

Natural language processing in legal document analysis software: A systematic review of current approaches, challenges, and opportunities

^DMd Mostafijur Rahman^{1*}; ^DNajmul Gony Md²; ^DMd Mashfiquer Rahman³; ^DMd Mostafizur Rahman⁴; ^DMaria Khatun Shuvra SD⁵

¹Department of Computer Science & Engineering, RUET, Bangladesh, Rajshahi University of Engineering & Technology (RUET), Bangladesh.

²Department: Master of Science in Business Analytics, Grand Canyon University, USA.

³Department of Computer Science, Louisiana State University in Shreveport, USA.

⁴College of Technology & Engineering, Westcliff University, Irvine, California, USA.

⁵Department of Master of Science in Business Analytics, Grand Canyon University, USA.

Corresponding author: Md Mostafijur Rahman (Email: rubelruet@gmail.com)

Abstract

Natural Language Processing (NLP) techniques have been integrated into legal software systems to address the increasing volume and complexity of legal content in regulatory compliance, litigation, and contract management. This systematic study examines the current advancements in NLP applications for legal document analysis, focusing on critical tasks such as contract appraisal, case law summarization, legal question answering, and compliance verification. Ten fundamental research papers published between 2019 and 2025 were selected from academic sources such as IEEE Xplore, ACM Digital Library, SpringerLink, arXiv, and others utilizing the PRISMA methodology. The paper highlights the evolution from rule-based and statistical models to deep learning architectures and large language models (LLMs) tailored for legal text, such as Legal-BERT and GPT-based systems. Despite the potential of NLP in legal practice to automate monotonous tasks and enhance legal reasoning, significant challenges persist. These include the absence of annotated legal datasets, difficulties in interpreting domain-specific terminology, model bias, insufficient output transparency, and ethical concerns over automation in critical sectors. Numerous systems also exhibit a deficiency in explainability, undermining regulatory approval and trust. This work encapsulates current achievements, evaluates model performance on common legal NLP tasks, and highlights significant gaps and future research paths. It facilitates the development of legally competent, auditable, domain-adaptive NLP systems that seamlessly integrate into judicial and commercial legal procedures.

Keywords: Artificial Intelligence in law, Information extraction, Legal document analysis, Legaltech, Natural language processing (NLP), Text mining.

Funding: This study received no specific financial support.

History: Received: 15 April 2025 / Revised: 19 May 2025 / Accepted: 21 May 2025 / Published: 10 June 2025

Copyright: @ 2025 by the authors. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (https://creativecommons.org/licenses/by/4.0/).

Competing Interests: The authors declare that they have no competing interests.

Authors' Contributions: All authors contributed equally to the conception and design of the study. All authors have read and agreed to the published version of the manuscript.

DOI: 10.53894/ijirss.v8i3.7702

Transparency: The authors confirm that the manuscript is an honest, accurate, and transparent account of the study; that no vital features of the study have been omitted; and that any discrepancies from the study as planned have been explained. This study followed all ethical practices during writing. Publisher: Innovative Research Publishing

1. Introduction

In the advancing realm of legal technology, the quantity, intricacy, and speed of legal paperwork have increasingly posed challenges for human practitioners to oversee [1]. Legal practitioners frequently engage with extensive statutes, contracts, case law, regulations, and regulatory documents that are replete with specialized terminology, intricate syntax, and contextual nuances. With the increasing demand for efficiency, accuracy, and cost reduction in legal services, Natural Language Processing (NLP) has emerged as a transformative technology that may automate and enhance numerous document-intensive legal processes. Natural Language Processing (NLP), a branch of artificial intelligence (AI) and computational linguistics, encompasses the automated analysis and creation of human language [2]. NLP tools in the legal area can facilitate several applications, such as automated contract analysis, legal question answering, document classification, case law summarization, and compliance verification. These technologies assist attorneys in identifying terms, duties, risks, precedents, and inconsistencies—tasks that would otherwise require considerable human labor. The emergence of deep learning and the current proliferation of large language models (LLMs), including GPT and BERT variations refined on legal texts, have elevated NLP-driven computers to unprecedented levels of semantic comprehension and contextual reasoning. Notwithstanding considerable progress, the utilization of NLP in the study of legal documents continues to be fraught with distinct obstacles [3]. Legal texts frequently exhibit ambiguity, require interpretation, and are context-dependent, markedly deviating from common language corpora. Furthermore, ethical and practical considerations-such as transparency, explainability, responsibility, and the risk of algorithmic bias-highlight the necessity for prudence, particularly in critical contexts such as judicial processes or regulatory adherence. The lack of high-quality annotated datasets exacerbates these issues, restricting the efficacy of supervised learning methods and the generalizability of pre-trained models. While numerous separate studies have investigated particular tasks and techniques in legal NLP, a thorough, systematic synthesis of the topic remains in development [4, 5]. A comprehensive study is necessary to unify disparate knowledge, evaluate existing technological competencies, and pinpoint significant deficiencies. This can function as a reference for future study and practical implementation by delineating the accomplishments, ongoing obstacles, and the most promising potential.

1.1. Understanding NLP and its Importance

Comprehending Natural Language Processing (NLP) and its importance in legal document analysis necessitates examining the convergence of technology, artificial intelligence (AI), data science, and contract interpretation. Instruments such as LexCheck are essential for automating the review and negotiation processes. NLP models in legal document analysis employ sophisticated algorithms to extract essential data points, sentences, and provisions from extensive text, facilitating lawyers in their tasks. Utilizing AI-driven solutions like LexCheck, legal practitioners can optimize contract evaluations and guarantee adherence to legal norms. The automated features of NLP accelerate the process and reduce human error, hence improving the efficiency and precision of legal document analysis. This groundbreaking technology has transformed the legal sector by offering innovative solutions that enhance workflow, lower costs, and diminish risks [6].

1.2. Data Science in Legal Document Examination

Data science in legal document analysis employs advanced technologies, such as NLP, AI, and machine learning, along with tools like LexCheck, to derive significant insights from intricate legal documents. The integration of data science techniques and advanced algorithms improves the efficiency and precision of document inspection and management procedures. LexCheck utilizes cutting-edge data science methodologies to transform legal document analysis. LexCheck employs AI and machine learning to automate the review of contracts, briefs, and other legal documents with accuracy and efficiency. These tools not only conserve time but also enhance the overall quality of work by minimizing human errors and ensuring adherence to legal norms. The incorporation of NLP facilitates the extraction of essential information, empowering legal practitioners to make prompt, informed choices. The integration of data science and sophisticated algorithms in legal document analysis enhances efficiency and accuracy [7].

This study aims to systematically examine Natural Language Processing applications in legal document analysis systems. It specifically examines what is explained in the given table:

Table 1.

Table 1	•
Natural	Language Processing applications in legal document analysis systems.
•	The range of NLP applications utilized in legal settings;

- The concepts and methodologies employed-from traditional algorithms to large language models;
- Accessibility and utilization of domain-specific data;
- Performance metrics and evaluation standards;

Technical, legal, and ethical challenges hindering acceptance; unresolved issues and prospective avenues for investigation.

The research guarantees rigor in selection, screening, and literary analysis through the PRISMA methodology. The findings should inform academic researchers and business executives, thereby bridging the divide between computational advancements and real legal requirements.

2. Methodology

To guarantee the validity and completeness of its conclusions, a systematic review calls for an open and repeatable research approach. Following the PRISMA (Preferred Reporting Items for Systematic Reviews and Meta-Analyses) guidelines, we defined, carried out, and reported our literature search, screening, inclusion criteria, and analytical methods for this study. The aim was to find and evaluate peer-reviewed scholarly works on the use of Natural Language Processing (NLP) in legal document analysis, namely those including software implementations, legal corpora, or pragmatic applications.

Table 1.

The following research questions (RQs) was developed to direct the review.							
RQ1	In legal document analysis, which kinds of NLP tasks are most often used?						
RQ2	Which models and approaches apply to carry out these chores?						
RQ3	What are the restrictions of which datasets facilitate legal NLP research?						
RQ4	Which main obstacles exist in implementing NLP into legal document processes?						
RQ5	Which research gaps and patterns can be found to direct next investigations?						

2.1. Search Strategy

We conducted a thorough search across five primary academic and scientific publication databases, utilizing Boolean operators in conjunction with keyword phrases.

- IEEE
- Xplore
- SpringerLink
- ACM Digital Library
- arXiv (computer science: artificial intelligence, computing, and linguistics)
- Google Books

Table 3.

Relevant keywords used for research. Natural Language Processing, Legal Document Analysis, Legal NLP, Contract Review, Legal AI, Systematic Review, Survey

Inclusion			xclusion
•	Articles must focus on the applications of NLP in the analysis of legal documents.	•	Articles beyond the legal domain (including general NLP and associated fields).
•	Must demonstrate either an innovative perspective, a comprehensive analysis, or an evaluation of the dataset/model.	•	Documents devoid of technical substance or just intellectual material that lacks practical application or assessment.
•	Either disseminated in reputable preprint repositories (e.g., arXiv with citations) or subjected to peer review.	•	Opinion pieces, white papers, or promotional content.
•	It is essential to include technical discussions, model architecture, or explanations of use cases.	•	Numerous studies or their more contemporary counterparts.

2.2. Inclusion And Exclusion Criteria

The established inclusion criteria ensured the pertinence and quality of the selected studies:

2.3. Selection Methodology

Studies were selected according to the PRISMA 2020 framework to ensure methodological rigor and transparency. The PRISMA flow diagram (Fig. #1) illustrates that the approach comprised four successive phases: Identification, Screening, Eligibility, and Inclusion. Initially, a search of primary academic databases, including IEEE Xplore, SpringerLink, ACM Digital Library, arXiv, and Google Scholar, yielded a total of 340 items. Grey literature searches and backward citation tracking yielded an additional 28 records, resulting in a total of 340 potential publications.

2.4. Data Extraction and Coding

Every chosen paper underwent hand coding depending on:

- 1. NLP chore (such as classification, summarizing, named entity recognition, etc.)
- 2. Legal use case (contract analysis, litigation support, etc.).
- 3. Technical model—e.g., BERT, GPT, SVM)
- 4. Ethical topics including justice, explainability, and prejudice

2.5. Co-Citation Mapping and Citations.

Reference lists were examined for citation connections and subsequently verified on Google Scholar. Quantifying the frequency of specific themes, such as "LLMs" and "Ethics" co-occurring in cited literature generated subject co-citation networks. Python's NetworkX library enables the visualization of networks.

2.6. Visualization Tools

- Visualizations generated with Python (Pandas, Matplotlib, NetworkX)
- PRISMA Framework: for the selection process of articles
- Algorithm for spring layouts applicable to reference and co-citation network graphs

2.7. Quality and Relevance Screening

Papers included:

- Specifically focused on NLP in legal domains.
- Precisely delineated methodologies and frameworks
- Implemented or proposed sets that provided either qualitative or quantitative evaluation

3. Results

This section consolidates the insights gathered from 60 selected peer-reviewed articles (2019–2025) [4-6, 8-94] analyzed using task distributions, model usage, datasets, ethical features, temporal trends, thematic evolution, and co-citation networks

3.1. Task Distribution in Legal NLP

Examining the sixty chosen studies reveals that thirty percent of the research focus is on contract analysis; hence it is the most prevalent chore. Usually, these studies consist of clause detection, risk extraction, and compliance validation. Following closely and reflecting needs in case filtering and legal information retrieval, legal document classification, and legal question answering (QA) because they can distill complicated case law and extract relevant legal phrases, tasks like legal summarization and named entity recognition (NER) are becoming more popular. Legal text generation and compliance verification are lesser-known but developing fields. These trends point to a field still oriented towards information extraction and automation of demanding legal review chores.

Distribution of NLP Tasts in Legal Document Analysis (n=60)



Distribution of NLP Tasts in Legal Document Analysis (n=60)

Figure 1.

Task Distribution in Legal Natural Language Processing.

3.2. Models and Techniques Used

The reviewed research predominantly utilized transformer-based models such as BERT and its legal counterparts (Legal-BERT, CaseLawBERT). These models are selected for their capacity for contextual embedding and their ability to be finetuned for domain-specific corpora. Legal summarization and the development of legal chatbots have gained popularity among GPT-style language models. Despite the increasing popularity of deep learning, traditional machine learning approaches such as SVM and decision trees remain prevalent, particularly in scenarios requiring interpretability and with limited datasets. In roles necessitating legal and regulatory validation with mandated explainability, rule-based systems continue to be relevant.



Number of Studies

Figure 2.

While advanced models offer power, rule-based and classical ML models remain essential for legal interpretability and traceability.

3.3. Data Use

Third of all studies depended on proprietary or closed legal corpora, suggesting a serious data access constraint. Among public databases, EUR-Lex and LexGLUE were the most often used ones, offering organized legal content fit for classification and search. Notable too were the CaseLaw corpora from the US and EU courts. However, relatively few studies presented or published fresh datasets, and generalizability is challenging due in large part to the lack of multilingual or low-resource datasets. This reflects a desire for increased cooperative data efforts among legal tech stakeholders as well as a significant field constraint.



Number of Studies Using Dataset

Figure 3. Number of Studies Using Dataset.

3.4. Analytical Characteristics Among Research

Despite their relevance in legal environments, only 15% of the studies directly address the explainability of the models, according to a thorough investigation of their analytical aspects. About twenty-three percent of the publications raised ethical questions including bias, openness, or artificial intelligence abuse in high-stakes choices. Dominant in the methodological

terrain, supervised learning highlighted the reliance of the discipline on labeled training data. Only a fraction of research conducted comparative model evaluations; few presented new standards or datasets. These shortcomings draw attention to the requirement for improved methodological rigor and openness in future legal NLP studies.



Features And Number Of Studies Using Them

Figure 3.

Feature and relative number of studies.

3.5. Temporal Mapping of Publications (2019–2025)

Publication temporal distribution shows a notable increase in legal NLP research beginning in 2020. The field peaked in 2024, most likely due to increased legal industry interest and the general acceptance of LLMs such as GPT-3 and GPT-4. While recent years progressed toward more complicated, ethical, and scalable solutions, earlier years included conventional NLP and rule-based techniques. Stabilization and field maturity help to explain the modest decline in 2025, implying a shift from exploratory research to implementation-oriented, policy-aware study.

Temporal Mapping of Legal NLP Publications (2019–2025)

Temporal Mapping of Legal NLP Publications (2019–2025)



Figure 5.

Temporal trend of publications on Legal Natural Language Processing (NLP) from 2019 to 2025, showing a steady rise from 2019 to 2024 with a peak of 13 publications, followed by a slight decline in 2025.

3.6. Thematic Evolution

From rule-based systems and categorization (2019–2020) to BERT-based modeling and information retrieval (2021–2022), and lastly, generation tasks and socio-ethical consequences (2023–2025), themes in legal NLP have emerged. With important connections to subjects including summarization, ethics, explainability, and multilingual legal data, LLMs have become a central theme in recent years. Thematic development reflects both technological advances and growing entanglement with legal, regulatory, and justice problems. It also shows how research has progressed from basic ideas to more integrated, practical legal solutions.





3.7. Topic Co-Citation Network

Clearly, thematic clustering is shown by the co-citation network built from the 60 studies. Often referenced together as a technology-methodology cluster were LLMs, summarization, and classification. Frequently found in articles addressing risks and governance, ethics and explainability formed the strongest conceptual duo. Contract analysis highlighted common data and task design alongside both QA systems and NER. This network offers a knowledge framework of the discipline, therefore enabling the identification of major themes and intersections across lines of methodology, ethics, and application.



Topic Co-Citation Network in Legal NLP Research



Bias Mitigation -	0	0	0	0	0	3	0	0	0	0	0	2	0	0		6
Case Law Summarization -	0	0	0	0	0	0	0	0	0	3	0	0	0	0		
Clause Extraction -	0	0	0	6	0	0	0	2	0	0	0	0	0	5	-	5
Contract Analysis -	0	0	6	0	0	0	0	0	0	0	0	0	4	0		
Ethical AI -	0	0	0	0	0	5	6	0	0	0	0	0	0	0	-	4
Explainability -	3	0	0	0	5	0	0	0	0	0	0	0	0	0		
GPT Models -	0	0	0	0	6	0	0	0	4	0	0	0	0	0		С
Legal Datasets -	0	0	2	0	0	0	0	0	3	0	0	0	0	0		2
Legal QA -	0	0	0	0	0	0	4	3	0	0	4	0	0	0		
Legal Question Answering -	0	3	0	0	0	0	0	0	0	0	0	0	0	0	-	2
Legal Summarization -	0	0	0	0	0	0	0	0	4	0	0	0	0	0		
Multilingual Legal NLP - 1	2	0	0	0	0	0	0	0	0	0	0	0	0	0	-	1
Named Entity Recognition -	0	0	0	4	0	0	0	0	0	0	0	0	0	0		
Regulatory Compliance -	0	0	5	0	0	0	0	0	0	0	0	0	0	0		0
Case Law Sun Clause Contract Analysis Contract A																

Heatmap of Topic Co-Citation Correlations in Legal NLP Research

Figure 8.

Heatmap of topic cocitation corelations.

3.8. Reference Historiograph (Conceptual Lineage)

From 2019 to 2025, the citation historiography charts conceptual dependencies in the field. Foundational articles from 2019 and 2020 set the stage with early categorization problems and rule-based systems. Important books between 2021 and 2023 bridged basic work with the new powers of big language models. By 2024–2025, the emphasis has turned to compliance, fairness, benchmarking, and LLM audits. This shows how the industry is gradually moving toward legally robust and explainable AI models and reflects a cumulative intellectual conversation.

Historiograph Citation Network (2019-2025)



Citation Historiograph – Conceptual Lineage of Legal NLP Research.

Cluster Map of Topic Co-Citation Relationships in Legal NLP

Figure 10.

Cluster Map of co citation relationship of legal NLP.

4. Discussions

The results of this systematic review reveal an increasingly complex and established body of research on the use of natural language processing (NLP) in the study of legal documents. By integrating information from the 60 papers that were reviewed, this discussion takes into account evolving methods, strategies, problems, and gaps that influence the discipline's current course.

The sector has seen a steady expansion in NLP operations tailored to legal situations. The goal of the early efforts, which focused mostly on contract analysis and document classification, was to reduce the amount of manual labor required to tag and organize legal documents. As the capabilities of NLP models increased, so did the tasks, with increasingly sophisticated applications such as compliance verification, case law summarization, and legal question answering (QA) becoming more prevalent. This diversity draws attention to an important trend: the legal field is not a monolith, and NLP systems need to adjust to the unique linguistic, structural, and interpretive challenges that various legal activities provide. For example, clause extraction in contract review is completely different from quality assurance in a court judgment setting. Thus, we show increasing task specialization driven by both dataset availability and use-case requirements. The evolution and advancement of transformer-based models, including BERT, Legal-BERT, and GPT-series models, have transformed legal text processing. These models dominate the reviewed literature not only for their accuracy but also because they enable transfer learning from general to legal-specific corpora, which is a significant advantage in a low-resource setting. While GPT-2 and GPT-3 began to enter legal QA and summarizing systems in 2021–2022, GPT-4 and other instruction-tuned LLMs were particularly subject to increased scrutiny in 2024 and 2025. The efficacy of these models was praised, but they were continuously criticized for

their opacity, potential bias, and vulnerability to hallucinations, which were especially problematic in the legal arena. Concurrently, older models such as SVMs, decision trees, and rule-based systems continue to be in demand, especially in situations where interpretability and explainability are crucial. Their ongoing use suggests that model complexity is not necessarily consistent with legal acceptability, especially in domains where results must be clearly justified. Throughout the reviewed literature, the lack of publicly available, high-quality legal datasets appears frequently. Despite the fact that many studies rely on databases like EUR-Lex, LexGLUE, and CaseLaw, nearly one-third of all the research that was reviewed relied on proprietary or private corpora, which limited benchmarking and repeatability.

Geographical and jurisdictional biases are also introduced by this reliance on custom datasets, as many of them are sourced from the legal systems of the United States or the European Union. The absence of multilingual and culturally diverse corpora hinders the development of globally adaptable legal natural language processing systems. The shortage of labeled training data has led to the underutilization of supervised learning in a number of high-value sectors (such as privacy regulations and regulatory monitoring), despite the growing availability of unsupervised and self-supervised methodologies. The ethical implications of legal artificial intelligence continue to receive insufficient funding, with only 15% of studies directly addressing explainability and 23% posing concerns about bias and fairness. The requirement that judicial judgments be auditable and interpretable makes this a critical oversight, particularly in nations where explainability is required by law. Despite their strengths, models such as GPT-3 and GPT-4 have epistemic dangers since they may produce information that is legally flawed but sounds plausible. This raises concerns about their suitability for positions involving drafting contracts or making recommendations to judges. Additionally, bias in training data, such as skewed case findings or an overrepresentation of particular jurisdictions, may have detrimental real-world consequences, such as perpetuating systemic prejudice. More precise model evaluation that goes beyond F1-score or accuracy is required by the recent discussion on fair artificial intelligence in law. Fairness audits, counterfactual testing, and assessments of legal compliance ought to be among the metrics. The trend from 2019 to 2025 suggests a shift away from exploratory, rule-based research and toward more indepth, interdisciplinary work that incorporates ethical, policy, and regulatory frameworks. Publications surged in 2024, propelled by stronger LLMs and a heightened awareness of AI governance in society. This approach is confirmed by the thematic evolution map: earlier years focused on classification and information extraction, whereas recent studies stress LLM integration, multilingual compliance, and legal AI governance. Additionally, exposing clusters of interest in issues like contract QA, LLMs and ethics, and explainability in compliance is theme co-citation analysis. A comprehensive conceptual framework supporting legal NLP is revealed via the topic co-citation network and citation historiograph. Foundational publications from 2019-2021 continue to influence subsequent work, particularly those that provided domain-specific BERT models or structured datasets.

Issues like LLMs and ethics, or summarizing and QA, are frequently mentioned simultaneously, indicating methodological and conceptual overlap, according to co-citation analysis. These intersections suggest that future research may benefit from composite benchmarks that represent real-world legal procedures or multitask models. Despite growing sophistication, there are still a number of significant gaps in the field, including a lack of formal benchmarks for justice and legal compliance, a lack of cross-jurisdictional and multilingual datasets, inadequate interpretability frameworks appropriate for legal use cases, and underdeveloped techniques for real-time legal decision support. The development of auditable LLMs with modular architectures, the production of multilingual legal datasets (for example, through expert annotation and synthetic generation), the use of NLP for legislative forecasting in legal policy simulation, and the use of causal inference in legal prediction tasks are all examples of future research opportunities. The field of legal NLP has undergone rapid yet uneven change. Despite the incredible powers granted by technological advancements, questions of justice, validity, and trust still need to be addressed. Instead of pursuing complexity for its own sake, this methodical evaluation identifies a path forward in the development of legally informed, human-centered, auditable NLP systems tools that serve law not just as text but also as a domain of societal consequence.



Figure 11.

This macro visual map shows how all important elements taken from a systematic review are interdependent. Every node relates to a fundamental idea such a task, model, dataset, analytical feature, year of publication, research theme, co-citation, or citation flow. The colors help distinguish these categories. Arrows between nodes show the conceptual or methodological relationships that structure the development of the legal NLP field over time.

5. Conclusions

The research terrain surrounding the use of Natural Language Processing (NLP) in legal document analysis between 2019 and 2025 is thoroughly explored in this systematic review. Across 60 peer-reviewed publications, it reveals a fastexpanding yet unevenly distributed discipline formed by developments in machine learning, legal informatics, and ethical artificial intelligence. From contract analysis and legal question answering to summarizing and compliance checking, the evaluation notes that the legal industry has embraced an expanding range of NLP tasks. These advances align with actual legal difficulties, including access to justice, document overload, and regulatory complexity. Legal NLP systems' capabilities have been transformed by the general acceptance of transformer-based models, including Legal-BERT and GPT versions, therefore allowing contextual understanding and domain adaptability. The assessment does, however, also highlight significant congestion. Particularly in non-Western and multilingual legal environments, dataset shortages seriously restrict advancement. Many studies depend on proprietary or non-reproducible corpora, which reduces benchmarking and generalizability. Furthermore, underrepresented in both model design and evaluation are ethical dimensions—bias, fairness, and explainability. Given the great weight of legal decisions, when trust, openness, and responsibility are not only desirable but legally required, this is a major issue. Temporal and thematic studies show a positive trend: the discipline is moving from technical feasibility studies to more sophisticated research addressing governance, control, and multidisciplinary concerns. Growing co-citation of topics including ethics, LLM auditing, and legal artificial intelligence governance reflects this progress. In essence, even though NLP has great potential to transform legal procedures, responsible and efficient application of this tool calls for much more than modern models. It requires cooperative models among legal academics, technologists,

ethicists, and legislators. Future initiatives should give top priority to building open, diverse datasets; constructing legally interpretable models; and incorporating justice and transparency into every level of the artificial intelligence pipeline.

References

- B. Alshemali and J. Kalita, "Improving the reliability of deep neural networks in NLP: A review," *Knowl-Based Syst*, vol. 191, p. 105210, 2020. https://doi.org/10.1016/j.knosys.2019.105210
- [2] O. Baclic, M. Tunis, K. Young, C. Doan, H. Swerdfeger, and J. Schonfeld, "Artificial intelligence in public health: Challenges and opportunities for public health made possible by advances in natural language processing," *Canada Communicable Disease Report*, vol. 46, 2020. https://doi.org/10.14745/ccdr.v46i06a02
- [3] E. Blanzieri and A. Bryl, "A survey of learning-based techniques of email spam filtering," *Artificial Intelligence Review*, vol. 29, no. 1, pp. 63–92, 2008. https://doi.org/10.1007/s10462-009-9109-6
- [4] N. Choudhary, "LDC-IL: The Indian repository of resources for language technology," *Lang Resources & Evaluation*, vol. 55, pp. 553-566, 2021. https://doi.org/10.1007/s10579-020-09523-3
- [5] H. Chouikhi, H. Chniter, and F. Jarray, Arabic sentiment analysis using BERT model. In international conference on computational collective intelligence. Cham: Springer, 2021.
- [6] Y. Fan, F. Tian, Y. Xia, T. Qin, X. Y. Li, and T. Y. Liu, "Searching better architectures for neural machine translation," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 28, no. 1, pp. 1-11, 2020. https://doi.org/10.1109/TASLP.2020.2995270
- [7] H. Fang *et al.*, "Topic aspect-oriented summarization via group selection," *Neurocomputing*, vol. 149, pp. 1609-1616, 2015. https://doi.org/10.1016/j.neucom.2014.08.031
- [8] H. Aa, K. J. Balder, F. M. Maggi, and A. Nolte, Say it in your own words: Defining declarative process models using speech recognition. In D. Fahland, C. Ghidini, J. Becker, & M. Dumas (Eds.), Business Process Management Forum. Cham: Springer, 2020.
- H. Aa, C. Ciccio, H. Leopold, and H. A. Reijers, *Extracting declarative process models from natural language. In P. Giorgini* & B. Weber (Eds.), Advanced Information Systems Engineering. Cham: Springer. https://doi.org/10.1007/978-3-030-21290-2_23, 2019.
- [10] F. Ariai and G. Demartini, "Natural language processing for the legal domain: A survey of tasks, datasets, models, and challenges," *arXiv preprint arXiv:2410.21306*, 2024. https://doi.org/10.48550/arXiv.2410.21306
- [11] S. Arts, J. Hou, and J. C. Gomez, "Natural language processing to identify the creation and impact of new technologies in patent text: Code, data, and new measures," *Res Policy*, vol. 50, p. 104144, 2021. https://doi.org/10.1016/j.respol.2020.104144
- [12] M. Barrientos, K. Winter, J. Mangler, and S. Rinderle-Ma, "Verification of quantitative temporal compliance requirements in process descriptions over event logs," in *Proceedings of the Business Process Management Workshops, Cham, Springer. https://doi.org/10.1007/978-3-031-34560-9_25, 2023.*
- [13] R. E. Beaty and D. R. Johnson, "Automating creativity assessment with semdis: An open platform for computing semantic distance," *Behav Res Methods*, vol. 53, pp. 1064–1077, 2021. https://doi.org/10.3758/s13428-020-01453-w
- [14] P. Bellan, M. Dragoni, and C. Ghidini, Extracting business process entities and relations from text using pre-trained language models and in-context learning. In J. P. A. Almeida, D. Karastoyanova, G. Guizzardi, M. Montali, F. M. Maggi, & C. M. Fonseca (Eds.), Enterprise Design, Operations, and Computing. Cham: Springer, 2022.
- [15] S. Casola and A. Lavelli, "Summarization, simplification, and generation: The case of patents," *Expert Systems with Applications*, vol. 205, p. 117627, 2022. https://doi.org/10.1016/j.eswa.2022.117627
- [16] T. Chakraborti, Y. Rizk, V. Isahagian, B. Aksar, and F. Fuggitti, From natural language to workflows: Towards emergent intelligence in robotic process automation. In A. Marrella (Ed.), Business Process Management: Blockchain, Robotic Process Automation, and Central and Eastern Europe Forum. Cham: Springer, 2022.
- [17] I. Chalkidis *et al.*, "LexGLUE: A benchmark dataset for legal language understanding in English," *arXiv preprint arXiv:2110.00976*, 2021. https://doi.org/10.48550/arXiv.2110.00976
- [18] H. Chen and W. Deng, "Interpretable patent recommendation with knowledge graph and deep learning," *Scientific Reports,* vol. 13, 2023. https://doi.org/10.1038/s41598-023-28766-y
- [19] L. Chen, S. Xu, L. Zhu, J. Zhang, X. Lei, and G. Yang, "A deep learning based method for extracting semantic information from patent documents," *Scientometrics*, vol. 125, no. 3, pp. 2341–2365, 2020. https://doi.org/10.1007/s11192-020-03634-y
- [20] S. Choi, H. Lee, E. Park, and S. Choi, "Deep learning for patent landscaping using transformer and graph embedding," *Technol Forecast Soc Chang*, vol. 175, p. 121413, 2022. https://doi.org/10.1016/j.techfore.2021.121413
- [21] D. Christofidellis, M. M. Lehmann, T. Luksch, M. Stenta, and M. Manica, "Automated patent classification for crop protection via domain adaptation," *Applied AI Letters*, vol. 4, p. e80, 2023. https://doi.org/10.1002/ail2.80
- [22] J. M. Chu, H. C. Lo, J. Hsiang, and C. C. Cho, "From paris to le-paris: Toward patent response automation with recommender systems and collaborative large language models," *Artificial Intelligence and Law*, 2024. https://doi.org/10.1007/s10506-024-09409-7
- [23] P. Chung and S. Y. Sohn, "Early detection of valuable patents using a deep learning model: Case of semiconductor industry," *Technol Forecast Soc Chang*, vol. 158, p. 120146, 2020. https://doi.org/10.1016/j.techfore.2020.120146
- [24] M. Cosler, C. Hahn, D. Mendoza, F. Schmitt, and C. Trippel, "nl2spec: Interactively translating unstructured natural language to temporal logics with large language models," in *Proceedings of the International Conference on Business Process Management, Cham Springer. https://doi.org/10.1007/978-3-031-37703-7_18*, 2023.
- [25] W. Du, Y. Wang, W. Xu, and J. Ma, "A personalized recommendation system for high-quality patent trading by leveraging hybrid patent analysis," *Scientometrics*, vol. 126, no. 2, pp. 1713–1734, 2021. https://doi.org/10.1007/s11192-021-04180-x
- [26] V. Etikala, Z. Veldhoven, and J. Vanthienen, *Text2Dec: Extracting decision dependencies from natural language text for automated DMN decision modelling*. Cham: Springer. https://doi.org/10.1007/978-3-030-66498-5_27, 2020.
- [27] F. R. Golra, F. Dagnat, J. Souquières, I. Sayar, and S. Guerin, *Bridging the gap between informal requirements and formal specifications using model federation. In E. B. Johnsen & I. Schaefer (Eds.), Software Engineering and Formal Methods.* Cham: Springer, 2018.

- [28] A. Haghighian Roudsari, J. Afshar, W. Lee, and S. Lee, "Patentnet: multi-label classification of patent documents using deep learning based language understanding," *Scientometrics*, vol. 127, no. 2, pp. 603-621, 2022. https://doi.org/10.1007/s11192-021-04179-4
- [29] D. S. Hain, R. Jurowetzki, T. Buchmann, and P. Wolf, "A text-embedding-based approach to measuring patent-to-patent technological similarity," *Technol Forecast Soc Chang*, vol. 177, p. 121559, 2022. https://doi.org/10.1016/j.techfore.2022.121559
- [30] L. Helmers, F. Horn, F. Biegler, T. Oppermann, and K. R. Müller, "Automating the search for a patent's prior art with a full text similarity search," *PLoS ONE*, vol. 14, p. e0212103, 2019. https://doi.org/10.1371/journal.pone.0212103
- [31] K. Higuchi and K. Yanai, "Patent image retrieval using transformer-based deep metric learning," *World Patent Information*, vol. 74, p. 102217, 2023. https://doi.org/10.1016/j.wpi.2023.102217
- [32] Z. Hu, X. Zhou, and A. Lin, "Evaluation and identification of potential high-value patents in the field of integrated circuits using a multidimensional patent indicators pre-screening strategy and machine learning approaches," *Journal of Informetrics*, vol. 17, p. 101406, 2023. https://doi.org/10.1016/j.joi.2023.101406
- [33] D. Huang, C. Yan, Q. Li, and X. Peng, "From large language models to large multimodal models: A literature review," *Applied Sciences*, vol. 14, p. 5068, 2024. https://doi.org/10.3390/app14125068
- [34] G. Izacard *et al.*, "Atlas: Few-shot learning with retrieval augmented language models," *Journal of Machine Learning Research*, vol. 24, pp. 1-24, 2023. https://doi.org/10.5555/12345678
- [35] H. Jang, S. Kim, and B. Yoon, "An explainable ai (xai) model for text-based patent novelty analysis," *Expert Systems with Applications*, vol. 231, p. 120839, 2023. https://doi.org/10.1016/j.eswa.2023.120839
- [36] D. Jeon, J. M. Ahn, J. Kim, and C. Lee, "A doc2vec and local outlier factor approach to measuring the novelty of patents," *Technological Forecasting and Social Change*, vol. 174, p. 121294, 2022. https://doi.org/10.1016/j.techfore.2021.121294
- [37] H. Jiang, S. Fan, N. Zhang, and B. Zhu, "Deep learning for predicting patent application outcome: the fusion of text and network embeddings," *Journal of Informetrics*, vol. 17, p. 101252, 2023. https://doi.org/10.1016/j.joi.2023.101252
- [38] S. Jiang, J. Luo, G. Ruiz-Pava, J. Hu, and C. L. Magee, "Deriving design feature vectors for patent images using convolutional neural networks," *Journal of Mechanical Design*, vol. 143, p. 101301, 2021. https://doi.org/10.1115/1.4049214
- [39] S. Jiang, S. Sarica, B. Song, J. Hu, and J. Luo, "Patent data for engineering design: A critical review and future directions," *Journal of Computational and Information Science in Engineering*, vol. 22, p. 011001, 2022. https://doi.org/10.1115/1.4054802
- [40] J. Just, "Natural language processing for innovation search-reviewing an emerging non-human innovation intermediary," *Technovation*, vol. 129, p. 102883, 2024. https://doi.org/10.1016/j.technovation.2023.102883
- [41] E. Kamateri, M. Salampasis, and K. Diamantaras, "An ensemble framework for patent classification," *World Patent Information*, vol. 75, p. 102233, 2023. https://doi.org/10.1016/j.wpi.2023.102233
- [42] Y. Kang, Z. Cai, C. W. Tan, Q. Huang, and H. Liu, "Natural language processing (NLP) in management research: A literature review," *Journal of Management Analytics*, vol. 7, pp. 1-17, 2020. https://doi.org/10.1080/23270012.2020.1756939
- [43] S. Kim and B. Yoon, "Multi-document summarization for patent documents based on generative adversarial network," *Expert Systems with Applications*, vol. 207, p. 117983, 2022. https://doi.org/10.1016/j.eswa.2022.117983
- [44] R. Krestel, R. Chikkamath, C. Hewel, and J. Risch, "A survey on deep learning for patent analysis," *World Patent Information*, vol. 65, p. 102035, 2021. https://doi.org/10.1016/j.wpi.2021.102035
- [45] A. Larroyed, "Redefining patent translation: The influence of chatgpt and the urgency to align patent language regimes in Europe with progress in translation technology," *GRUR International*, vol. 72, 2023. https://doi.org/10.1093/grurint/ikad099
- [46] J. S. Lee, "Evaluating generative patent language models," *World Patent Information*, vol. 72, p. 102173, 2023. https://doi.org/10.1016/j.wpi.2023.102173
- [47] J. S. Lee and J. Hsiang, "Patent claim generation by fine-tuning openai gpt-2," *World Patent Information*, vol. 62, p. 101983, 2020. https://doi.org/10.1016/j.wpi.2020.101983
- [48] J. S. Lee and J. Hsiang, "Patent classification by fine-tuning bert language model," *World Patent Information*, vol. 61, p. 101965, 2020. https://doi.org/10.1016/j.wpi.2020.101965
- [49] C. Li, J. Chang, X. Wang, L. Zhao, and W. Mao, Formalization of natural language into PPTL specification via neural machine translation. In S. Liu, Z. Duan, & A. Liu (Eds.), Structured Object-Oriented Formal Language and Method SOFL+MSVL 2022. Cham: Springer, 2022.
- [50] W. Lin, J. Xiao, and Z. Cen, "Exploring bias in NLP models: Analyzing the impact of training data on fairness and equity," *Journal of Industrial Engineering and Applied Science*, vol. 2, no. 5, pp. 24-28, 2024.
- [51] W. Lin, W. Yu, and R. Xiao, "Measuring patent similarity based on text mining and image recognition," *Systems*, vol. 11, p. 294, 2023. https://doi.org/10.3390/systems11060294
- [52] W. Liu, S. Li, Y. Cao, and Y. Wang, "Multi-task learning based high-value patent and standard-essential patent identification model," *Information Processing & Management*, vol. 60, p. 103327, 2023. https://doi.org/10.1016/j.ipm.2023.103327
- [53] H. A. López, R. Strømsted, J. M. Niyodusenga, and M. Marquard, *Declarative process discovery: Linking process and textual views*. Cham: Springer, 2021.
- [54] K. Manas and A. Paschke, Semantic role assisted natural language rule formalization for intelligent vehicle. In A. Fensel, A. Ozaki, D. Roman, & A. Soylu (Eds.), Rules and Reasoning. Cham: Springer, 2023.
- [55] A. Mansouri and M. Mohammadpour, "Determining technology life cycle prediction based on patent bibliometric data," *International Journal of Information Science and Management*, vol. 21, no. 3, pp. 161-185, 2023.
- [56] B. Min et al., "Recent advances in natural language processing via large pre-trained language models: A survey," ACM Computing Surveys, vol. 56, no. 2, pp. 1-40, 2023. https://doi.org/10.1145/3605943
- [57] Y. A. Mohamed, A. H. H. M. Mohamed, A. Khanan, M. Bashir, M. A. E. Adiel, and M. A. Elsadig, "Navigating the ethical terrain of AI-generated text tools: A review," *IEEE Access*, vol. 12, pp. 197061-197120, 2024. https://doi.org/10.1109/ACCESS.2024.3521945
- [58] H. Mustroph, M. Barrientos, K. Winter, and S. Rinderle-Ma, Verifying resource compliance requirements from natural language text over event logs. Cham: Springer. https://doi.org/10.1007/978-3-031-41620-0_15, 2023.
- [59] A. Nayak, H. Timmapathini, V. Murali, K. Ponnalagu, V. G. Venkoparao, and A. Post, *Req2spec: Transforming software requirements into formal specifications using natural language processing*. Cham: Springer. https://doi.org/10.1007/978-3-030-98464-9_8, 2022.

- [60] T. Novotná and T. Libal, *An evaluation of methodologies for legal formalization*. Cham: Springer. https://doi.org/10.1007/978-3-031-15565-9_12, 2022.
- [61] D. W. Otter, J. R. Medina, and J. K. Kalita, "A survey of the usages of deep learning for natural language processing," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, 2020. https://doi.org/10.1109/TNNLS.2020.2979670
- [62] J. Pan, G. Chou, and D. Berenson, "Data-efficient learning of natural language to linear temporal logic translators for robot task specification," *IEEE International Conference on Robotics and Automation*, 2023. https://doi.org/10.1109/ICRA48891.2023.10161125
- [63] Y. Peng and J. Chi, "Unsupervised cross-media retrieval using domain adaptation with scene graph," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 6, pp. 1839-1850, 2019. https://doi.org/10.1109/TCSVT.2019.2953692
- [64] X. Pi, J. Shi, Y. Huang, and H. Wei, *Automated mining and checking of formal properties in natural language requirements*. Cham: Springer. https://doi.org/10.1007/978-3-030-29563-9_8, 2019.
- [65] Q. Plantec, P. Masson, and B. Weil, "Impact of knowledge search practices on the originality of inventions: A study in the oil & gas industry through dynamic patent analysis," *Technological Forecasting and Social Change*, vol. 168, p. 120782, 2021. https://doi.org/10.1016/j.techfore.2021.120782
- [66] G. Puccetti, V. Giordano, I. Spada, F. Chiarello, and G. Fantoni, "Technology identification from patent texts: A novel named entity recognition method," *Technological Forecasting and Social Change*, vol. 186, p. 122160, 2023. https://doi.org/10.1016/j.techfore.2022.122160
- [67] L. Quishpi, J. Carmona, and L. Padró, *Extracting decision models from textual descriptions of processes*. Cham: Springer, 2021.
- [68] C. Raffel *et al.*, "Exploring the limits of transfer learning with a unified text-to-text transformer," *Journal of Machine Learning Research*, vol. 21, no. 40, pp. 1-67, 2020.
- [69] J. Sànchez-Ferreres, A. Burattin, J. Carmona, M. Montali, and L. Padró, *Formal reasoning on natural language descriptions of processes*. Cham: Springer, 2019.
- [70] S. Sarica, J. Luo, and K. L. Wood, "Technet: technology semantic network based on patent data," *Expert Systems with Applications*, vol. 142, p. 112995, 2020. https://doi.org/10.1016/j.eswa.2019.112995
- [71] S. Sarica, B. Song, J. Luo, and K. L. Wood, "Idea generation with technology semantic network," AI EDAM, vol. 35, pp. 557-567, 2021. https://doi.org/10.1007/s10115-020-01492-4
- [72] V. J. Schmitt and N. M. Denter, "Modeling an indicator for statutory patent novelty," *World Patent Information*, vol. 78, p. 102283, 2024. https://doi.org/10.1016/j.wpi.2024.102283
- [73] V. J. Schmitt, L. Walter, and F. C. Schnittker, "Assessment of patentability by means of semantic patent analysis-a mathematicallogical approach," *World Patent Information*, vol. 73, p. 102182, 2023. https://doi.org/10.1016/j.wpi.2023.102182
- [74] D. Seal, U. K. Roy, and R. Basak, Sentence-level emotion detection from text based on semantic rules. In M. Tuba, S. Akashe, & A. Joshi (Eds.), Information and Communication Technology for Sustainable Development. Singapore: Springer, 2020.
- [75] S. Shibayama, D. Yin, and K. Matsumoto, "Measuring novelty in science with word embedding," *PLoS ONE*, vol. 16, p. e0254034, 2021. https://doi.org/10.1371/journal.pone.0254034
- [76] L. Siddharth, L. T. Blessing, K. L. Wood, and J. Luo, "Engineering knowledge graph from patent database," *Journal of Computing and Information Science in Engineering*, vol. 22, no. 5, p. 051003, 2022. https://doi.org/10.1115/1.4052293
- [77] L. Siddharth, G. Li, and J. Luo, "Enhancing patent retrieval using text and knowledge graph embeddings: A technical note," *Journal of Engineering Design*, vol. 33, no. 10, pp. 789-805, 2022. https://doi.org/10.1080/09544828.2022.2144714
- [78] L. Siddharth, N. Madhusudanan, and A. Chakrabarti, "Toward automatically assessing the novelty of engineering design solutions," *Journal of Computing and Information Science in Engineering*, vol. 20, no. 4, p. 041009, 2020. https://doi.org/10.1115/1.4044318
- [79] M. Siino, M. Falco, D. Croce, and P. Rosso, "Exploring LLMs applications in law: A literature review on current legal NLP approaches," *IEEE Access*, vol. 13, pp. 18253-18276, 2025. https://doi.org/10.1109/ACCESS.2025.3533217
- [80] J. Son *et al.*, "Ai for patents: A novel yet effective and efficient framework for patent analysis," *IEEE Access*, vol. 10, 2022. https://doi.org/10.1109/ACCESS.2022.3176877
- [81] C. M. Souza, M. R. Meireles, and P. E. Almeida, "A comparative study of abstractive and extractive summarization techniques to label subgroups on patent dataset," *Scientometrics*, vol. 126, no. 1, pp. 135-156, 2021. https://doi.org/10.1007/s11192-020-03732-x
- [82] X. Sun, N. Chen, and K. Ding, "Measuring latent combinational novelty of technology," *Expert Systems with Applications*, vol. 210, p. 118564, 2022. https://doi.org/10.1016/j.eswa.2022.118564
- [83] A. Trappey, C. V. Trappey, and A. Hsieh, "An intelligent patent recommender adopting machine learning approach for natural language processing: A case study for smart machinery technology mining," *Technological Forecasting and Social Change*, vol. 164, p. 120511, 2021. https://doi.org/10.1016/j.techfore.2020.120511
- [84] A. J. Trappey, C. V. Trappey, J. L. Wu, and J. W. Wang, "Intelligent compilation of patent summaries using machine learning and natural language processing techniques," *Advanced Engineering Informatics*, vol. 43, p. 101027, 2020. https://doi.org/10.1016/j.aei.2019.101027
- [85] K. Vowinckel and V. D. Hähnke, "Searchformer: Semantic patent embeddings by siamese transformers for prior art search," World Patent Information, vol. 73, p. 102192, 2023. https://doi.org/10.1016/j.wpi.2023.102192
- [86] H. Wang, H. Wu, Z. He, L. Huang, and K. W. Church, "Progress in machine translation," *Engineering*, vol. 18, pp. 143-153, 2022. https://doi.org/10.1016/j.eng.2021.03.023
- [87] J. Wang and Y. J. Chen, "A novelty detection patent mining approach for analyzing technological opportunities," *Advanced Engineering Informatics*, vol. 42, p. 100941, 2019. https://doi.org/10.1016/j.aei.2019.100941
- [88] M. Wang, H. Sakaji, H. Higashitani, M. Iwadare, and K. Izumi, "Discovering new applications: Cross-domain exploration of patent documents using causal extraction and similarity analysis," *World Patent Information*, vol. 75, p. 102238, 2023. https://doi.org/10.1016/j.wpi.2023.102238
- [89] X. Wang, G. Li, C. Li, L. Zhao, and X. Shu, *Automatic generation of specification from natural language based on temporal logic*. Cham: Springer, 2020.
- [90] Y. Wang, *Towards automated contract analysis: Applying language models to risk identification in the context of public-private partnerships.* United States -- Maryland: University of Maryland, College Park, 2024.

- [91] T. Wei, D. Feng, S. Song, and C. Zhang, "An extraction and novelty evaluation framework for technology knowledge elements of patents," *Scientometrics*, vol. 129, pp. 1-26, 2024. https://doi.org/10.1007/s11192-024-04990-9
- [92] Z. Wen and Y. Peng, "Multi-level knowledge injecting for visual commonsense reasoning," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, 2020. https://doi.org/10.1109/TCSVT.2020.2991866
- [93] T. Xia, "A constant time complexity spam detection algorithm for boosting throughput on rule-based filtering Systems," *IEEE Access*, vol. 8, pp. 82653-82661, 2020. https://doi.org/10.1109/ACCESS.2020.2991328
- [94] X. Yan, Y. Ye, Y. Mao, and H. Yu, "Shared-private information bottleneck method for cross-modal clustering," *IEEE Access*, vol. 7, pp. 36045-36056, 2019. https://doi.org/10.1109/ACCESS.2019.2904554