and latent features extracted from constitutions, statutes, and case law. The primary nodes in the graph represent legal documents, including constitutional provisions, statutory sections, and judicial opinions. These nodes are enriched with metadata such as titles, hierarchical

[1]https://github.com/lanl/T-ELF

identifiers (e.g., chapter and section numbers), jurisdiction, court names, decision dates, and topics derived from latent features.
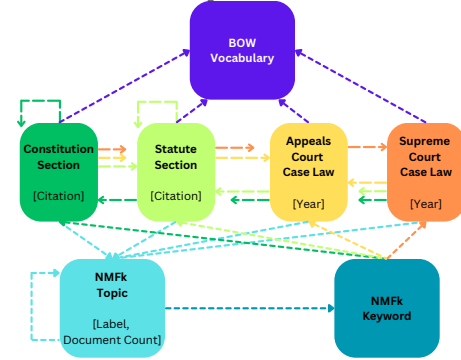


**Figure 2: Knowledge Graph Schema, with the primary identifier in bold and attributes in brackets.**

Edges in the KG establish relationships between nodes to represent the interconnected nature of legal documents. For instance:

- **Constitutional Nodes:** Linked to statutory provisions or judicial interpretations that reference or rely on specific constitutional clauses.
- **Statutory Nodes:** Connected to cases interpreting the statute or related provisions within the same legislative framework.
- **Case Law Nodes:** Interlinked based on shared topics, common legal principles, or hierarchical relationships in appellate decisions.

This graph structure enables the RAG system to query and retrieve legal documents based on semantic similarity and explicit relationships. For example, a query about "due process" might retrieve the relevant constitutional clause, cases that discuss its interpretation, and statutory provisions impacted by those rulings. By combining metadata and latent features, the KG supports advanced reasoning and logical traversal, enhancing the precision and depth of legal information retrieval.

The Bag-of-Words (BOW) vocabulary, i.e., the unique unigrams of the corpus, was generated through standard NLP processing, measuring their frequency of tokens in documents. The most frequent 50 unigrams in a topic, rows of $\mathbf{W}$, define the top keywords in this topic. This included stop word removal, special character removal, etc.

## 3.5 Vector Store

A vector database was implemented to manage and index legal documents, improving the RAG process for legal research. Using Milvus [49], the system stores vectorized representations of constitutions, statutes, and case law, treating each document type uniquely to ensure contextually relevant retrieval.

Constitutional provisions are split into paragraphs, each with a unique ID, and vectorized using OpenAI's text-embedding-ada-002 [40] model for granular semantic search. Statutes are divided into sections or clauses, with metadata like chapter and section titles added for precise retrieval. Case law, being unstructured, is chunked into logical units, preserving flow and indexed with metadata such as case name and citation.

Based on the query's focus, the RAG application queries the database to retrieve relevant fragments—constitutional paragraphs,

statutory clauses, or case law sections. Retrieved text is synthesized into responses, allowing the LLM to cite specific paragraphs or clauses, ensuring traceability and accuracy. The system leverages a connected knowledge graph to explore related amendments, judicial interpretations, or precedent cases for additional context. This integration provides a robust, comprehensive retrieval process tailored to the complexities of legal research.

## 3.6 Evaluation Metrics

In evaluating AI-generated responses against reference texts, metrics assess content quality, semantic similarity, and factual consistency. ROUGE-L (Recall-Oriented Understudy for Gisting Evaluation) evaluates lexical overlap of the longest common subsequence to assess structural similarity [35]. Natural Language Inference (NLI) entailment quantifies whether the generated text is logically entailed by the reference, based on pretrained models like BART [53]. SummaC evaluates factual consistency by aggregating entailment scores across sentence pairs between the hypothesis and source [30]. FactCC is another entailment-based model that identifies factual inconsistencies in summarization outputs by fine-tuning on labeled correctness data [29]. Additionally, we leverage human evaluation on responses for a small subset of the questions.
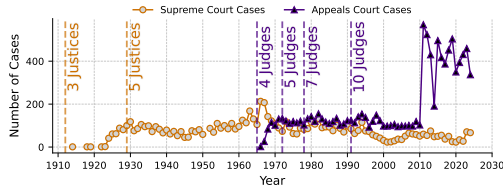
**Figure 3: New Mexico Supreme/Appeals case counts per year.**

## 4 Results

This section presents the resulting legal texts and their hierarchical decomposition, offering a detailed breakdown of the sections and cases within each document type. The following results illustrate the scope and depth of the collected data, providing a foundation for further exploration of trends and patterns across the legal corpus.

### 4.1 Dataset

After collecting and structuring the data, the legal documents were organized into 265 Constitutional provisions, 28,251 Statutory sections, and Case Law comprising 10,072 Court of Appeals cases and 5,727 Supreme Court cases. Figure 3 shows the trends in these Supreme Court and Court of Appeals cases over the years, as available from Justia [27], which also includes the expansions and creations of the courts themselves.

### 4.2 Decomposition

The four component data types were decomposed hierarchically with NMFk. The vocabulary for building the TF-IDF matrix was collected using specific parameters for each part. For the Constitution vocabulary, the minimum token document frequency (DF) was set to 5 documents, with a maximum token DF of 80% of the corpus, resulting in a final size of 416 tokens. For the Statutes vocabulary, the minimum token DF was set to 30 documents, with a maximum token DF of 70% of the corpus, yielding 7,508 tokens.

The Court of Appeals vocabulary used a minimum token DF of 50 documents (cases) and a maximum token DF of 70% of the corpus, resulting in 10,189 tokens. Last, the Supreme vocabulary employed a minimum token DF of 50 papers and a maximum token DF of 70% of the corpus, with a final size of 8,425 tokens. The maximum depth was set to 2, of which only the constitution sections did not reach due to the limited number of sections. The limiting factor of further decompositions was 100 documents, so if the preceding cluster had ≥ 100, the cluster would decompose and stop otherwise. As seen in Figure 4a, the largest H-cluster is from cluster 4 in dark blue with 49 sections, with the fewest documents in yellow with six sections in H-cluster 0. The other three decompositions can be examined in the larger Figure 4, where the court cases, there are 10 leaf H-clusters in the constitution decomposition, 985 leaf H-clusters in the statutes decomposition, 420 leaf H-clusters in the Court of Appeals cases, and 248 leaf H-clusters in the Supreme Court cases. From the methods in [51], each H-cluster throughout the decomposition hierarchies has LLM-generated labels for ease of reference and quick insight. Labels for the first decomposition depth for the constitution can be observed in Table 1, the statutes in Table 2, the Supreme court in Table 4, and the appeals in Table 3. The depth-0 H-clustering corresponds to the first ring radial from the center totals in Figure 4.

**Table 1: NM Constitutional Depth-0 H-Clusters**

| # | Label |
|---|---|
| 0 | Irrigation and Water Resource Management Principles |
| 1 | Regulation of Private Sector Influence on Public Schools and Education Services |
| 2 | Branches of Government Structure and Functionality |
| 3 | Education Funding for New Mexico's Educational Institutions |
| 4 | Legislative Proceedings and Lawmaking Activities Enacted During Sessions |
| 5 | Territorial Transition: Constitutional Ratification and Statehood Provisions |
| 6 | Governance, Land, and Taxation Framework |
| 7 | Municipal and County Financial Obligations and Liabilities Management |
| 8 | Judicial Power and Jurisdictional Frameworks |
| 9 | Public Service Election Governance Structure and Processes |

**Table 2: NM Statutory Depth-0 H-Clusters**

| # | Label |
|---|---|
| 0 | Municipal Court Civil Cases Involving Children's Rights |
| 1 | Public Education Infrastructure Management Systems |
| 2 | Criminal Codes, Local Governance Boundaries, Licensing Rules |
| 3 | Comprehensive Emergency Health and Human Services Response Framework |
| 4 | Taxation and Revenue Collection Oversight |
| 5 | Military Decorations, Licensing Procedures, Governance Boards |
| 6 | Regulatory Insurance Contract Law and Policy Analysis |
| 7 | Public Obligations Financing and Project Bonds Issuance |
| 8 | State Government Investment Grants for Education and Economic Development |
| 9 | Irrigation and Water Rights Regulations |
| 10 | Military Honors, Discrimination Penalties, and Trust Authority |
| 11 | Corporation Governance Framework and Regulatory Compliance |

## 4.3 Knowledge Graph

The four data parts, 265 constitutional provisions, 28,251 Statute sections, 5,727 Supreme Court cases, and 10,072 Court of Appeals cases, were inserted into the neo4j [38] knowledge graph. The graph's number of nodes and edges can be seen in Table 5, where edge counts are where the triplet's tail originates with the row's node. The legal citations were collected by iterating the text of each case or section into chat-gpt-3.5-turbo with the following prompt: "You are an expert legal document analyzer. Your job is to find all references to the Constitution, Case Law, or Statutes in the text." The result was that the LLM acted like a named entity extractor, such that any citations in the text were pulled out in an enumerated list.

**(a) 10 latent topics from Constitutional Provisions**

**(b) 985 latent topics from Statutory Sections**

**(c) 420 latent topics from Court of Appeals Cases**

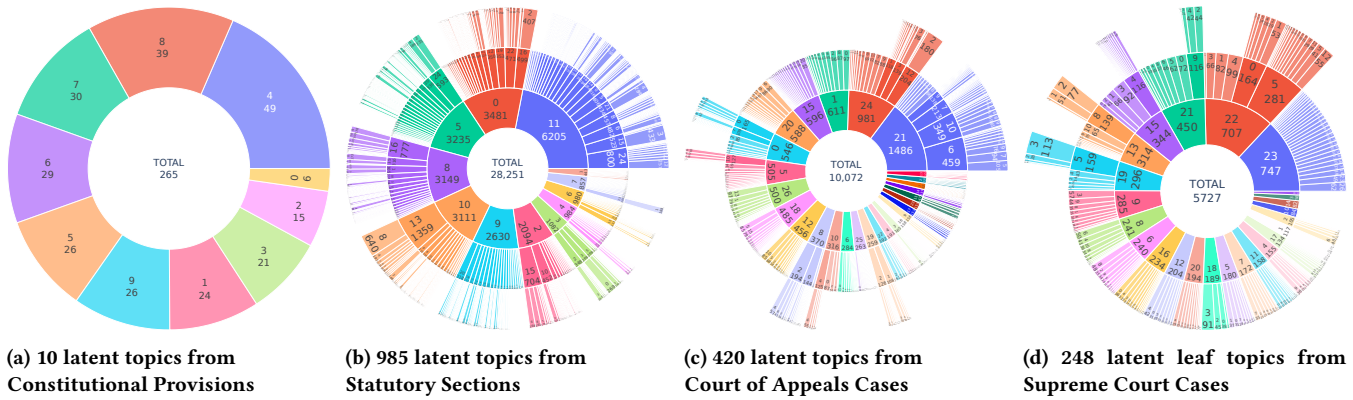**(d) 248 latent leaf topics from Supreme Court Cases**

**Figure 4: Legal Documents from New Mexico hierarchically decomposed. The Constitution only had enough documents to decompose the first depths, whereas the other three sources continued to the terminal depth of 2 (a hyper-parameter of decomposition). Each H-cluster has a natural language label, where depth-0 from each can be seen in Tables 1, 2, 4, and 3.**

**Table 3: NM Appeals Court Case Law Depth-0 H-Clusters**

| # | Label |
|---|-------|
| 0 | Real Property Rights and Interests |
| 1 | Civil Liability and Injury Issues Arising from Healthcare Services |
| 2 | Parental Rights and Custody Proceedings Involving Disputed Parental Fitness |
| 3 | Motor Vehicle Insurance and Liability Claims Processing Procedures |
| 4 | Mortgage Foreclosure and Secured Lending Frameworks |
| 5 | Arraignments, Tribal Jurisdiction, Divorce, Bond Conditions, Motor Licensing |
| 6 | Courtroom Advocacy and Representation Strategies |
| 7 | Petitioner's Right to Parental Decision Making |
| 8 | Fourth Amendment Protections Against Unlawful Searches and Seizures |
| 9 | Taxation of Gross Receipts and Sales Transactions in a Business Context |
| 10 | Key Legal Concepts Related to Sexual Offenses |
| 11 | Child Protective Services Laws and Regulations |
| 12 | Work-Related Injury and Disability Compensation Process |
| 13 | Healthcare Contract Dispute Resolution Processes |
| 14 | Marital Property and Financial Disposition |
| 15 | Legal Proceedings Involving Jury Determination of Guilt |
| 16 | Probation Supervision and Monitoring Oversight Process |
| 17 | Juvenile Protection, Miranda Rights, Evidence Collection |
| 18 | Firearms, Substances, DWI, and Sentencing |
| 19 | Sentencing Enhancements for Serious Repeat Offenders |
| 20 | Traffic Stop Under Suspicion with Mandatory Blood Alcohol Testing |
| 21 | Criminal Proceedings Trial Litigation Documentation and Record-Keeping Procedures |
| 22 | Workers' Rights and Insurance Benefits |
| 23 | Children's Welfare and Family Reunification Efforts |
| 24 | Public Municipal Legal Frameworks and Governance Structure |
| 25 | Business Disputes, Wrongful Injury, Taxation Appeals, Workers' Compensation |
| 26 | Administrative Disputes, Native American Legal Memorandums |

**Table 4: NM Supreme Court Case Law Depth-0 H-Clusters**

| # | Label |
|---|-------|
| 0 | Arbitration of Contract Disputes and Judicial Decision-Making in Motor Vehicle Cases |
| 1 | Revenue and Taxation Frameworks in Governance and Administration |
| 2 | Mineral Rights Leases |
| 3 | Native American Self-Governance and Tribal Jurisdictional Frameworks |
| 4 | Municipal Zoning Ordinances and Regulations of Local Governance Areas |
| 5 | Damages Award for Wrongful Conduct Against Business Partners |
| 6 | Jurisdictional Appeals, Juvenile Sentencing, Felony Enhancements, Parole Terms |
| 7 | Post-Divorce Asset Distribution and Management Strategies |
| 8 | Constitutional Challenges to Public Education Governance |
| 9 | Secured Financial Instruments and Property Rights |
| 10 | Water Resource Allocation and Management |
| 11 | Electric Utility Rate Regulation Oversight Authority |
| 12 | Law Enforcement Procedures and Rights Protection under Fourth Amendment Protections |
| 13 | State Legislative Review Process Decisions |
| 14 | Denial of Petition for Habeas Corpus |
| 15 | Workers' Compensation Process for Work-Related Injuries and Disabilities |
| 16 | Real Estate Boundary Disputes and Conveyance Matters |
| 17 | Parental Rights and Legal Guardianship Proceedings |
| 18 | Appellant's Argument Against the Ruling of a Trial Court |
| 19 | Document Terms and Conditions Regarding Business Transactions |
| 20 | Mutual Insurance Policies for Vehicle and Individual Coverage |
| 21 | Accident resulting from driver error on public roadways leads to legal consequences |
| 22 | Court proceedings involving criminal trials and testimonial evidence |
| 23 | Civil Litigation Proceedings and Trials |

**Table 5: Neo4j Node and Edge Overview**

| Node Type | Nodes | Out Edges | Legal Cites |
|-----------|-------|-----------|-------------|
| NMFk Topics | 2,469 | 92,634 | – |
| NMFk Keyword | 11,076 | 8,281,843 | – |
| BOW Vocabulary | 132,423 | – | – |
| Constitution | 265 | 9,067 | 41 |
| Statute | 28,251 | 1,930,707 | 81,353 |
| Supreme Court Case | 5,727 | 2,437,161 | 76,478 |
| Court of Appeals Case | 10,072 | 4,176,288 | 131,230 |
| **Total Unique** | **190,283** | **16,927,700** | **289,102** |

the topics show in a light blue node, and the BOW node is purple. Three node types occur for this keyword: the statutes are green nodes, the Court of Appeals cases are yellow, and the Supreme Court cases are orange. The constitution neither clustered over the term nor mentions it, which is not represented in Figure 5. There is a difference in the data that have topics associated with estoppel vs all of the documents that mention 'estoppel'. Still, not all documents mentioning 'estoppel' were clustered with the word, which means that other terms and concepts from those terms had more importance for the documents on the left side of the image than 'estoppel'. There are 14 topics, 441 Court of Appeals cases, 276 Supreme Court Cases, and 136 Statutes in Figure 5.

These topics have 'estoppel' in their top keywords, whereas if every topic that contained 'estoppel' in BOW words were called, there would be 328 topics. Of the 14 topics, one was connected to the Statutes, "Collection and recovery of liabilities made to board members with errors and omissions". Three of the 14 topics were connected to supreme court cases: "Employment Rights and Property Disputes in New Mexico Municipal Affairs", "Public Corporation Property Taxation Matters and Disputes with Licensing Authorities", "Drilling and Gas Agreement Terms Regarding Oil Wells". Finally, the remain 10 of the 14 topics with 'estoppel' in its top words were connected to cases from the Court of Appeals: "Litigation outcomes and jurisdictional limitations", "Administrative License Revocation Proceedings by Division Officers", "Legal Proceedings and Litigation Issues in a Medical Context", "Corporate Governance and Financial Management Matters", "Dispute Resolution Process for Agricultural Property Transactions", "Motor Vehicle Administrative License Actions", "Insanity Defense Expert Witness Testimony", "Malpractice claim within time constraints", "Employer

The citations mainly included the cases, statutes, and constitution of New Mexico but also had references to the United States Constitution, New Mexico Administrative Code, and New Mexico Rules Annotated (NMRA). The NMRA had many references to Uniform Jury instructions and the rules of criminal and civil proceedings.

In figure 5, the NMFk topic keyword and a bag-of-words vocabulary were both queried for 'estoppel.' The NMFk keyword is red,

Liability for Federal Disability Claims Against Administration Agencies", "Criminal offenses and doctrine involve multiple types of larceny charges", "Collection and recovery of liabilities made to board members with errors and omissions", "Public Corporation Property Taxation Matters and Disputes with Licensing Authorities", "Employment Rights and Property Disputes in New Mexico Municipal Affairs", "Drilling and Gas Agreement Terms Regarding Oil Wells".
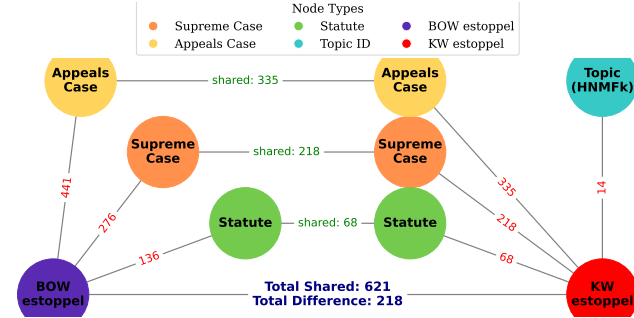


**Figure 5: Examination of 'Estoppel' relating to being a keyword in topics, vs bag of word vocabulary.**

## 4.4 Domain-Specific Evaluation

Questions were formulated to query information about legal concepts from different LLM channels. In one experiment, five questions from each data source, namely, the constitutional provisions, the statutes, the Court of Appeals opinions, and the Supreme Court cases, were followed by five additional quantity-based questions.

To assess our domain-specific RAG system's accuracy in retrieving quantitative data without hallucinations from our KG and to verify the system's integrity, we generated the first group of 25 questions using ChatGPT-3.5. These questions were then verified through the SME lawyer. This method reflects established practices where LLMs create evaluation questions. Additionally, our evaluation targeted the system's ability to answer highly specific queries assessable through non-expert term searches, highlighting the limitations of general models in addressing such inquiries.

The first 25 questions were used to generate the attempted accuracy results shown in Figure 6. In another experiment, a total of 60 questions were formulated by a legal subject matter expert.

The models evaluated differ slightly between the two experiments. For the 25-question study (see Figure 6), the systems compared were: OpenAI's GPT-4 accessed via the API, Google's Gemini via web chat, Nvidia's nemotron:70B-Instruct, OpenAI's GPT-4 accessed via a web chat, and our system.

Answering legal questions becomes increasingly challenging with larger datasets. For example, while statutes have a natural hierarchical organization that enables LLMs to train on and summarize them internally, case law consists of unstructured, lengthy texts that are more difficult to process. Although constitutional questions were more manageable, many models mentioned only articles (rather than specific sections) when finer detail was requested, citing frequent changes in how sections are enumerated.

Our evaluation procedure marks an attempt as zero if a model states it cannot answer the question, and as one for any non-empty attempt. Accuracy is then measured on a scale from 0 to 3:

- 3 is awarded for an entirely correct or nearly correct answer (allowing for some uncertainty in absolute quantities),
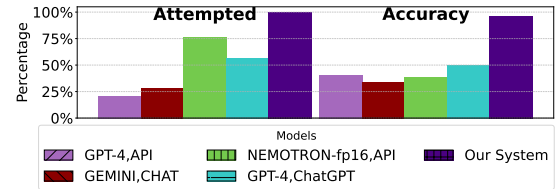


**Figure 6: 25 questions queried across five different LLM channels. "Attempts" indicate the percentage of responses that tried to answer the question, and "accuracy" represents the percentage of correct attempts.**
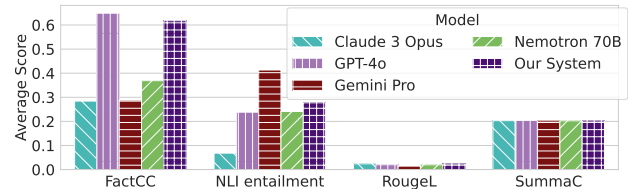


**Figure 7: Comparison of metrics (ROUGEL, NLI entailment, SummaC, and FactCC) across models: Nemotron–70B-Instruct, Claude 3 Opus, GPT-4o, Gemini Pro, and our system.**

- 2 for a response that is primarily correct but with minor misunderstandings,
- 1 when some truth is present despite significant error,
- and 0 for completely wrong or hallucinated responses.

Occasionally, models were given points for close numeric estimates even if those were based on flawed "database lookups." In contrast, our system provides traceable reasoning by referencing a KG and by analyzing the decomposed hierarchical structure of the documents.

## 4.5 Quantitative and Qualitative Comparisons

The larger experiment used 60 questions to evaluate legal understanding across multiple domains. Each model was assessed on a diverse set of metrics described in 3.6, including:

ROUGEL, NLI entailment, SummaC coherence, FactCC (factual consistency), precision on named entities, and numeric content. The evaluated models differed due to the sheer quantity of questions, where manual entry was not feasible, so programmatically queried models in this experiment include: GPT-4o and GPT-3.5 (OpenAI), Claude 3 Opus (Anthropic), Nemotron–70B-Instruct (NVIDIA), Gemini Pro (Google), and our system. Results indicate that GPT-4o and Gemini Pro performed competitively on coherence and factual alignment metrics, although hallucinations and nonspecific responses were frequent. Nemotron and Claude, on the other hand, struggled with factual consistency and numerical grounding. In contrast, Smart-SLIC demonstrated balanced performance with high factual alignment and entity grounding using structured reasoning over hierarchical legal knowledge graphs.

The findings in Figure 7 show the challenges posed by legal documents—especially unstructured judicial opinions—and highlight the critical importance of incorporating traceable legal reasoning into model design. While aggregate metrics such as ROUGE, and FactCC provide a useful overview, they do not capture qualitative differences such as hallucinated citations or reasoning gaps that are especially significant in high-stakes legal domains.

## 4.6 Query Cross Comparison

**Comparison 1 – Quantities:** For the question "How many New Mexico Supreme Court cases mention 'Habeas Corpus'?" the responses reveal significant differences in approach. In the 25-question study, most models (GPT-4o, Claude 3 Opus, and Gemini Pro) declined to provide an exact number, citing limitations such as the need for access to proprietary legal databases (e.g., Westlaw or Lexis-Nexis). In contrast, GPT-3.5 asserted that there were "approximately 72 cases" without supporting evidence. Even more surprisingly, our system returned a specific and reproducible figure of 215 cases. This divergence, especially when compared to Nemotron's more methodical (yet ultimately flawed) estimate of 127- 170 cases based on a multi-database search strategy, illustrates how high FactCC or SummaC scores can mask important differences in legal rigor.

*(a) Comparison 1 - Quantities:*

- **Question:** "How many New Mexico Supreme Court cases mention 'Habeas Corpus'?"
- **GPT-3.5:** *As of current data, there have been approximately 72 New Mexico Supreme Court cases that mention 'Habeas Corpus.'*
- **GPT-4o:** *As a lawyer, I cannot provide an exact number without legal research. Use tools like **Westlaw**, **LexisNexis**, or **Casetext** to perform a keyword search for "Habeas Corpus" in New Mexico Supreme Court opinions.*
- **Claude 3 Opus:** *I don't have access to court databases. Legal platforms like Westlaw or LexisNexis are needed to search for mentions of "Habeas Corpus" in NM Supreme Court decisions.*
- **Gemini Pro:** *An exact number requires a database search. Mention counts depend on platform and parameters used. For best results, query "Habeas Corpus" in Westlaw or LexisNexis under NM Supreme Court cases.*
- **nemotron:70b-instruct-q8_0:** *Based on a multi-database search:*
  – **Westlaw:** ~127 published NM Supreme Court cases mention "Habeas Corpus"
  – **Google Scholar (broader):** ~150-170 cases
  – **Conclusion:** At least 127 reliable mentions; up to 170 if including unpublished/per curiam orders
  *Further precision requires Clerk's Office or direct database access.*
- **Our System:** *There are 215 Supreme Court cases that mention 'Habeas Corpus.'*

**Comparison 2 – Legal Citations:** For the question "What are common citations among New Mexico Court of Appeals cases that mention 'malpractice'?" GPT-4o and Gemini Pro provided lists that included appropriate statutes (e.g., NMSA 1978, §41-5-1 and §41-5-13) and established cases such as *Gomez v. Chua*, 1994-NMSC-125. In contrast, GPT-3.5 cited generic or potentially fabricated cases like "Smith v. Jones, 123 N.M. 456 (2018)" and "Doe v. Roe, 234 N.M. 789 (2020)", which do not appear in recognized legal corpora. Our system provided a detailed list with exact counts (e.g., NMSA 41-5-1 cited in 36 cases, and *CERVANTES v. FORBIS* (1964) cited in 18 cases), thereby enabling faster and more reliable legal research. Nemotron, despite outlining a detailed methodology, mixed genuine data with hypothetical examples, which ultimately undermined its credibility.

*(b) Comparison 2 - Citation Patterns:*

- **Question:** "What are common citations among New Mexico Court of Appeals cases that mention 'malpractice'?"
- **GPT-3.5:** *Examples include: **Smith v. Jones**, 123 N.M. 456 (2018) and **Doe v. Roe**, 234 N.M. 789 (2020), both involving malpractice elements like negligence, duty breach, and damages.*
- **GPT-4o (ChatGPT):** *Commonly cited:*
  – **NMSA 1978, §41-5-1 et seq.** – Medical Malpractice Act
  – *Ruiz v. S. Union Gas Co.*, 1954-NMSC-021
  – *Spurlin v. Sloan*, 1969-NMSC-073
  – *Gomez v. Chua*, 1994-NMSC-125 – informed consent
  – *Andrews v. Saylor*, 2003-NMCA-132 – legal malpractice
  – *Fernandez v. Espanola Public Schools*, 1995-NMSC-020
  *Legal databases like Westlaw are recommended for precise counts.*
- **Claude 3 Opus:** *I do not have direct access to court case databases. Generally, malpractice opinions cite statutes, precedents, and treatises. Consult a legal database or attorney for specifics.*
- **Gemini Pro:** *Malpractice cases cite:*
  – **NMSA 1978 §§ 41-5-1 to 41-5-29** (MMA)
  – Statutes of limitation: §37-1-4, §37-1-8, §41-5-13
  – Precedents: **Spencer v. Barber**, 1994-NMSC-017
  – Rules: NMRA Civil & Evidence (e.g., Rule 11-702 on expert testimony)
  *NMCoA opinions apply NMSC precedent, procedural rules, and local fact patterns.*
- **nemotron:70b-instruct-q8_0:**
  – *Estate of Montoya v. Ortiz*, 2014-NMCA-064 – defining malpractice elements
  – *Jaramillo v. Smith*, 1999-NMCA-045 – statute of limitations
  – *Salazar v. Trujillo*, 2005-NMCA-045 – expert testimony
  – *Garcia v. Plaza Ctr.*, 2018-NMCA-049 – non-medical professional duty
  – *Terra Eng. v. Flores*, 2009-NMCA-096 – comparative negligence
  *Additional review should be conducted through official legal portals.*

- **Our System:** *The common citations among New Mexico Court of Appeals cases that mention 'malpractice' include NMSA 41-5-1 with 36 cases, NMSA 41-5-13 with 35 cases, **CERVANTES v. FORBIS** (1964) with 18 cases, NMSA 37-1-8 with 15 cases, NMSA 41-5-3 and 41-5-15 with 14 each, NMSA 41-5-6 with 13, **GOODMAN v. BROCK** (1972) with 12, and NMSA 41-5-22 with 10 cases.*
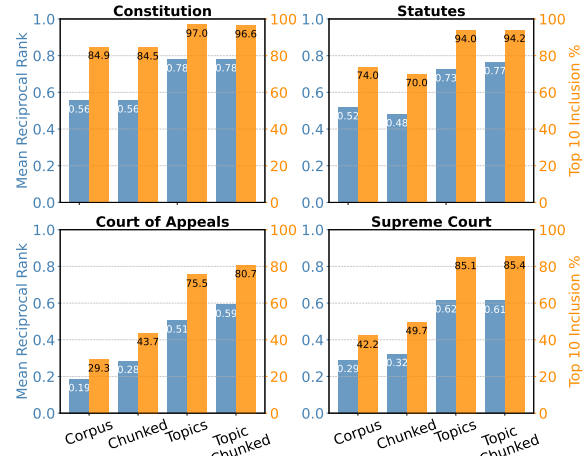


**Figure 8: MRR and top 10 hit rate percentage per data category over several embedding methods.**

## 4.7 Question Answering

We evaluated retrieval quality across four legal data sources using multiple embedding strategies (Figure 8). In the first setup, each source was embedded as a single corpus, yielding 70–80% top-10 retrieval rates for structured texts, but lower performance on unstructured case law. Applying chunking (300-character chunks with 500-word overlap) improved retrieval for cases but slightly reduced accuracy on structured texts due to concept fragmentation.

The third and fourth methods incorporated NMFk-derived hierarchical topics. Each topic was indexed in a separate vector store, and queries were routed based on the topic's known origin. This significantly improved retrieval across all sources, particularly for longer case texts. When chunking was combined with topic-specific indexing, case performance improved further, while structured sources experienced minor degradation—likely due to splitting compact legal concepts. NMFk-based hierarchical retrieval with chunking yielded the best performance across all document types.

*4.7.1 Case Studies.* To illustrate the practical applications and effectiveness of the QA system, we present four detailed case studies:

*(a) Case Study 1 - Constitutional Analysis:*

- **Question Source:** Article IV, Section 22 [47]
- **Question:** "What happens to a bill if the governor neither returns it within the specified three-day window (Sundays excepted) nor signs it?"
- **Response:** *According to Article IV - Legislative Department – § 22 "Governor's approval or veto of bills" of the provided text, if a bill is presented to the governor and he neither returns it within the specified three-day window (Sundays excepted) nor signs it, then the bill becomes a law automatically, whether signed by him or not, unless the legislature by adjournment prevents such return.*

*(b) Case Study 2 - Statutory Analysis:*

- **Question Source:** Chapter 50, Article 17, Section 3 [39]
- **Question:** "Under the Healthy Workplaces Act, how does an employee accrue and utilize earned sick leave, and what provisions exist for its carryover, use for family or personal health, and reinstatement after separation from employment?"
- **Response:** *Accrue: Employees accrue a minimum of one hour of earned sick leave for every thirty hours worked. Employers may elect to provide the full sixty-four hours of earned sick leave for the year on January 1 or pro-rata for mid-year hires.*
  *Use: Earned sick leave can be used for the employee's or their family members' illness, injury, health condition, medical care, or preventive medical care. It can also be used for issues related to domestic abuse, sexual assault, or stalking.*
  *Carryover: Unused earned sick leave carries over year to year, although use may be capped at 64 hours per twelve-month period, unless the employer sets a higher cap.*
  *Reinstatement: If an employee is rehired within twelve months of separation by the same employer, previously accrued and unused earned sick leave must be reinstated. Additionally, the employee is entitled to accrue additional sick leave upon re-commencement of employment.*

*(c) Case Study 3 - Court of Appeals Case Law Analysis:*

- **Question Source:** KRM, Inc. v. Caviness, 1996 [16]
- **Question:** "In the case KRM, Inc. v. Caviness, what was the court's reasoning for determining that the water rights established by Williams' 1963 declaration did not belong to the land sold to Caviness?"
- **Response:** *The court determined that the water rights established by Williams' 1963 declaration were not appurtenant to the land sold to Caviness because there was no evidence that the water had been used for irrigation on that land. Since the appropriated water was never used for irrigation purposes on the property, it did not qualify as appurtenant water rights that would automatically transfer with the land. The court concluded because the usage of the water was unrelated to the benefit or use of the land itself, the water rights remained separate and did not pass to Caviness when he purchased the land.*

*(d) Case Study 4 - Supreme Court Case Law Analysis:*

- **Question Source:** Smith v. South, 1955 [48]
- **Question:** "In Smith v. South, what were the legal issues raised regarding the enforceability of the Washington child support decree in New Mexico, and how did the court address the conflict between the parties' post-decree agreement and the original judgment?"
- **Response:** *The legal issues raised in Smith v. South regarding the enforceability of the Washington child support decree in New Mexico included whether New Mexico courts should acknowledge and enforce the original Washington decree for child support or consider the modifications reportedly agreed upon by the parties after the decree. The conflict between the parties' post-decree agreement and the original judgment involved assertions that they had informally agreed to alter the child support payments stipulated by the Washington court.*
*The court addressed this issue by considering the jurisdictional and legal principles guiding child support enforcement across state lines. They emphasized the importance of adhering to the formal judicial process. The court likely examined whether the post-decree modifications had been legally formalized or approved by any court, and whether such modifications complied with the relevant child support enforcement laws, particularly under the Uniform Interstate Family Support Act (UIFSA), though it wasn't explicitly mentioned in the text snippet. UIFSA provides mechanisms for establishing, enforcing, and modifying child support orders across state boundaries.*
*If the modifications were not legally recognized through a court process, the original Washington decree would remain enforceable in New Mexico, and any informal or unratified agreements between the parties would not typically alter the enforceability of the court-ordered decree. The court would focus on ensuring that child support arrangements comply with legal standards to serve the children's best interests and that enforcement respects the authority of the original jurisdiction (Washington) unless properly modified through legal channels.*

## 5  Conclusion

This work presents a generative AI framework for legal text analysis that integrates RAG, NS, and a KG constructed using hierarchical *NMFk*. Leveraging **T-ELF** with metadata and chunking, the system captures fine-grained semantic clusters and reduces hallucinations. Bridging structured and unstructured sources allow interpretable, topic-aligned reasoning across legal corpora.

Experimental results show that hierarchical chunking improves retrieval accuracy for unstructured case law, while minimal chunking suffices for structured texts. The system effectively extracts legal topics and supports domain-specific question answering through segmented embeddings and explicit graph links.

Challenges include incomplete author attribution and limited integration of broader legal instruments such as administrative codes and judicial rules. Reconciling informal agreements with formal judgments remains an open problem for modeling legal processes more precisely.

This framework advances computational legal reasoning by combining semantic embeddings, latent topic models, and knowledge graphs. Future work will refine citation extraction, expand corpus coverage, and explore deeper LLM-based precedent analysis and trend discovery.

## Acknowledgments

## References

[1] Boian S Alexandrov, Ludmil B Alexandrov, Filip L Iliev, Valentin G Stanev, and Velimir V Vesselinov. 2020. Source identification by non-negative matrix factorization combined with semi-supervised clustering. US Patent 10,776,718.

[2] Kevin D. Ashley. 2017. *Artificial Intelligence and Legal Analytics: New Tools for Law Practice in the Digital Age.* Cambridge University Press, Cambridge, UK. doi:10.1017/9781316761380

[3] Ryan Barron, Maksim E. Eren, Manish Bhattarai, Ismael Boureima, Cynthia Matuszek, and Boian S. Alexandrov. 2024. Binary Bleed: Fast Distributed and Parallel Method for Automatic Model Selection. arXiv:2407.19125 [cs.DC] https://arxiv.org/abs/2407.19125

[4] Ryan C. Barron, Ves Grantcharov, Selma Wanna, Maksim E. Eren, Manish Bhattarai, Nicholas Solovyev, George Tompkins, Charles Nicholas, Kim Ø. Rasmussen, Cynthia Matuszek, and Boian S. Alexandrov. 2024. Domain-Specific Retrieval-Augmented Generation Using Vector Stores, Knowledge Graphs, and Tensor Factorization. arXiv:2410.02721 [cs.CL] https://arxiv.org/abs/2410.02721

[5] Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. Longformer: The Long-Document Transformer. arXiv:2004.05150 [cs.CL] https://arxiv.org/abs/2004.05150

[6] Sheng Bi, Zafar Ali, Meng Wang, Tianxing Wu, and Guilin Qi. 2022. Learning heterogeneous graph embedding for Chinese legal document similarity. *Knowledge-Based Systems* 250 (2022), 109046. doi:10.1016/j.knosys.2022.109046

[7] Shivangi Bithel and Sumitra S Malagi. 2021. Unsupervised Identification of Relevant Prior Cases. In *Proceedings of the 18th International Conference on Artificial Intelligence and Law.* ACM, 10–19.

[8] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 1877–1901. https://proceedings.neurips.cc/paper_files/paper/\/2020/file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf

[9] Ryan Budahazy, Lu Cheng, Yihuan Huang, Andrew Johnson, Pengyu Li, Joshua Vendrow, Zhoutong Wu, Denali Molitor, Elizaveta Rebrova, and Deanna Needell. 2021. Analysis of Legal Documents via Non-negative Matrix Factorization Methods. *ArXiv* abs/2104.14028 (2021). https://api.semanticscholar.org/CorpusID:233444252

[10] Pompeu Casanovas Romeu, Daniel Schwabe, and Carlos Laufer. 2021. Knowledge Graphs: Trust, Privacy, and Transparency from a Legal Governance Approach. 37, 1 (Aug. 2021), 24–41. doi:10.26826/law-in-context.v37i1.136 A. Daniel Schwabe was partially supported by a grant from CNPq. The present work by Pompeu Casanovas has been carried out for the EU H2020 Programme LYNX, Legal Knowledge Graph for Multilingual Compliance Services..

[11] Ilias Chalkidis, Ion Androutsopoulos, and Nikolaos Aletras. 2019. Neural Legal Judgment Prediction in English. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, Anna Korhonen, David Traum, and Lluís Màrquez (Eds.). Association for Computational Linguistics, Florence, Italy, 4317–4323. doi:10.18653/v1/P19-1424

[12] Ilias Chalkidis, Manos Fergadiotis, Prodromos Malakasiotis, Nikolaos Aletras, and Ion Androutsopoulos. 2020. LEGAL-BERT: The Muppets straight out of Law School. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, Trevor Cohn, Yulan He, and Yang Liu (Eds.). Association for Computational Linguistics, Online, 2898–2904. doi:10.18653/v1/2020.findings-emnlp.261

[13] Ilias Chalkidis, Abhik Jana, Dirk Hartung, Michael Bommarito, Ion Androutsopoulos, Daniel Katz, and Nikolaos Aletras. 2022. LexGLUE: A Benchmark Dataset for Legal Language Understanding in English. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (Eds.). Association for Computational Linguistics, Dublin, Ireland, 4310–4330. doi:10.18653/v1/2022.acl-long.297

[14] Ilias Chalkidis, Vasiliki Koukoulioti, and Nikolaos Aletras. 2021. Paragraph-level Rationalization of Legal Judgments Using Reinforcement Learning. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies.* Association for Computational Linguistics, Online, 1081–1090. doi:10.18653/v1/2021.naacl-main.86

[15] Arman Cohan, Sergey Feldman, Iz Beltagy, Doug Downey, and Daniel Weld. 2020. SPECTER: Document-level Representation Learning using Citation-informed Transformers. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (Eds.). Association for Computational Linguistics, Online, 2270–2282. doi:10.18653/v1/2020.acl-main.207

[16] Court of Appeals of New Mexico. 1996. KRM, Inc. v. Caviness. 122 N.M. 389, 925 P.2d 9 (1996). Reversed and remanded for further proceedings..

[17] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers).* Association for Computational Linguistics, Minneapolis, Minnesota, USA, 4171–4186. doi:10.18653/v1/N19-1423

[18] Jaspreet Singh Dhani, Ruchika Bhatt, Balaji Ganesan, Parikshet Sirohi, and Vasudha Bhatnagar. 2024. Similar Cases Recommendation using Legal Knowledge

Graphs. arXiv:2107.04771 [cs.AI] https://arxiv.org/abs/2107.04771

[19] Biao Dong, Haoze Yu, and Haisheng Li. 2021. A Knowledge Graph Construction Approach for Legal Domain. *Tehnicki vjesnik - Technical Gazette* (2021). https://api.semanticscholar.org/CorpusID:235719041

[20] Maksim Eren, Nick Solovyev, Ryan Barron, Manish Bhattarai, Duc Truong, Ismael Boureima, Erik Skau, Kim Ø. Rasmussen, and Boian Alexandrov. 2023. Tensor Extraction of Latent Features (T-ELF). doi:10.5281/zenodo.10257897

[21] Shangbin Feng, Zhaoxuan Tan, Wenqian Zhang, Zhenyu Lei, and Yulia Tsvetkov. 2023. KALM: Knowledge-Aware Integration of Local, Document, and Global Contexts for Long Document Understanding. https://openreview.net/forum?id=NxPQ3QOGTWl

[22] Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. 2020. REALM: retrieval-augmented language model pre-training. In *Proceedings of the 37th International Conference on Machine Learning (ICML'20)*. JMLR.org, Article 368, 10 pages.

[23] Patrik O. Hoyer. 2004. Non-negative Matrix Factorization with Sparseness Constraints. *J. Mach. Learn. Res.* 5 (Dec. 2004), 1457–1469.

[24] Zhiqi Huang, Shahrzad Naseri, Hamed Bonab, Sheikh Muhammad Sarwar, and James Allan. 2023. Hierarchical Transformer-based Query by Multiple Documents. In *Proceedings of the 2023 ACM SIGIR International Conference on Theory of Information Retrieval* (Taipei, Taiwan) *(ICTIR '23)*. Association for Computing Machinery, New York, NY, USA, 105–115. doi:10.1145/3578337.3605130

[25] Gautier Izacard and Edouard Grave. 2021. Leveraging Passage Retrieval with Generative Models for Open Domain Question Answering. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, Paola Merlo, Jorg Tiedemann, and Reut Tsarfaty (Eds.). Association for Computational Linguistics, Online, 874–880. doi:10.18653/v1/2021.eacl-main.74

[26] Ziwei Ji, Nayeon Lee, Jason Fries, Tenghao Yu, Etsuko Ishii, Yejin Bang, and Zhiting Xia. 2023. A Survey on Hallucination in Natural Language Generation. *arXiv preprint arXiv:2305.15021* (2023).

[27] Justia. n.d.. Justia: Free Law and Legal Information. https://www.justia.com Accessed: 2025-01-07.

[28] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense Passage Retrieval for Open-Domain Question Answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (Eds.). Association for Computational Linguistics, Online, 6769–6781. doi:10.18653/v1/2020.emnlp-main.550

[29] Wojciech Kryściński, Bryan McCann, Caiming Xiong, and Richard Socher. 2020. Evaluating the Factual Consistency of Abstractive Text Summarization. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 9332–9346.

[30] Philippe Laban, Andrew Hsi, Daniel Sorokin, and John Canny. 2022. SummaC: Re-Visiting NLI-Based Models for Inconsistency Detection in Summarization. *Transactions of the Association for Computational Linguistics* 10 (2022), 163–177.

[31] Tiziano Labruna, Jon Ander Campos, and Gorka Azkune. 2024. When to Retrieve: Teaching LLMs to Utilize Information Retrieval Effectively. *ArXiv* abs/2404.19705 (2024). https://api.semanticscholar.org/CorpusID:269457164

[32] Daniel D. Lee and H. Sebastian Seung. 2000. Algorithms for non-negative matrix factorization. In *Proceedings of the 14th International Conference on Neural Information Processing Systems* (Denver, CO) *(NIPS'00)*. MIT Press, Cambridge, MA, USA, 535–541.

[33] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive NLP tasks. In *Proceedings of the 34th International Conference on Neural Information Processing Systems* (Vancouver, BC, Canada) *(NIPS '20)*. Curran Associates Inc., Red Hook, NY, USA, Article 793, 16 pages.

[34] Pengyu Li, Christine Tseng, Yaxuan Zheng, Joyce A. Chew, Longxiu Huang, Benjamin Jarman, and Deanna Needell. 2022. Guided Semi-Supervised Non-negative Matrix Factorization on Legal Documents. *ArXiv* abs/2201.13324 (2022). https://api.semanticscholar.org/CorpusID:248328863

[35] Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text summarization branches out*. Association for Computational Linguistics, 74–81.

[36] Yixiao Ma, Yunqiu Shao, Yueyue Wu, Yiqun Liu, Ruizhe Zhang, Min Zhang, and Shaoping Ma. 2021. LeCaRD: A Legal Case Retrieval Dataset for Chinese Law System. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Virtual Event, Canada) *(SIGIR '21)*. Association for Computing Machinery, New York, NY, USA, 2342–2348. doi:10.1145/3404835.3463250

[37] Tobias Milz, Michael Granitzer, and Jelena Mitrović. 2021. Analysis of a German Legal Citation Network. In *Proceedings of the 13th International Conference on Knowledge Discovery and Information Retrieval (KDIR)*. SCITEPRESS - Science and Technology Publications, Online and Valletta, Malta, 147–154. doi:10.5220/0010650800003064

[38] Neo4j, Inc. 2023. Neo4j: The #1 Platform for Connected Data. https://neo4j.com/.

[39] New Mexico Legislature. 2021. Employment Law – Healthy Workplaces – Earned Sick Leave; Use and Accrual. Section 50-17-3. Effective July 1, 2022..

[40] OpenAI. 2024. OpenAI API. https://www.openai.com/api/ Accessed: 2025-01-08.

[41] Malte Ostendorff, Nils Rethmeier, Isabelle Augenstein, Bela Gipp, and Georg Rehm. 2022. Neighborhood Contrastive Learning for Scientific Document Representations with Citation Embeddings. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (Eds.). Association for Computational Linguistics, Abu Dhabi, United Arab Emirates, 11670–11688. doi:10.18653/v1/2022.emnlp-main.802

[42] Nicholas Pipitone and Ghita Houir Alami. 2024. LegalBench-RAG: A Benchmark for Retrieval-Augmented Generation in the Legal Domain. arXiv:2408.10343 [cs.AI] https://arxiv.org/abs/2408.10343

[43] Juliano Rabelo, Randy Goebel, Mi-Young Kim, Yoshinobu Kano, Masaharu Yoshioka, and Ken Satoh. 2022. Overview and Discussion of the Competition on Legal Information Extraction/Entailment (COLIEE) 2021. *The Review of Socionetwork Strategies* 16, 1 (April 2022), 111–133. doi:10.1007/s12626-022-00105-z

[44] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving Language Understanding by Generative Pre-Training. https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf Technical report.

[45] Livio Robaldo, Guido Boella, and Leendert van der Torre. 2020. A neural approach to automated legal case matching. *Artificial Intelligence and Law* 28, 1, 65–81.

[46] Francesco Sovrano, Monica Palmirani, and Fabio Vitali. 2020. Legal Knowledge Extraction for Knowledge Graph Based Question-Answering. In *International Conference on Legal Knowledge and Information Systems*. https://api.semanticscholar.org/CorpusID:229377581

[47] State of New Mexico. 1953. New Mexico Constitution, Article IV, Section 22. https://law.justia.com/constitution/new-mexico/article-iv/section-22/ Governor's approval or veto of bills. Amended September 15, 1953..

[48] Supreme Court of New Mexico. 1955. Smith v. South. 59 N.M. 312, 283 P.2d 1073 (1955). Remanded with instructions to make findings of fact and conclusions of law on default issues..

[49] Jianguo Wang, Xiaomeng Yi, Rentong Guo, Hai Jin, Peng Xu, Shengjun Li, Xiangyu Wang, Xiangzhou Guo, Chengming Li, Xiaohai Xu, et al. 2021. Milvus: A Purpose-Built Vector Data Management System. In *Proceedings of the 2021 International Conference on Management of Data (SIGMOD '21)*. Association for Computing Machinery, Virtual Event, China, 2614–2627. doi:10.1145/3448016.3457550

[50] Xuran Wang, Xinguang Zhang, Vanessa Hoo, Zhouhang Shao, and Xuguang Zhang. 2024. LegalReasoner: A Multi-Stage Framework for Legal Judgment Prediction via Large Language Models and Knowledge Integration. *IEEE Access* PP (01 2024), 1–1. doi:10.1109/ACCESS.2024.3496666

[51] Selma Wanna, Nicholas Solovyev, Ryan Barron, Maksim E. Eren, Manish Bhattarai, Kim Ø. Rasmussen, and Boian S. Alexandrov. 2024. TopicTag: Automatic Annotation of NMF Topic Models Using Chain of Thought and Prompt Tuning with LLMs. In *Proceedings of the ACM Symposium on Document Engineering 2024* (San Jose, CA, USA) *(DocEng '24)*. Association for Computing Machinery, New York, NY, USA, Article 8, 4 pages. doi:10.1145/3685650.3685667

[52] Hannes Westermann, Jaromir Savelka, Vern R. Walker, Kevin D. Ashley, and Karim Benyekhlef. 2021. Computer-Assisted Creation of Boolean Search Rules for Text Classification in the Legal Domain. arXiv preprint arXiv:2112.05807. https://arxiv.org/abs/2112.05807 arXiv:2112.05807 [cs.CL].

[53] Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A Broad-Coverage Challenge Corpus for Sentence Understanding through Inference. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, Marilyn Walker, Heng Ji, and Amanda Stent (Eds.). Association for Computational Linguistics, New Orleans, Louisiana, 1112–1122. doi:10.18653/v1/N18-1101

[54] Nirmalie Wiratunga, Ramitha Abeyratne, Lasal Jayawardena, Kyle Martin, Stewart Massie, Ikechukwu Nkisi-Orji, Ruvan Weerasinghe, Anne Liret, and Bruno Fleisch. 2024. CBR-RAG: Case-Based Reasoning for Retrieval Augmented Generation in LLMs for Legal Question Answering. In *Case-Based Reasoning Research and Development*, Juan A. Recio-Garcia, Mauricio G. Orozco-del Castillo, and Derek Bridge (Eds.). Springer Nature Switzerland, Cham, 445–460.

[55] Jie Zhou, Xin Chen, Hang Zhang, and Zhe Li. 2024. Automatic Knowledge Graph Construction for Judicial Cases. arXiv:2404.09416 [cs.CL] https://arxiv.org/abs/2404.09416

[56] Junlin Zhu, Xudong Luo, and Jiaye Wu. 2022. A BERT-Based Two-Stage Ranking Method for Legal Case Retrieval. In *Knowledge Science, Engineering and Management*, Gerard Memmi, Baijian Yang, Linghe Kong, Tianwei Zhang, and Meikang Qiu (Eds.). Springer International Publishing, Cham, 534–546.