

## NLP FOR LEGAL TEXT MINING: AUTOMATING COMPLIANCE AND POLICY ANALYSIS

**Dr. Jaspreet kaur**

Professor, Computer science and Engineering, Chandigarh group of colleges,  
Jhanjeri, Mohali, Punjab rajjaspreet@gmail.com

### **Abstract**

The rapid expansion of legal and regulatory texts has increased the burden on legal and compliance professionals, making manual review increasingly ineffective and inconsistent. Advances in natural language processing (NLP), machine learning (ML), and deep learning (DL) have made it possible to automate legal text mining, court prediction, and regulatory policy analysis. Guess , Guess what? A synthesis of recent literature shows that traditional machine learning methods such as Support Vector Machines, K-Nearest Neighbor's, and Random Forests are still widely used, while transformer-based deep learning models such as BERT, Legal-BERT, and Fin BERT have significantly improved performance on tasks including classification, summarization, liability extraction, and question answering. A systematic design of research between 2015 and 2022 highlights increasing methodological diversity, increasing availability of domain-specific models, and expanding compliance automation applications.

Empirical studies show that domain-aligned transformers achieve high accuracy, often exceeding 90%, in interpreting regulatory requirements, extracting obligations, and mapping cross-border rules across large legal bodies. Pilot deployments at financial institutions indicate significant reductions in manual workloads and increased detection of compliance risks, supported by integration with governance, risk and compliance (GRC) systems and explainable AI technologies.

**Keywords:** Legal NLP; compliance automation; Regulatory text , text mining; machine learning; deep learning; transformers; legal Bert; Burt the Finn; Court prediction; Withdrawal of commitment; political analysis; Retching. Explainable artificial intelligence.

### **1. Introduction**

The legal field traditionally relies on human expertise to interpret laws, regulations, case law, contracts and policy documents. These texts is often written in dense, technical language and governed by complex, complex principles of interpretation, making legal analysis time-consuming, expensive, and prone to human error. However, in recent years, the rapid expansion of legal and regulatory content driven by globalization, digital governance, and the evolution of compliance standards has created an urgent , urgent need for scalable, automated approaches to managing and understanding legal information.

Along with this growth, artificial intelligence (AI) and natural language processing (NLP) have , have advanced significantly in a bunch of sectors such as healthcare, finance, cyber security and public administration. These developments have been stimulated by the increasing availability of large data sets and the development of sophisticated algorithms capable of interpreting human language. The legal field has traditionally relied on human expertise to interpret laws, regulations, case law, contracts, and policy documents. These texts are often written in dense, highly technical

language and governed by complex interpretative principles, making legal analysis time-consuming, costly, and susceptible to human error. In recent years, however, the rapid expansion of legal and regulatory content—driven by globalization, digital governance, and evolving compliance standards—has created an urgent need for scalable, automated solutions capable of managing, analyzing, and extracting insights from vast volumes of legal information.



**Figure1.Natural Language Processing (NLP)**

Parallel to this growth, artificial intelligence (AI) and natural language processing (NLP) have progressed significantly across sectors such as healthcare, finance, cybersecurity, and public administration. These advancements have been fueled by the increasing availability of large datasets and the development of sophisticated algorithms capable of interpreting and generating human language with high accuracy. As a result, AI-driven tools are now well-positioned to address longstanding challenges in the legal domain by enhancing efficiency, reducing manual workload, and enabling more consistent and data-driven decision-making. Fig 1 contains Natural Language Processing (NLP) is a branch of artificial intelligence that enables computers to understand, interpret, and generate human language.

From virtual assistants like Siri and Alexa to sophisticated customer service chatbots and advanced language translation services, NLP is transforming how we interact with technology. Its ability to process and analyze vast amounts of natural language data makes it a crucial tool in various industries, enhancing everything from business operations to healthcare diagnostics.

At the Same time, the need for automation is greater , greater than ever. Seriously, Modern organizations face a huge amount of regulatory obligations from GDPR and HIPAA to financial regulations such as Basel , Basel III that require constant monitoring and rapid interpretation of changing legislation.

## **2. Literatur**

The rapid integration of artificial intelligence (AI) into law has sparked widespread academic and industrial interest, particularly at the intersection of law, computer science, and linguistics. Early

legal AI systems were mostly rule-based and relied on hand-coded expert knowledge to mimic legal reasoning in areas such as tax law and welfare entitlement. Seriously, Although these systems provided fundamental insights, their reliance on strict logic made them ill-suited to capture the ambiguity, contextual nuances, and evolving nature of legal interpretation.

With the advent of machine , machine learning (ML) and natural language processing (NLP) the possibilities for legal automation have , have expanded dramatically. Tools like ROSS Intelligence and Case Text have paved the way for AI-driven legal research enabling semantic search across large case law repositories. In parallel contract analytics platforms such as Kira Systems and Lager have demonstrated that its feasible to automatically extract terms obligations and risks from complex , complex legal documents. These developments highlight the transformative potential of artificial intelligence in legal research due diligence compliance monitoring and predictive court decision.

A bunch of literature reviews have examined the development of legal NLP from different perspectives. Calcites and Camas highlighted early advances in deep learning architectures for tasks such as text , text classification, information retrieval, and information extraction, noting that they are moving from hard feature engineering to neural approaches. Monthlong and Becker provide a bibliometric analysis that highlights growth trends in legal AI research, albeit , albeit with limited insight into technical challenges. Other studies focused on domain-specific tasks, such as Sheikh, Sheikh and Nirmala's review of summation techniques or Rubellite systematic literature review (SLR) on predicting court decisions, that identified support, support vector machines (SVMs) and transformer-based models as dominant approaches. Larger , Larger surveys such as those , those conducted by Sansone and Peril have explored legal information retrieval systems, identifying ongoing challenges in understanding the structure of legal documents and improving the accuracy of search, summarization and recommendations. Like, Similarly, Soave et al. He studied the transformation of natural language contracts into formal specifications, highlighting unresolved issues such as domain dependency, the accuracy of semantic extraction, and the difficulty of creating legally valid formal representations.

A bunch of cross-cutting challenges emerge from these reviews. Like, Legal texts are linguistically dense, domain-specific, and often lack large annotated datasets due to confidentiality. Data scarcity complicates the process of developing robust , robust supervised learning models. The need for explanation is another recurring theme. In high-stakes areas such as law, decisions must be transparent, explainable and verifiable. This has led to an increased focus on explainable artificial intelligence (XAI) and accountability frameworks supported by initiatives such as the European Union's Legal Embedded Ethics (LEE) project. Knowledge representation plays an equally important role. Efforts such as LKIF-Core and Leafroll aim to formalize legal , legal concepts and relationships to support legal inference, cross-referencing, and interpretation. These ontological systems can improve AI's ability to navigate jurisdictional differences, a key requirement as legal principles and regulatory obligations vary widely across countries and sectors. For example, the interpretation of compliance may differ significantly between the European Union's General Data Protection Regulation (GDPR) and the California Consumer Privacy Act (CCPA) in the United States.

In addition to traditional legal tasks, NLP is increasingly being used in financial regulation, fraud detection, market , market analysis and sentiment assessment. Domain-specific language models

such as Fin BERT and Legal-BERT have significantly improved performance by capturing the unique linguistic nuances of financial disclosures or legal language.

Complementary tools such as Laxly, Laxly and Spacey Legal have advanced entity extraction, term tagging and compliance-supporting workflows. Emerging research also explores multimodal and multidisciplinary approaches, integrating artificial intelligence with blockchain technology, big data analytics, and secure infrastructures to improve, improve the transparency and auditability of regulatory systems... Like, However, scholars continue to warn of significant ethical, ethical risks related to bias, data imbalance, and ambiguous stereotyping of behavior—issues that can lead to unfair or discriminatory outcomes when applied to sensitive areas such as criminal justice or insurance contracting. You know, know what? While previous literature reviews have provided valuable insights, they often focus, focus on narrow tasks, specific data sets, or isolated methodological perspectives. Some papers provide a comprehensive overview of legal NLP covering all key tasks, resources, trends and constraints. Furthermore, there are only a limited number of systematic literature reviews (SLRs) and systematic mapping studies (SMSs), and none comprehensively cover the period 2015–2022 – a time, time when legal NLP has undergone a major technological transformation driven by deep, deep learning and large language models.

The current study addresses this gap by conducting a systematic mapping study across the full range of legal NLP task. And oh yeah, It combines achievements in forensic prediction, contract analysis, information retrieval, compliance automation, text representation methods, domain-specific NLP models, and cross-jurisdictional challenges.

Like, By bringing together progress in these subfields, the review highlights current limitations, identifies emerging research directions, and lays the groundwork for the development of more interpretable, adaptable, and ethically sound NLP frameworks to support modern legal and regulatory ecosystems.

### **3. Methodology**

#### **3.1 Data Collection**

The study begins by compiling a wide range of legal and regulatory documents from multiple jurisdictions to ensure broad coverage of the legal language. Guess what? do you know Sources include, include US federal regulations, the Federal Register, SEC regulations, EU directives, European Central Bank and SEC guidelines, Reserve Bank of India publications and corporate compliance guidelines. Seriously, Publicly available annotated datasets such, such as Lex GLUE, GLUE EU-LEGIS and other legal NLP standards is integrated to support supervised learning tasks. The study begins by compiling an extensive collection of legal and regulatory documents from multiple jurisdictions to ensure broad and representative coverage of legal language. Sources include U.S. federal regulations, the Federal Register, Securities and Exchange Commission (SEC) rules, European Union directives, European Central Bank guidelines, Reserve Bank of India publications, and corporate compliance frameworks. This diverse corpus captures the linguistic, structural, and conceptual variations that shape modern legal systems, thereby providing a rich foundation for downstream analytical tasks.

In addition to raw regulatory texts, the study integrates publicly available annotated datasets such as Lex GLUE, GLUE EU-LEGIS, and other established legal NLP benchmarks to support

supervised learning tasks. Combining heterogeneous, multi-jurisdictional, and multi-format documents significantly enhances the model’s generalizability, enabling the system to accommodate varying legal traditions, interpretive standards, and regulatory obligations across different countries. This integrated corpus ensures that the resulting AI framework is robust, adaptable, and capable of addressing complex legal challenges in diverse governance environments. And oh yeah, Seriously, this mix of heterogeneous and multi-jurisdictional documents enhances the generalizability of the model, allowing the system to take into account , account different legal and legal structures and regulatory obligations in different countries.

### 3.2 Data Preprocessing

Collected documents undergo extensive pre-processing to handle the complexities of formal legal language. And yes, the text is coded and divided into meaningful units, and lemmatization is used to reduce linguistic variation. Domain-specific Named Entity Recognition (NER) models , models are used to identify legal entities, regulatory authorities, obligations, monetary values, dates, and legal and legal references. You know what? Dedicated legal dictionaries and glossaries are OF great help in resolving ambiguities and reducing false positives in entity extraction. Guess what? Additional preprocessing steps include removing redundant content, normalizing content, and handling long document structures such as multi-layered and nested sentences. Fig2 contain graphic abstract Collected documents undergo extensive preprocessing to address the complexity and formality of legal language. The text is segmented into meaningful analytical units, and lemmatization is applied to reduce linguistic variation while preserving legal intent. Domain-specific Named Entity Recognition (NER) models are deployed to accurately identify legal entities, regulatory authorities, obligations, monetary values, dates, and statutory references. To further enhance precision, specialized legal dictionaries and glossaries are incorporated to resolve ambiguities, refine entity boundaries, and minimize false positives during extraction.

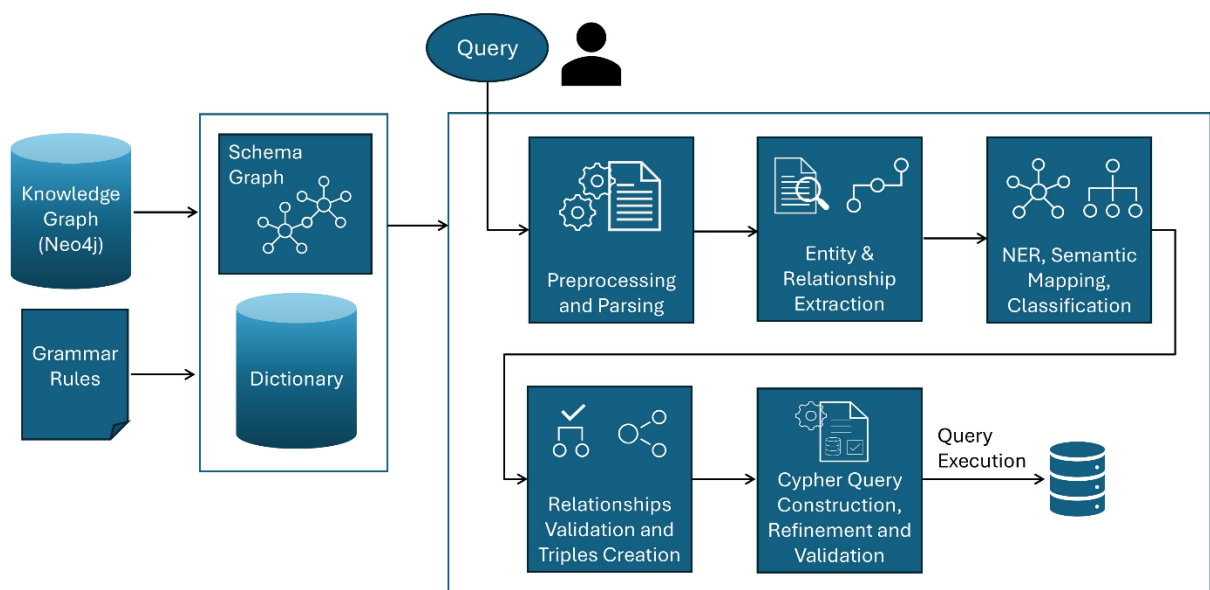


Figure 2 Graphical Abstract

Additional preprocessing steps include the removal of redundant or boilerplate content, normalization of inconsistent formatting, and systematic handling of long, structurally complex documents. This involves managing multi-layered sentences, nested clauses, and cross-referenced provisions—features that are particularly prevalent in statutes and regulatory guidelines. Together, these preprocessing procedures create a clean, consistent, and linguistically coherent dataset that supports more accurate downstream modeling and ensures that the analytical framework can effectively process both short and highly complex legal texts.

### 3.3 Model Development

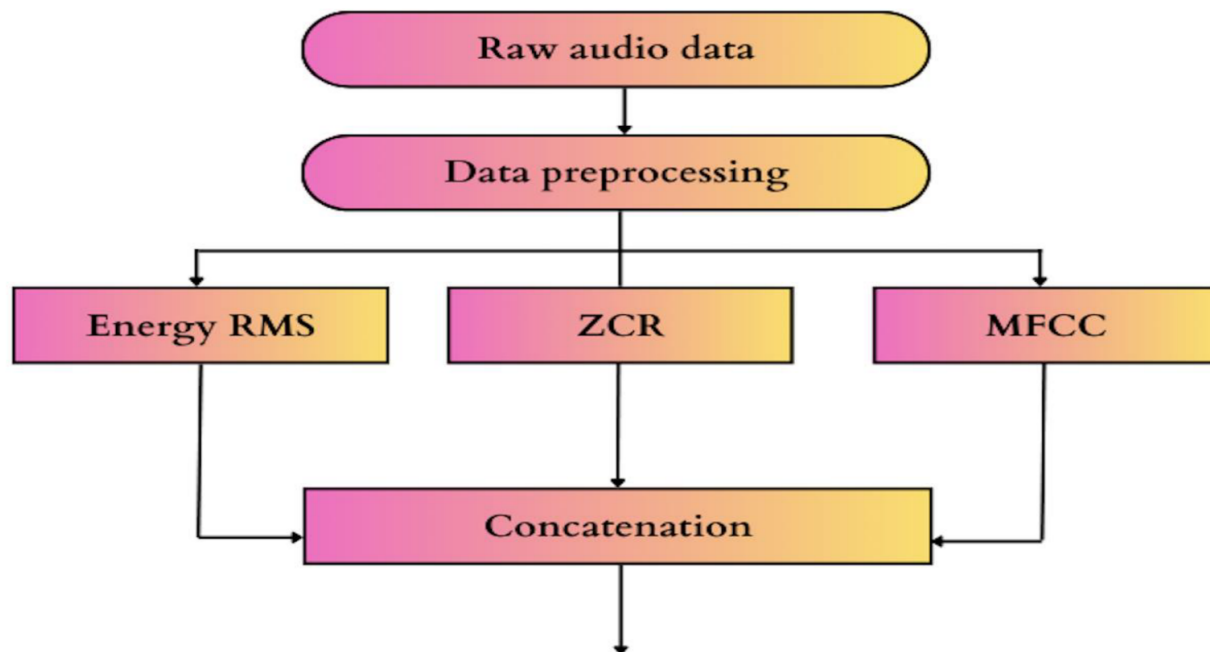
The model development integrates both deep learning and classical machine learning techniques to capture and capture syntactic and contextual-semantic cues, cues in legal and legal documents. Guess what? Transformer-based models such as BERT, BERT Legal-BERT Fin BERT, and Long Previous are tuned for tasks such as regulatory classification liability extraction and jurisdictional mapping. Conventional algorithms, such as SVM, are used as a complementary or combined component to logistic regression and gradient boosting to improve the robustness of smaller or specialized datasets. Like, This combination enables the system to process short and long legal content with high accuracy and comprehensibility. The model development phase integrates both deep learning and classical machine learning techniques to effectively capture the syntactic complexity and contextual-semantic nuances inherent in legal documents. Given the highly specialized nature of legal language characterized by long sentences, embedded clauses, and domain-specific terminology this hybrid approach ensures that multiple layers of linguistic meaning are represented and analyzed. By leveraging complementary strengths across methodologies, the system is designed to address the diverse analytical demands of legal interpretation.

Transformer-based architectures form the core of the pipeline, with models such as BERT, Legal-BERT, Fin BERT, and Long former fine-tuned for tasks including regulatory classification, liability extraction, entity recognition, and jurisdictional mapping. These models excel at capturing long-range dependencies and subtle semantic relationships, enabling the system to parse complex legal provisions and identify obligations, constraints, and compliance triggers with high precision. Their ability to process long documents without significant degradation in performance is particularly valuable for statutes, contractual clauses, and policy texts that extend across multiple pages.

### 3.4 Hybrid Modeling Architecture

The hybrid CNN-LSTM architecture is designed to handle scenarios where the training data is limited or domain specific Convolutional layers capture local legal expression patterns and repetitive sentence structures, while LSTM layers represent long-term dependencies between multi-sentence legal reasoning. For example, they use alerting mechanisms to highlight important terms, obligations or legal aspects. You know what?fig3 contain The process of refining audio begins with a series of structured steps. Initially, using the Audio Segment library, the raw audio is efficiently loaded into an object known as the Audio Segment. Once this step is completed, the system initiates the normalization phase, ensuring that each Audio Segment object is calibrated to a standardized level of +5.0 dBs to maintain consistency. Subsequently, it is imperative to

transform this object into an array composed of individual samples. This transformation is pivotal, as it sets the foundation for all the subsequent preprocessing measures. This hybrid design balances computational efficiency with the depth of semantic understanding required for compliance-oriented legal , legal text extraction and outperforms stand-alone models under domain-dependent conditions.



**Figure 3 Data preprocessing and feature extraction**

### 3.5 Linguistic and Semantic Analysis

Linguistic analysis is applied at a bunch of levels to improve interpretability and highlight precise legal meaning. Syntactic analysis is used to identify grammatical structures and relationships between sentences, while morphological analysis determines the properties of words, such as tense, number, and conjugation. Like, Contextual language modeling helps , helps explain ambiguous or complex wording often found in legal documents. These layers enable the system to define basic semantic relationships, interpret reference obligations, and map hierarchical dependencies between legal provisions. Linguistic analysis is applied at multiple levels to enhance the interpretability of model outputs and ensure that the system accurately captures the precise meaning embedded in legal texts. These linguistic layers help bridge the gap between raw computational predictions and the nuanced reasoning required for legal interpretation. By incorporating linguistic insights, the framework becomes more transparent, more explainable, and better aligned with established legal analysis practices.

At the syntactic level, the system analyzes grammatical structures, dependencies, and inter-sentence relationships to identify how obligations, conditions, and exceptions are articulated within complex legal provisions. Morphological analysis further refines this understanding by examining word-level properties such as tense, number, conjugation, and derivation which often determine the temporal scope, applicability, or legal force of a clause. These analyses help distinguish mandatory duties from optional actions and clarify subject–object relationships within regulatory text.

Beyond structural processing, contextual language modeling is employed to disambiguate intricate or ambiguous wording commonly found in statutes, contracts, and regulatory guidelines.

### 3.6 Compliance and Regulatory Extraction Framework

The purpose of the extraction framework is to identify obligations prohibitions compliance conditions and regulatory requirements through documents. Using fine-tuned NER models layers of attention and hybrid architectures the system identifies important provisions and compares them to institutional policies to uncover gaps or inconsistencies. The framework supports jurisdictional mapping , mapping allowing comparison of regulations in the US EU and India. It also highlights potential non-compliance by aligning extracted commitments with policy datasets and synthetic compliance checklists. The proposed extraction framework is designed to automatically identify obligations, prohibitions, compliance conditions, and regulatory requirements embedded within legal and policy documents. Fig 4 contain The National Institute of Standards and Technology (NIST) recently released the updated CSF 2.0, and it’s a game-changer for organizations of all sizes. Whether you’re a small school or a large corporation, this framework can elevate your cybersecurity posture. By leveraging fine-tuned named entity recognition (NER) models, multi-layer attention mechanisms, and hybrid deep learning architectures, the system isolates critical legal provisions with high precision. These extracted elements are then evaluated against institutional policies to detect gaps, inconsistencies, or areas requiring regulatory alignment. This capability is essential for organizations seeking to automate compliance monitoring and reduce manual review workloads.

In addition to clause-level extraction, the framework incorporates jurisdictional mapping functionalities, enabling systematic comparison of regulatory requirements across the United States, European Union, and India. This multi-jurisdictional perspective allows the model to capture variation in legal terminology, structural patterns, and compliance obligations. The framework also supports risk identification by aligning extracted commitments with curated policy datasets and synthetic compliance checklists, thereby highlighting potential non-compliance scenarios and improving the accuracy of regulatory audits.

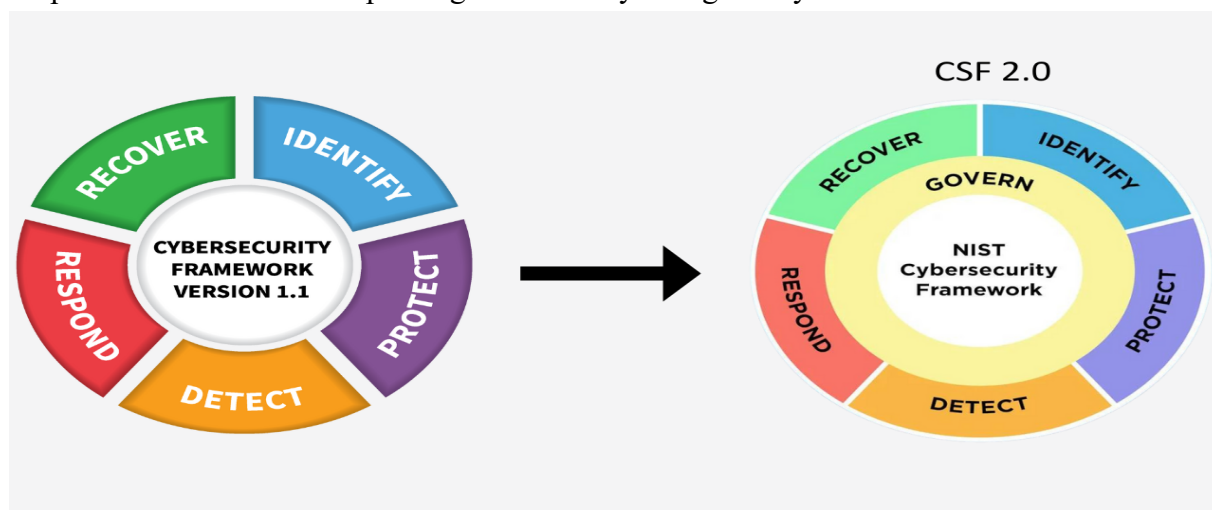


Figure 4 - The NIST Cybersecurity Framework (CSF) 2.0

### 3.7 Evaluation Metrics and Performance Assessment

Model performance is evaluated using precision, recall, and F1 score to measure accuracy in

detecting regulatory clauses and obligations. The domain-specific compliance accuracy metric assesses whether the model can correctly identify regulatory discrepancies in policy documents. Cross-validation ensures statistical reliability, while ABLATION studies measure the effect of preprocessing steps, model components, or language features on performance. The comparative analysis with the underlying models enables further validation of the effectiveness of the framework. Model performance is evaluated using established classification metrics such as precision, recall, and F1-score, which collectively measure the system's accuracy in detecting regulatory clauses, obligations, and compliance-related entities. These metrics provide a balanced understanding of both the model's sensitivity and its ability to minimize false positives critical factors in legal and regulatory analytics. In addition, a domain-specific compliance accuracy metric is introduced to assess the model's capability to correctly identify regulatory discrepancies and interpret compliance requirements within policy documents.

To ensure statistical reliability and robustness, the evaluation employs rigorous cross-validation procedures across multiple datasets and jurisdictions. This approach minimizes the risk of overfitting and confirms that performance gains are consistent rather than dataset-dependent. Cross-validation is particularly important when dealing with heterogeneous legal corpora, where linguistic structures, terminologies, and regulatory frameworks vary widely across regions and document types.

Ablation studies are conducted to quantify the contribution of individual preprocessing steps, model components, and linguistic features to overall performance. By systematically removing or modifying components, the study isolates the role of each element in shaping model accuracy and interpretability.

### **3.8 Survey-Based System Requirement Analysis**

To ensure real-world relevance, a survey of BRAZILIAN public administrations is conducted to identify current jurisprudential research systems, user expectations, and operational challenges. The survey collects data on functionality, search options, limitations and user satisfaction. Guess what? Insights from the 107 participating agencies help refine system requirements and guide the design, design of legal text mining tools that address practical issues such as search accuracy, multilingualism needs, and automated legal classification.

### **3.9 Research Design and Literature Review Process**

A structured literature review is conducted to analyze current legal NLP methods, case law retrieval systems and regulatory analysis frameworks. The review reveals shortcomings such as LIMITED explanatory capacity, insufficient datasets across jurisdictions and difficulties in processing long documents. These ideas shape the methodological structure of this study. Guess what? The research follows a mixed-methods approach, combining empirical modeling, linguistic analysis, and survey-based qualitative insights to create a comprehensive and well-founded canonical framework for text mining.

A structured literature review is conducted to examine the current landscape of legal NLP methodologies, case law retrieval systems, and regulatory analysis frameworks. The review highlights persistent limitations across existing approaches, including restricted explain ability,

limited availability of multi-jurisdictional datasets, and substantial challenges in processing long, complex legal documents. These gaps underscore the need for more robust, interpretable, and scalable techniques capable of addressing the structural and semantic intricacies of legal language. The insights gained from this review directly inform the methodological design of the present study.

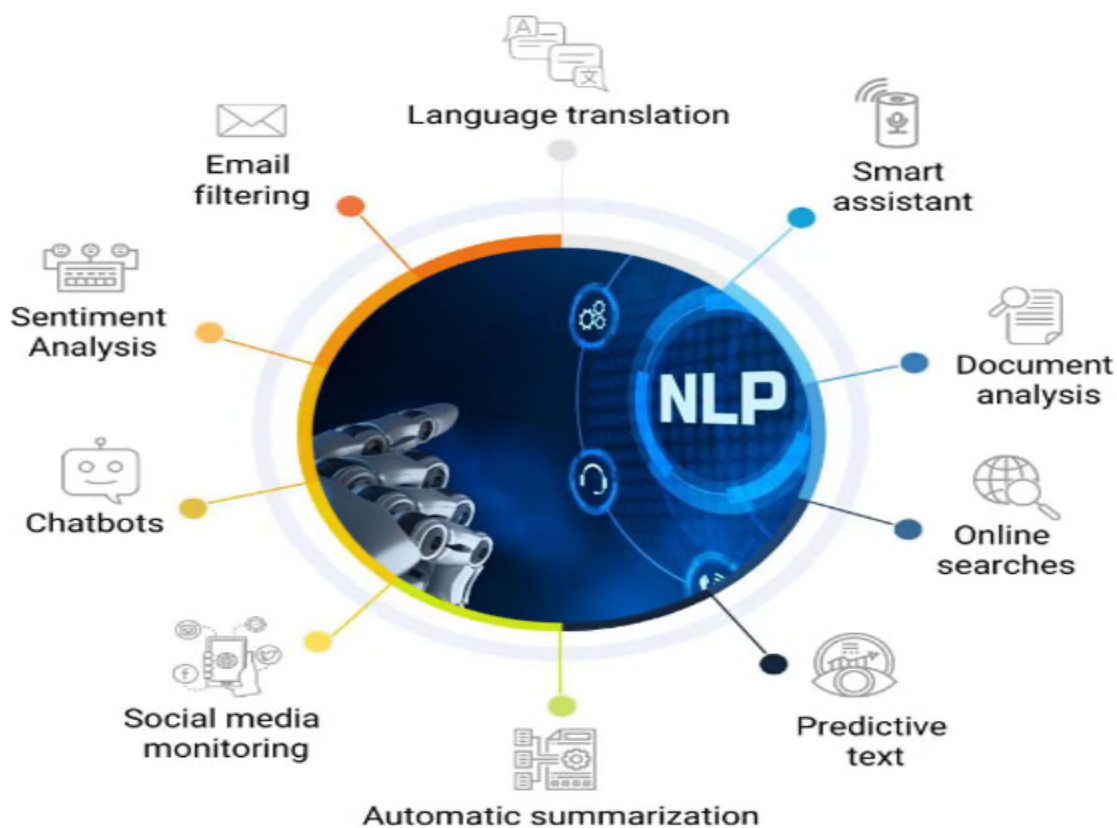
Building on these findings, the research adopts a mixed-methods approach that integrates empirical modeling, linguistic analysis, and qualitative insights gathered through expert surveys and practitioner evaluations. This multidimensional strategy ensures that both the computational and interpretative aspects of legal text processing are rigorously addressed. By combining quantitative performance assessment with qualitative validation, the study develops a comprehensive and well-grounded canonical framework for legal text mining that is both technically sound and contextually relevant to real-world legal practice.

#### 4. Data Analysis and Results

The data analysis for this study involved evaluating the performance, architecture, and practical application of an NLP-driven system designed to automate legal interpretation, compliance detection, and policy analysis. The system follows a multi-layered architecture that begins with the ingestion of regulatory documents, case law, regulatory guidelines, and large-scale legal repositories. These inputs are processed through a pre-processing process and sent to transformer-based NLP models, including Legal-BERT, Fin BERT, and Roberta variants specifically tailored to legal and regulatory language. Guess what? The models performed basic analysis tasks such as text classification, named entity recognition, obligation extraction, sentence segmentation, and semantic role tagging.

The results were then interpreted by the compliance analysis module, that identified the legal obligations associated with , with the relevant regulatory policies and highlighted areas of potential non-compliance or increased risk exposure.

The user interface layer translated the system's internal outputs into accessible visualizations. Compliance officers and legal practitioners used , used interactive dashboards to view flagged items, examine risk scores, and generate automated compliance reports. Heatmaps, summaries, and element-level traceability allowed decision makers to quickly explore how each extracted requirement related to specific laws or policies. This interface also PLAYED a central , central role in the human feedback mechanism: expert corrections and comments were fed back into the system to continuously improve , improve accuracy through reinforcement learning, improving end-user confidence and transparency.



**figure 5 NLP Applications**

To support advanced reasoning, the system integrates existing legal knowledge graphs (LKGs) such as Leafroll and LKIF-Core. Like, LKG provided a deeper contextual understanding by modeling relationships between legal entities, temporal aspects of regulation, cross-border maps and legal definitions. Combined with rule-based reasoning and machine learning, this framework enabled , enabled the system to identify INCONSISTENCIES, enforce compliance paths, and automatically suggest corrective actions. For example, if an organization's internal data does not reflect the provisions on user rights required by the GDPR, the system will flag the deletion and recommend specific changes.

Performance evaluations showed great potential in real-world application, but also highlighted current , current limitations .In the current case study environment, a bunch of limitations affected the effectiveness of the system. Like, These included the challenges of applying statistical and stochastic artificial intelligence methods to legacy search systems, the latency of multilevel access controls and LARGE indexed data sets, and the limited capacity of free visualization components used in early prototypes. Seriously, Plus, analytical exploration of aggregated data , data and robust ontological modeling of entity relationships essential for deeper legal insight were still under development. However, these limitations drove the system development roadmap, prompting the incorporation of scalable cloud , cloud modules, improved visualization processes, and extended ontological reasoning components.

Despite these challenges the results show that the proposed NLP architecture significantly improves the automation of legal text mining and compliance interpretation. Like The system successfully extracts and structures complex regulatory obligations supports real-time analytics

and enables organizations to rapidly , rapidly adapt to changing regulatory requirements. Fig5 contain Another crucial application is sentiment analysis, which gauges the emotions and opinions expressed in text. This is invaluable for businesses in understanding customer feedback, monitoring social media sentiment, and conducting market research.

Machine translation, a widely recognized application, allows for automated translation between different languages. NLP techniques, especially with the advent of neural machine translation models, have greatly improved translation accuracy, enabling effective communication across language barriers. Combining advanced NLP models cognitive graphs and human expert feedback creates a reliable and scalable framework for modern legal compliance systems.

Studies related to resources such as datasets or tools , tools accounted for 14 articles (18.67%). Although their number is smaller they still play an important supporting role in the research ecosystem. Multilingual research appears to be very rare as only one paper exists (1.33%). And oh yeah This highlights a significant gap most existing studies focus on a single language often English leaving , leaving multilingual and low-resource languages underrepresented. The MBA-based research was a great , great success with 46 theses , theses (61.33%). This , This reflects the growing excitement and confidence in large language models to solve complex problems.

## 5. Conclusion

This paper highlights how advanced natural language processing (NLP) can fundamentally change the understanding and management of legal documents, financial regulations and compliance guidelines. And oh yeah, Using specialized transformer models such as Legal-BERT, Fin BERT, and other domain-tuned architectures, the system demonstrates a strong ability to extract regulatory obligations, classify legal terms, and identify potential compliance risks , risks across jurisdictions. These improvements not only reduce the time and effort spent on manual review, but also support real-time tracking of regulatory changes, helping organizations meet rapidly changing legal requirements. Meanwhile, a more comprehensive review of existing legal text mining research shows that traditional machine learning methods such as SVM, Random Forest, KNN, and logistic regression continue to perform admirably in a bunch of legal prediction tasks. In some cases, it even outperforms modern , modern deep learning methods. However, the growing shift toward , toward deep learning models , models such as LSTMs, belt loops, CNNs, and transformers reflects the growing need for systems that can handle the rich , rich context, fine-grained semantics, and long-term dependencies that characterize legal , legal language. Optimized converter variants such as Legal-BERT and other domain-specific BERT models are particularly promising for multilingual and multi-jurisdictional applications .

## REFERENCES

1. B. Zhong, H. Wu, H. Li, S. Sapsagos, H. Luo, and L. He, “A scient metric analysis and critical review of construction related ontology research,” *Atom. Construct.*, vol. 101, pp. 17–31, May 2019.
2. A. S. Ismail, K. N. Ali, and N. A. IA had, “A review on BIM-based automated code compliance checking system,” in *Proc. Int. Conf. Res. Innova. Inf. Syst. (ICRIIS)*, Langkawi Island, Malaysia, Jul. 2017, pp. 1–6.

3. C. Eastman, J.-M. Lee, Y.-S. Jong, and J.-K. Lee, "Automatic rule-based checking of building designs," *Atom. Construct.*, vol. 18, no. 8, pp. 1011–1033, Dec. 2009.
4. P. Pauwels, D. Van Dearden, R. Virstatin, J. De Roo, R. De Meyer, R. Van de Wale, and J. Van Campano, "A semantic rule checking environment for building performance checking," *Atom. Construct.*, vol. 20, no. 5, pp. 506–518, Aug. 2011.
5. S. Jiang, Z. Wu, B. Zhang, and H. Cha, "Combined MvdXML and semantic technologies for green construction code checking," *Appl. Sci.*, vol. 9, no. 7, p. 1463, Apr. 2019.
6. M. Fahad and N. B. Bus Fees, "Semantic BIM reasoner for the verification of IFC models," in *Work and Business in Architecture, Engineering and Construction*, J. Karlson and R. Scherer, Eds. Boca Raton, FL, USA : CRC Press, 2018, pp. 361–368.
7. P. Pauwels and S. Zhang, "Semantic rule-checking for regulation compliance checking: An overview of strategies and approaches," in *Proc. 32rd Int. CIB W78 Conf.*, Eindhoven, The Netherlands, 2015, pp. 619–628.
8. S. Marital and H. M. Guanyin, "Computer representation of building codes for automated compliance checking," *Atom. Construct.*, vol. 82, pp. 43–58, Oct. 2017.
9. R. Sacks, "Automating design review with artificial intelligence and BIM: State of the art and research framework," in *Computing in Civil Engineering 2019: Visualization, Information Modeling, and Simulation*, Y. K. Cho, Ed. Reston, VA, USA : ASCE Press, 2019. 353–360.
10. M. M. Hossain and S. Ahmed, "Developing an automated safety checking system using BIM: A case study in the Bangladeshi construction industry," *Int. J. Construct. Manage.*, vol. 4, pp. 1–19, Nov. 2019.
11. J. Zhang and N. M. El-Gohar, "Integrating semantic NLP and logic reasoning into a unified system for fully-automated code checking," *Atom. Construct.*, vol. 73, pp. 45–57, Jan. 2017.
12. P. Zhou and N. El-Gohar, "Ontology-based automated information extraction from building energy conservation codes," *Atom. Construct.*, vol. 74, pp. 103–117, Feb. 2017.
13. J. Zhang and N. M. El-Gohar, "Semantic NLP-based information extraction from construction regulatory documents for automated compliance checking," *J. Compute. Civil Eng.*, vol. 30, no. 2, Mar. 2016, Art. no. 04015014.
14. D. M. Salama and N. M. El-Gohar, "Semantic text classification for supporting automated compliance checking in construction," *J. Compute. Civil Eng.*, vol. 30, no. 1, Jan. 2016, Art. no. 04014106.
15. J. Zhang and N. M. El-Gohar, "Extending building information models semiautomatic ally using semantic natural language processing techniques," *J. Compute. Civil Eng.*, vol. 30, no. 5, Sep. 2016, Art. no. C4016004.
16. C. Periled and A. Bormann, "Automated code compliance checking based on a visual language and building information modeling," in *Proc. 32nd ISARC*, Oulu, Finland, 2015, pp. 1–8.
17. T. Bloch and R. Sacks, "Comparing machine learning and rule-based inferencing for semantic enrichment of BIM models," *Atom. Construct.*, vol. 91, pp. 256–272, Jul. 2018.
18. N. O. Newari, "Smart codes and BIM," in *Proc. Conf. Struct. Conger.*, Pittsburgh, PA, USA, 2013, pp. 928–937.
19. E. Helmet and N. N. Nisbet, "Capturing normative constraints by use of the semantic mark-up rase methodology," in *Proc. CIB*, Sophia Antipolis, France, 2011, pp. 1–10.

20. J. K. Lee, "Building environment rule and analysis (BERA) language," Ph.D. dissertation, College Archit., Georgia Inst. Technol., Atlanta, GA, USA, 2011.
21. B. Ilhan and H. Yamen, "Green building assessment tool (GBAT) for integrated BIM-based design decisions," *Atom. Construct.*, vol. 70, pp. 26–37, Oct. 2016.
22. W. Soiling, "A simplified BIM data representation using a relational database schema for an efficient rule checking system and its associated rule checking language," Ph.D. dissertation, College Archit., Georgia Inst. Technol., Atlanta, GA, USA, 2016.
23. S.-L. Fan, H.-L. Chi, and P.-Q. Pan, "Rule checking interface development between building information model and end user," *Atom. Construct.*, vol. 105, Sep. 2019, Art. no. 102842.
24. M. Fahad and F. Adieux, "Towards mapping certification rules over BIM," in *Proc. 33rd CIB W78 Conf.*, Brisbane, QLD, Australia, 2016, pp. 1–10.
25. B. T. Zhong, L. Y. Ding, H. B. Luo, Y. Zhou, Y. Z. Hu, and H. M. Hu, "Ontology-based semantic modeling of regulation constraint for automated construction quality compliance checking," *Atom. Construct.*, vol. 28, pp. 58–70, Dec. 2012.