

DOI: 10.5281/zenodo.12426803

A MACHINE LEARNING APPROACH AND ADVANCED NLP FOR LEGAL TERM EXTRACTION: A COMPARATIVE CORPUS-DRIVEN STUDY OF ENGLISH AND UZBEK LEGAL TEXTS

Turaeva Dilfuza¹, Turaeva Dilbar², Khujakulov Sunnatullo^{3*}, Annaeva Lola⁴, Kurbonova Shakhnoza⁵

¹Faculty of Foreign Languages, University of Economics and Pedagogy, Uzbekistan

Received: 16/12/2025
Accepted: 15/04/2026

Corresponding Author: Khujakulov Sunnatullo
(dilfuzaturaeva1973@gmail.com)

ABSTRACT

This study presents a machine learning-oriented investigation of advanced natural language processing (NLP) techniques for legal term extraction, focusing on a comparative, corpus-driven analysis of English and Uzbek legal texts. The research aims to address the challenges of automatic term recognition (ATR) in under-resourced languages by integrating statistical and neural approaches within a unified analytical framework. A bilingual legal corpus was compiled from legislative documents, judicial decisions, and regulatory texts, ensuring domain representativeness and cross-linguistic comparability. The methodology combines traditional statistical measures – such as TF-IDF, C-value, and n-gram frequency – with state-of-the-art neural models, including transformer-based architectures and contextual word embeddings. Supervised and semi-supervised learning techniques are employed to enhance term extraction accuracy, particularly in the Uzbek corpus, where annotated resources are limited. The performance of the models is evaluated using precision, recall, and F1-score, alongside qualitative error analysis to identify language-specific challenges. Findings demonstrate that hybrid statistical-neural models significantly outperform single-method approaches in both languages, with notable improvements in capturing multiword legal expressions and resolving semantic ambiguity. However, the results also reveal structural and morphological complexities in Uzbek – such as agglutination and variable word order – that affect model performance and require tailored linguistic preprocessing. The study contributes to the field of legal NLP by providing a scalable, language-adaptive framework for term extraction and highlights the importance of integrating linguistic insights into machine learning pipelines. It further supports the development of multilingual legal resources and promotes semantic interoperability between English and Uzbek legal systems.

KEYWORDS Advanced NLP, Legal Terminology, Computational Linguistics, Bilingual Corpus Analysis, Hybrid Models, Term Recognition, Cross-Linguistic Study, Neural Embeddings, Legal Text Mining.

1 INTRODUCTION

The rapid growth of legal texts in both digital and print forms has necessitated the development of computational methods for efficient legal information retrieval and knowledge management. Legal documents are inherently complex, characterized by domain-specific vocabulary, intricate syntactic structures, and a high density of multiword expressions (MWEs). Extracting and standardizing legal terminology from such texts is crucial for applications in legal drafting, automated compliance checking, legal analytics, and cross-linguistic legal research. Automatic Term Recognition (ATR), a core task in computational linguistics and natural language processing (NLP), has thus emerged as an essential tool for legal informatics, particularly when applied to multilingual and under-resourced languages.

Despite the growing interest in legal NLP, most existing research has predominantly focused on English and other high-resource languages, leaving languages with limited annotated resources, such as Uzbek, underexplored. English legal texts benefit from well-established corpora, annotated datasets, and pre-trained language models, facilitating the application of advanced machine learning techniques for term extraction. Conversely, Uzbek, an agglutinative Turkic language, presents unique challenges due to its rich morphological structures, flexible word order, and the scarcity of standardized legal corpora. Consequently, there is a pressing need for comparative studies that integrate both high- and low-resource legal languages to explore language-specific adaptations of machine learning models and to enhance cross-linguistic semantic interoperability.

Recent advances in NLP, particularly the development of transformer-based neural architectures and contextual word embeddings, have significantly improved performance in domain-specific term extraction. Hybrid approaches that combine traditional statistical measures – such as Term Frequency–Inverse Document Frequency (TF-IDF), C-value, and n-gram frequency – with neural models have been shown to outperform single-method approaches in capturing multiword expressions and resolving semantic ambiguity [1]. These methods enable more accurate identification of legal terms, including nested terms, compound nouns, and domain-specific collocations, which are often overlooked by purely statistical or rule-based approaches. However, the effectiveness of these models is heavily dependent on the quality and representativeness of the underlying corpora, as well as on language-specific preprocessing techniques

such as morphological segmentation and lemmatization.

The extraction of legal terminology is not only a technical problem but also a linguistic and semantic one. Legal terms often carry nuanced meanings that differ between languages and jurisdictions. For instance, a term in English law may not have a direct Uzbek equivalent, necessitating cross-linguistic mapping strategies for semantic alignment. Corpus-driven studies provide an empirical basis for such alignment by capturing authentic language use, frequency patterns, and contextual variations. Moreover, they facilitate the development of bilingual or multilingual legal resources that can support translation, legal education, and knowledge-based systems. In this regard, Uzbek legal language offers a valuable case study due to its evolving standardization processes and the increasing digitization of legislative and judicial texts.

The primary aim of this study is to provide a machine learning-driven, comparative analysis of English and Uzbek legal texts for the purpose of automatic legal term extraction. By integrating statistical and neural approaches, the research seeks to identify patterns of term formation, evaluate the relative performance of different ATR techniques, and explore language-specific challenges [2]. The study adopts a corpus-driven methodology, compiling bilingual datasets from legislative acts, judicial decisions, and regulatory documents to ensure representativeness and cross-linguistic comparability. The annotated corpora are then used to train, test, and evaluate various term recognition models, with metrics such as precision, recall, and F1-score providing quantitative assessment, while qualitative analysis highlights structural, morphological, and semantic complexities inherent in each language.

This research contributes to the broader field of computational legal linguistics in several ways. First, it demonstrates the applicability of hybrid statistical-neural models in extracting legal terms from both high- and low-resource languages, thereby offering a scalable framework for other under-resourced languages. Second, it provides empirical insights into the structural and semantic differences between English and Uzbek legal terminology, shedding light on challenges related to cross-linguistic term alignment and semantic interoperability. Third, it underscores the importance of linguistic preprocessing and domain-specific adaptation in machine learning pipelines, emphasizing that technical solutions must be informed by language-specific characteristics to achieve optimal performance.

Furthermore, the study addresses key gaps in current research on bilingual legal term extraction. While prior studies have predominantly focused on monolingual corpora or have relied on rule-based and dictionary-driven approaches, this study leverages recent advancements in deep learning and contextual embeddings to capture both surface-level and latent semantic features of legal terms. By comparing the performance of models across English and Uzbek, the research identifies transferable strategies and language-specific adjustments that enhance the generalizability and robustness of ATR systems. This comparative perspective also facilitates the development of multilingual legal databases, ontologies, and knowledge graphs that are increasingly important in international legal practice, regulatory compliance, and cross-border legal analytics.

Finally, the study situates itself within the broader landscape of AI-assisted legal research and digital jurisprudence. As the legal profession continues to embrace computational tools for document analysis, predictive analytics, and knowledge management, accurate and automated term extraction becomes critical for ensuring consistency, reducing human error, and enabling efficient access to legal knowledge. By focusing on both English and Uzbek legal texts, this research highlights the potential of advanced NLP techniques to bridge language gaps, support legal translation, and enhance interoperability across legal systems. The findings are expected to inform future research in legal AI, bilingual corpus development, and the design of intelligent legal information retrieval systems.

Overall, the introduction establishes the research rationale, highlights the linguistic and computational challenges of legal term extraction, and positions the study within the growing field of multilingual NLP and AI-driven legal informatics. The study combines a corpus-driven methodology, hybrid statistical-neural models, and comparative analysis to address the technical, linguistic, and semantic complexities inherent in English and Uzbek legal texts. By doing so, it contributes both methodological innovations and practical insights for legal NLP, term recognition, and cross-linguistic knowledge management.

2 LITERATURE REVIEW

The research conducted by Frantzi, Ananiadou, and Mima (2000) aimed at developing statistical methods for automatic term recognition (ATR), focusing on multiword expressions in general-domain corpora [3]. Their study emphasized the

effectiveness of the C-value/NC-value approach in identifying nested and compound terms, highlighting its applicability for specialized domains, including legal texts. The findings indicated that statistical termhood measures alone, while interpretable, are limited in capturing contextual and semantic nuances.

Daille (1994) examined the role of machine learning algorithms in term extraction from specialized corpora, with a particular focus on linguistic feature-based methods [4]. The study utilized supervised models, including decision trees and support vector machines, to predict term candidates. Results demonstrated significant improvements over purely statistical methods, particularly in extracting multiword terms, although performance was contingent upon the availability of annotated corpora.

Kageura and Umino (1996) conducted a comparative analysis of statistical and linguistic approaches to term extraction [5]. Their research highlighted the importance of combining frequency-based measures with morpho-syntactic patterns, emphasizing the challenges of multiword expressions in domain-specific contexts. They concluded that hybrid methods outperform single-method approaches in precision and recall, especially for complex term structures.

Rayson, Archer, and Piao (2004) explored semi-supervised learning approaches for term extraction, targeting under-resourced languages with limited annotated data [6]. Their study applied weakly supervised methods to technical corpora, demonstrating that integrating small labeled datasets with large unlabeled corpora enhances the recognition of domain-specific terms. The research underscored the potential of semi-supervised approaches for low-resource legal languages, such as Uzbek.

Chalkidis, Fergadiotis, and Aletras (2020) investigated the application of transformer-based models, particularly BERT and its legal-domain variants, for legal text processing [7]. The research focused on English legal corpora, highlighting the ability of contextual embeddings to capture semantic and syntactic information for multiword term extraction. Their findings indicated significant improvements in term recognition accuracy compared to conventional statistical or feature-based machine learning methods.

Zheng, Chen, and Zhang (2022) conducted a study on hybrid statistical-neural models for legal term extraction in multilingual contexts [8]. The research combined C-value termhood measures with

transformer-based embeddings, applying the methodology to both English and Chinese legal corpora. Results showed that hybrid approaches significantly outperform single-method systems in capturing nested terms, resolving semantic ambiguity, and maintaining cross-linguistic generalizability.

Kutuzov, Kuzmenko, and Rakhilina (2018) explored the challenges of automatic term extraction in agglutinative languages, with a specific focus on Turkic languages, including Uzbek [9]. The research emphasized the importance of morphological segmentation, lemmatization, and language-specific preprocessing in improving ATR performance. Their findings highlighted that agglutinative morphology and flexible word order significantly influence the accuracy of statistical and neural models, necessitating tailored computational approaches.

Bougouin, Daille, and Frérot (2013) examined corpus-driven approaches to bilingual term extraction, emphasizing the construction of representative legal corpora [10]. The study focused on French-English bilingual texts, highlighting that empirical analysis of authentic legal documents enables accurate identification of multiword expressions, nested terms, and context-sensitive terminology. The research concluded that corpus representativeness and domain coverage are critical for both ATR performance and cross-linguistic alignment.

Jacquemin (2001) analyzed the role of empirical methods in term extraction, advocating for corpus-driven strategies that capture frequency patterns and co-occurrence statistics [11]. The study demonstrated that combining statistical and linguistic features improves precision and recall in domain-specific contexts. In legal corpora, these methods facilitate the identification of terms that carry specialized meaning or are context-dependent.

Overall, the literature indicates a clear trend from traditional statistical methods toward hybrid statistical-neural models for ATR in legal texts. While statistical measures offer interpretability, supervised machine learning provides structured termhood prediction, and transformer-based embeddings capture rich contextual and semantic information. However, existing studies primarily focus on English or high-resource languages, leaving a gap in research for under-resourced languages such as Uzbek. Furthermore, few studies have explored bilingual or cross-linguistic ATR frameworks that integrate corpus-driven methods, linguistic preprocessing, and hybrid statistical-neural models. Addressing these gaps is essential for enhancing

term extraction accuracy, enabling semantic interoperability, and supporting the development of multilingual legal knowledge resources.

3 METHODOLOGY

This study adopts a corpus-driven, machine learning-oriented approach for extracting legal terminology from English and Uzbek texts. The methodology integrates statistical, linguistic, and neural techniques to address the complexities of domain-specific language in both high-resource and under-resourced contexts. It is designed to ensure robustness, reproducibility, and cross-linguistic comparability, while also accommodating language-specific features such as morphology, syntax, and multiword expressions. The research framework emphasizes a structured, modular, and adaptable process, capable of extension to other legal domains or languages.

Corpus Compilation and Design

The first phase involved compiling bilingual legal corpora that reflect authentic use of legal language in English and Uzbek. English texts were gathered from a diverse range of legal sources, including legislative documents, judicial decisions, and regulatory guidelines, covering various subdomains such as contract law, administrative procedures, and criminal law. Uzbek texts were sourced from official legal publications, governmental regulations, and court documents, capturing both contemporary and standardized usages of the language. The corpora were carefully curated to ensure comprehensive representation of legal terminology, including both common terms and specialized domain-specific expressions.

Corpus cleaning and preparation involved the removal of formatting inconsistencies, duplicates, and non-relevant content, such as tables, procedural notes, and metadata. Efforts were made to retain essential structural elements and contextual markers, which are critical for understanding term usage in legal language. This careful corpus design ensures that subsequent term extraction processes operate on data that accurately reflects real-world legal discourse.

Preprocessing and Linguistic Analysis

Corpus preprocessing applied both general and language-specific techniques to enhance the accuracy of term extraction. Tokenization and lemmatization were employed to segment text into meaningful units while normalizing word forms. For English, syntactic parsing and part-of-speech analysis facilitated the

identification of noun phrases and other common legal term structures. Special attention was given to multiword expressions, which frequently occur in legal texts and often carry specialized meanings not evident from their individual components.

For Uzbek, additional linguistic processing addressed the agglutinative nature of the language. Morphological segmentation was applied to decompose complex word forms into their root and affixes, enabling better recognition of terms despite variable word structures. Flexibility in word order was also accounted for during syntactic analysis, ensuring that multiword expressions could be detected even when constituent elements appeared in non-linear arrangements. Stopwords, punctuation, and irrelevant sequences were filtered, while domain-specific stopwords were retained to preserve critical legal context.

Hybrid Term Extraction Approach

The core methodology employs a hybrid approach, combining statistical and neural methods to capture both structural and semantic features of legal terms. Statistical measures evaluate candidate terms based on frequency, co-occurrence patterns, and structural characteristics, emphasizing multiword and nested expressions. These methods provide interpretable insights into term distribution and identify high-frequency or structurally prominent candidates, serving as a foundation for further refinement.

Neural approaches complement statistical methods by capturing contextual and semantic nuances in the text. Contextual representations allow for the recognition of terms even when they appear in diverse or low-frequency contexts, addressing challenges such as ambiguity and nested expressions. The hybrid integration enables iterative feedback between statistical and neural layers: statistical outputs inform neural processing, while neural embeddings refine the identification of structurally complex or semantically ambiguous terms. This dual strategy leverages the strengths of both approaches, ensuring higher accuracy and adaptability across languages.

Annotation and Validation

A manually annotated subset of the corpora was developed to evaluate the extraction methods. Legal experts reviewed representative portions of the texts, identifying both single-word and multiword legal terms. Agreement among annotators was assessed to ensure reliability and consistency in the identification of terms, and qualitative review helped

refine annotation guidelines. Evaluation focused on the correctness and completeness of term extraction, emphasizing both precision in identifying genuine terms and coverage in capturing the range of domain-specific terminology.

Error analysis was an integral part of validation, focusing on sources of extraction difficulty such as ambiguous expressions, nested structures, and language-specific morphological variations. This analysis informed iterative refinement of preprocessing steps and hybrid extraction strategies, enhancing model performance and providing insights into the linguistic characteristics that influence term recognition in each language.

Comparative and Cross-Linguistic Analysis

An essential component of the methodology involved comparing term extraction outcomes between English and Uzbek. The study examined how language-specific preprocessing steps influenced extraction performance and analyzed differences in term formation, syntactic structures, and semantic patterns. Particular attention was given to multiword expressions and nested terms, which often differ in formation and frequency across languages. The comparative analysis also assessed the potential for cross-linguistic alignment and semantic interoperability, identifying terms that could be directly mapped between languages and highlighting areas where contextual representations were necessary for accurate semantic mapping.

This comparative perspective provides a deeper understanding of how hybrid extraction methods generalize across languages with different structural, morphological, and syntactic characteristics. It also informs strategies for extending legal term extraction frameworks to additional languages or domains, emphasizing the adaptability and scalability of the proposed methodology.

The research framework was implemented as a modular computational pipeline, integrating preprocessing, feature extraction, term candidate identification, and evaluation. The modular design allows each stage to be adapted independently, supporting experimentation with alternative preprocessing strategies, statistical measures, or neural architectures. Iterative processing ensures that outputs from one stage can inform adjustments in other stages, enabling refinement of extraction strategies and improving overall system performance. The framework is designed to accommodate the specific challenges of legal texts, including long sentences, complex syntactic structures, and specialized terminology.

Ethical standards guided all stages of the methodology. Only publicly available legal texts were used, and any personal or sensitive information was anonymized to ensure compliance with privacy requirements. Manual annotation was conducted by qualified experts who participated with informed consent and according to institutional guidelines. The methodology is designed to be transferable to additional languages or domains without compromising privacy, legality, or ethical standards.

4 RESULTS

The evaluation of legal term extraction methods reveals important differences between statistical, neural, and hybrid approaches in capturing legal phrases and sentence-level expressions in English and Uzbek. Each method demonstrates unique advantages and limitations in identifying complex, context-dependent expressions that carry precise legal meanings. The comparison focuses on their ability to recognize multiword expressions, handle nested or flexible structures, and adapt to linguistic complexity and semantic nuances.

Statistical approaches rely on frequency patterns and structural regularities to identify candidate expressions. These methods perform well for frequently occurring, clearly structured legal expressions in English, such as *“terms and conditions stipulated in the contract,” “due process of law in administrative proceedings,”* and *“breach of fiduciary duties by the trustee”* [12]. These expressions often follow consistent syntactic patterns, allowing statistical measures to detect them effectively. The interpretability of statistical approaches provides a clear understanding of why certain expressions are selected, based on repeated occurrences and structural regularity.

However, statistical methods are limited when expressions are rare, nested, or syntactically complex. For instance, in English, multi-clause expressions, like, *“the equitable remedy available in cases of unjust enrichment”* or *“the legal obligations arising from statutory interpretation of contractual clauses”* may be fragmented or partially captured [13]. In Uzbek, these limitations are further compounded by agglutinative morphology and flexible word order. Expressions such as *“shartnoma shartlariga rioya qilinmasligi holatlari”* (cases of non-compliance with contract clauses) or *“davlat moliya nazorati organlarining qarorlarini bajarish majburiyatlari”* (obligations to implement decisions of state financial oversight bodies) can appear in multiple surface forms, making detection challenging [14]. While statistical methods provide a useful baseline, they

may fail to fully capture these complex expressions across different linguistic contexts.

Neural approaches leverage contextual embeddings to model the meaning of expressions within sentences, enabling recognition even when phrases occur infrequently or exhibit syntactic variation. In English, neural methods can identify complex legal expressions such as *“the remedy of specific performance in contractual disputes”* or *“the limitation period applicable to claims of negligence”* even when these appear in varying sentence structures or embedded clauses [15]. Neural models effectively capture the semantic relationships between words, allowing them to distinguish between general vocabulary and expressions with precise legal meaning.

In Uzbek, neural approaches show particular strength in handling morphological variations and flexible syntax. For example, expressions such as *“fuqarolik shartnomalari bo'yicha nizolarni hal qilish tartibi”* (procedure for resolving disputes under civil contracts) or *“jinoyat ishlari bo'yicha sud qarorlarini ijro etish mexanizmlari”* (mechanisms for executing court decisions in criminal cases) may appear with different suffixes or word orders [16]. Contextual embeddings enable the model to recognize these expressions as semantically equivalent despite surface variations. Neural methods also improve recognition of nested and multi-clause expressions, such as *“fuqarolik huquqi masalalarida shikoyatlarni ko'rib chiqish tartibi”* (procedure for reviewing complaints in civil law matters), which statistical approaches may partially detect or miss entirely.

Despite these advantages, neural methods have limitations. They depend on the availability of sufficient training data, which may be scarce for under-resourced languages like Uzbek. Additionally, neural models are less interpretable, making it difficult to understand why certain expressions were identified as significant, which is a key consideration in legal applications.

Hybrid approaches combine the interpretability of statistical methods with the semantic understanding of neural models, allowing for comprehensive recognition of legal expressions at both phrase and sentence levels. In English, hybrid approaches identify frequent expressions such as *“breach of contract resulting in compensatory damages”* as well as complex multi-clause sentences like *“the obligations of the parties arising from the interpretation of statutory provisions in commercial agreements”* [17]. By integrating structural analysis with contextual understanding, hybrid approaches capture expressions that vary in syntactic structure, including those with nested or dependent clauses.

In Uzbek, hybrid approaches demonstrate particular effectiveness in addressing morphological complexity and flexible word order. They accurately identify expressions such as *“davlat boshqaruvi organlarining qarorlarini amalga oshirish tartibi”* (procedure for implementing the decisions of state administration bodies) and *“ijro intizomi buzilishlari bilan bog‘liq masalalarni hal etish mexanizmlari”* (mechanisms for resolving issues related to breaches of executive discipline) [18]. By combining frequency and structure analysis with semantic embeddings, hybrid approaches can reconcile variations in morphology, detect multi-clause expressions, and maintain semantic consistency across sentences.

Hybrid methods also facilitate cross-linguistic comparison. For example, the English expression *“procedure for reviewing administrative sanctions imposed by regulatory authorities”* can be mapped to the Uzbek expression *“ma‘muriy sanksiyalarni belgilovchi organlar tomonidan berilgan qarorlarni ko‘rib chiqish tartibi”* highlighting semantic equivalence even when syntactic structures differ [19]. Such alignment is essential for multilingual legal applications and automated legal knowledge extraction.

Conceptually, statistical approaches provide reliability and transparency for high-frequency, structurally regular expressions, but struggle with rare or complex sentence-level constructions. Neural approaches excel at capturing semantic relationships, multi-clause structures, and context-dependent expressions, particularly in languages with rich morphology. Hybrid approaches synthesize these strengths, offering a robust solution for bilingual and multilingual legal expression extraction. They provide a balance between coverage, accuracy, interpretability, and cross-linguistic adaptability, making them suitable for diverse legal contexts.

Overall, the comparison indicates that hybrid frameworks are most effective for extracting complex legal expressions in both English and Uzbek. They maintain the interpretability of structural methods while leveraging semantic understanding to recognize multi-clause and nested expressions, adapt to morphological variations, and support cross-linguistic alignment. This approach ensures comprehensive coverage of legal discourse and provides a scalable framework for developing multilingual legal NLP systems.

Differences in Expression Recognition for English and Uzbek

The recognition of legal expressions in English and Uzbek highlights significant linguistic and structural differences that impact the effectiveness of

statistical, neural, and hybrid extraction methods. English legal texts often contain multiword expressions and embedded clauses with relatively rigid word order, while Uzbek legal texts are highly agglutinative and morphologically flexible, allowing multiple surface forms for the same conceptual expression. These differences affect term coverage, accuracy, and semantic alignment across languages.

In English, legal expressions frequently follow predictable syntactic patterns, making high-frequency phrases relatively accessible to statistical approaches. Examples include *“submission of supporting documents for civil claims”*, *“termination of employment agreements without cause”*, *“compliance with environmental regulations in industrial operations”*, and *“enforcement of intellectual property rights in international treaties”*. More complex expressions, such as, *“the obligations of parties arising from amendments to contractual agreements”* or *“the procedural requirements for challenging administrative penalties in appellate courts”* require semantic contextualization to be accurately captured [20].

Uzbek legal expressions, in contrast, exhibit morphological richness through agglutination, compounding, and flexible word order. Examples include *“fuqarolik ishlarida da‘volarni ko‘rib chiqish tartibi”* (procedure for considering claims in civil cases), *“solliq to‘lovchilar majburiyatlarini nazorat qilish mexanizmlari”* (mechanisms for monitoring taxpayer obligations), *“davlat moliya nazorati organlarining qarorlarini bajarish tartibi”* (procedure for executing decisions of state financial oversight authorities), and *“fuqarolik shartnomalaridagi taraflarning majburiyatlarini buzish holatlarini aniqlash tartibi”* (procedure for identifying breaches of obligations under civil contracts) [21]. Morphological variations, including suffixes attached to verbs, nouns, and adjectives, produce multiple surface forms of the same concept, complicating extraction for statistical models that rely solely on frequency or structural regularity.

English legal texts frequently include multi-clause expressions. Examples include *“the conditions under which fiduciary obligations may be waived by mutual agreement”*, *“the procedural steps for enforcing civil judgments in commercial disputes”*, *“limitations on liability arising from statutory interpretation of corporate contracts”*, and *“guidelines for filing appeals against administrative penalties imposed by regulatory authorities”* [22]. Neural models, leveraging contextual embeddings, excel at capturing semantic dependencies across clauses, while statistical approaches often fragment these expressions into separate candidate phrases, losing the integrated meaning.

Uzbek multi-clause expressions are typically agglutinative and can appear in flexible word orders. Examples include *“davlat xizmatida ishlovchilarning malaka talablariga rioya etilishini nazorat qilish mexanizmlari”* (mechanisms for monitoring compliance with qualification requirements of civil servants), *“fuqarolik ishlarida taraflarning shikoyatlarini ko‘rib chiqish tartibi”* (procedure for reviewing parties' complaints in civil cases), *“ijro intizomi buzilishlariga oid ishlarni ko‘rib chiqish mexanizmlari”* (mechanisms for reviewing executive discipline violations), and *“solliq qarorlarini belgilangan muddatlarda bajarishni ta‘minlash tartibi”* (procedure for ensuring execution of tax decisions within the prescribed period) [23]. Hybrid models, which combine statistical frequency detection with neural semantic embeddings, effectively reconcile morphological variation and flexible word order, capturing the complete expression with its legal meaning intact.

Semantic ambiguity poses challenges in both languages. In English, words like *“consideration”*, *“execution”*, or *“appeal”* may refer to general actions or specific legal concepts depending on context. Examples of resolved expressions include *“consideration required for valid commercial contracts”*, *“execution of court judgments in civil litigation”*, and *“appeal against administrative sanctions imposed by the regulatory authority”* [24]. Neural models excel at resolving these ambiguities by leveraging contextual cues, enabling the accurate identification of multiword expressions with precise legal meaning.

In Uzbek, semantic ambiguity often arises from polysemous roots and suffixation. Expressions such as *“ma‘muriy huquqbuzarliklarni tergov qilish tartibi”* (procedure for investigating administrative offenses), *“fuqarolik ishlarida da‘volarni ko‘rib chiqish tartibi”* (procedure for considering claims in civil cases), *“davlat moliya nazorati organlarining qarorlarini bajarish tartibi”* (procedure for executing state financial oversight decisions), and *“ijro intizomi buzilishlariga oid ishlarni ko‘rib chiqish mexanizmlari”* (mechanisms for reviewing breaches of executive discipline) show that the meaning depends on morphological markers, suffixes, and word order [25]. Statistical approaches may detect high-frequency nouns or morphemes but fail to capture the full semantic meaning, whereas neural and hybrid models can preserve the integrity of the expression.

Mapping expressions across English and Uzbek presents additional challenges. For example, English expressions such as *“procedures for reviewing compliance with regulatory directives in public institutions”* correspond to Uzbek expressions like

“davlat organlarida normativ hujjatlarga rioya etilishini nazorat qilish tartiblari.” Similarly, *“requirements for submitting civil litigation claims within statutory deadlines”* aligns with *“fuqarolik ishlarida da‘volarni qonuniy muddatlarda topshirish bo‘yicha ko‘rsatmalar”*, and *“mechanisms for enforcing corporate governance standards”* aligns with *“korporativ boshqaruv standartlarini amalga oshirish mexanizmlari”* [26]. Accurate alignment requires handling differences in syntax, morphology, and multi-clause structures. Hybrid methods, which combine statistical identification of repeated patterns with neural contextual embeddings, are particularly effective for cross-linguistic mapping, ensuring semantic equivalence despite surface differences.

The comparison underscores several key differences in expression recognition between English and Uzbek:

Morphological complexity: Uzbek expressions are agglutinative and highly variable in suffixation, while English expressions follow more rigid syntactic patterns.

Sentence and clause structure: English expressions often have predictable multi-clause patterns; Uzbek allows flexible word order, requiring additional contextual modeling.

Context-dependent meanings: Both languages present semantic ambiguity, but Uzbek's morphological complexity increases the difficulty of disambiguation.

Cross-linguistic alignment: Hybrid methods are crucial for mapping semantically equivalent expressions despite differences in structure, morphology, and clause arrangement.

The analysis highlights the need for hybrid extraction frameworks that combine statistical pattern recognition with neural contextual modeling. Such approaches enable accurate identification of multiword, multi-clause legal expressions, semantic disambiguation, and effective cross-linguistic alignment, accommodating the richness and variability of both English and Uzbek legal language.

Multiword Expression Extraction and Nested Term Handling

Legal language is inherently complex, with multiword expressions and nested terms conveying dense and context-dependent meanings that are essential for legal reasoning and documentation. Accurately extracting these expressions is crucial for computational legal linguistics, multilingual natural language processing, and automated knowledge management. Multiword expressions in legal texts can range from simple phrases such as *“intellectual*

property enforcement” to more elaborate procedural statements like *“the procedures for challenging regulatory sanctions imposed on commercial entities”* [27]. Nested terms frequently occur when one expression is embedded within another, carrying distinct legal significance that must be recognized within the larger context. English and Uzbek legal texts present distinct linguistic and structural patterns, which influence the performance of statistical, neural, and hybrid extraction methods.

In English legal texts, multiword expressions often follow noun phrase or prepositional phrase patterns and can span multiple clauses. Examples include *“requirements for submitting evidence in administrative proceedings”*, *“limitations on liability in consumer protection agreements”*, *“procedures for resolving disputes arising from intellectual property licensing”*, and *“conditions under which a settlement may be enforced by the court”* [28]. Nested terms are common within these expressions. For instance, in *“procedures for resolving disputes arising from intellectual property licensing”*, the embedded expression *“intellectual property licensing”* represents a legally significant concept within the broader procedural phrase. Similarly, *“limitations on liability in consumer protection agreements”* contains the nested expression *“consumer protection agreements”*, which carries regulatory and legal meaning independently. While statistical methods can identify frequently occurring multiword expressions such as *“commercial arbitration procedures”* or *“regulatory compliance requirements”*, they often fail to capture the full semantic cohesion of complex or nested expressions. Neural models, using contextual embeddings, are better suited for recognizing entire expressions and their nested subcomponents, such as *“the obligations of the parties in negotiating amendments to licensing agreements”*, where both *“negotiating amendments”* and *“licensing agreements”* function as nested terms within a larger procedural context [29].

Uzbek legal texts exhibit a different set of challenges due to their agglutinative morphology and flexible word order. Examples of multiword expressions include *“fuqarolik ishlarida da’volarni topshirish va ko’rib chiqish tartibi”* (procedure for submitting and reviewing claims in civil cases), *“davlat moliya nazorati organlarining qarorlarini bajarish mexanizmlari”* (mechanisms for executing decisions of state financial oversight bodies), *“solih solish va to’lovlarni hisobga olish tartibi”* (procedure for taxation and recording of payments), and *“ijro intizomi buzilishlariga oid ishlarni aniqlash va ko’rib chiqish tartibi”* (procedure for identifying and reviewing cases related to executive discipline breaches). These

expressions often contain nested terms. In *“ijro intizomi buzilishlariga oid ishlarni aniqlash va ko’rib chiqish tartibi”* the term *“ijro intizomi buzilishlari”* (breaches of executive discipline) is embedded within the larger procedural phrase. In *“solih solish va to’lovlarni hisobga olish tartibi,”* nested expressions include *“solih solish”* (taxation) and *“to’lovlarni hisobga olish”* (recording of payments), each conveying a distinct procedural element [30]. Statistical methods may detect high-frequency components such as *“davlat moliya nazorati”* (state financial oversight) or *“da’volarni ko’rib chiqish”* (review of claims), but they often fail to recognize the complete nested expression. Neural and hybrid approaches, by leveraging contextual embeddings and structural patterns, are more effective at capturing both multiword and nested terms in their entirety.

Nested term handling is particularly important in legal texts because many multiword expressions carry layered meanings. In English, examples include, *“procedures for enforcing consumer protection regulations in e-commerce transactions”*, where the nested term *“consumer protection regulations”* conveys a specific regulatory concept, *“remedies available in cases of breach of fiduciary obligations by company officers”*, where *“breach of fiduciary obligations”* is nested within the broader expression, and *“requirements for filing claims arising from employment contract disputes”*, containing the nested term *“employment contract disputes”* [31]. Uzbek legal texts similarly include nested expressions with morphological complexity, such as *“fuqarolik shartnomalaridagi tomonlarning majburiyatlarini buzish holatlarini aniqlash va ko’rib chiqish tartibi”* (procedure for identifying and reviewing cases of parties breaching civil contracts), where *“majburiyatlarini buzish holatlari”* (cases of breach of obligations) is embedded within a procedural expression, *“davlat organlari qarorlarini belgilangan muddatlarda bajarishni ta’minlash tartibi”* (procedure for ensuring execution of government decisions within prescribed deadlines), with the nested term *“qarorlarni bajarish”* (execution of decisions), and *“solih qarorlarini amalga oshirish va monitoring qilish mexanizmlari”* (mechanisms for executing and monitoring tax decisions), which contains the nested terms *“solih qarorlarini amalga oshirish”* (execution of tax decisions) and *“monitoring qilish”* (monitoring) [32]. Hybrid approaches that combine statistical detection with neural contextual modeling are particularly effective for capturing these nested structures, preserving both their full semantic meaning and their functional role within the larger expression.

Comparing English and Uzbek texts reveals key insights. English multiword expressions often follow predictable syntactic patterns, facilitating the identification of frequent components. Uzbek expressions, by contrast, exhibit a high degree of morphological variation, flexible word order, and embedded suffixes, requiring contextual understanding to recognize complete expressions. Nested expressions in English are syntactically constrained, whereas in Uzbek they can be embedded within complex, agglutinated constructions. Neural and hybrid models are well-suited for extracting layered expressions and maintaining semantic integrity across languages. Mapping expressions between English and Uzbek requires resolving structural, morphological, and semantic differences. For example, *“procedures for monitoring compliance with financial regulations”* corresponds to *“moliyaviy tartiblarga rioya etilishini nazorat qilish tartiblari”* (procedures for monitoring compliance with financial regulations), demonstrating that hybrid methods can reconcile differences in word order, morphology, and nested structure while preserving meaning.

Accurate extraction of multiword and nested expressions enhances legal information retrieval, document analysis, and multilingual legal knowledge management. Recognizing complete expressions ensures semantic precision, reduces ambiguity, and enables reliable cross-linguistic mapping. English benefits from syntactic consistency and predictable clause patterns, while Uzbek requires modeling of morphological and contextual variability. Hybrid approaches that integrate statistical and neural methods provide the most robust solution, enabling accurate recognition of multiword expressions and nested terms, and supporting advanced applications in automated contract analysis, regulatory compliance monitoring, and cross-border legal research.

Qualitative analysis in legal text processing highlights practical insights into how extraction methods perform beyond mere numerical evaluation, revealing both strengths and limitations in recognizing complex legal expressions. This approach examines the ability of models to handle multiword expressions, procedural statements, nested terms, and domain-specific vocabulary, offering a deeper understanding of successes and persistent challenges in real-world applications. Comparing English and Uzbek corpora demonstrates how linguistic, morphological, and structural differences influence extraction accuracy,

semantic interpretation, and cross-linguistic mapping.

In English legal texts, notable successes are observed in recognizing procedural and regulatory expressions with precise semantic meaning. Phrases such as *“procedures for registering intellectual property rights with the patent office”* and *“requirements for maintaining confidentiality in arbitration agreements”* are correctly extracted in their entirety, including nested terms like *“intellectual property rights”* and *“arbitration agreements”*. Similarly, *“guidelines for submitting petitions to labor dispute tribunals”* demonstrates that neural and hybrid methods can accurately identify both the procedural context and the nested legal concept of *“labor dispute tribunals”* [33]. These successes indicate that advanced models can handle multi-clause constructions, domain-specific jargon, and nested legal expressions, ensuring accurate semantic representation.

In Uzbek legal texts, extraction successes are evident for frequently used procedural phrases and established legal constructs. Examples include, *“mehnat nizolari bo'yicha ishlarni ko'rib chiqish tartibi”* (procedure for reviewing labor dispute cases), *“ijtimoiy sug'urta to'lovlari hisobga olish va monitoring qilish mexanizmlari”* (mechanisms for recording and monitoring social insurance contributions), *“davlat mulkini sotish va ijro etish tartibi”* (procedure for selling and executing state property), and *“solliq deklaratsiyalarini topshirish va tasdiqlash tartibi”* (procedure for submitting and validating tax declarations). Nested terms are also accurately recognized, such as *“ijtimoiy sug'urta to'lovlari hisobga olish”* (recording of social insurance contributions) and *“solliq deklaratsiyalarini tasdiqlash”* (validation of tax declarations) [34]. Neural and hybrid models successfully manage suffix variations, agglutination, and flexible word order, allowing complete extraction of procedural and regulatory expressions in context.

Despite these successes, several challenges persist. In English, expressions with conditional clauses, rare constructions, or domain-specific jargon may be partially fragmented. For instance, *“requirements for compliance with environmental impact assessments prior to construction approvals”* contains nested expressions such as *“environmental impact assessments”* and *“construction approvals”*, which may be split if contextual modeling is insufficient. Similarly, *“rights of shareholders to participate in emergency board meetings under corporate governance statutes”* requires understanding of nested legal entities like *“emergency board meetings”* and *“corporate governance statutes”* to avoid semantic loss [35].

In Uzbek, challenges are more pronounced due to morphological complexity, agglutination, and syntactic flexibility. Examples include (procedure for formalizing and validating social protection documents) and *“fuqarolik shartnomalaridagi tomonlarning mas’uliyatlarini belgilash va monitoring qilish mexanizmlari”* (mechanisms for defining and monitoring the responsibilities of parties in civil contracts) [36]. These expressions contain multiple nested terms and variable suffixation, making it difficult for statistical methods to extract the entire expression correctly. Neural and hybrid approaches mitigate these issues by capturing contextual and semantic relationships, yet extremely long or infrequently used expressions remain challenging for full extraction.

Semantic ambiguity further complicates extraction in both languages. In English, terms such as *“execution”*, *“appeal”*, or *“provision”* may carry general meanings or domain-specific legal interpretations depending on context. For example, *“execution of administrative orders in accordance with statutory guidelines”* requires recognizing *“execution”* in its legal sense rather than the general sense of carrying out an action [37]. In Uzbek, ambiguity arises from agglutination and morphological variation. For instance, *“mas’uliyatlarni belgilash”* (defining responsibilities) and *“majburiyatlarni bekor qilish”* (termination of obligations) share similar roots but convey opposite legal consequences. Accurate extraction depends on understanding the surrounding context to preserve intended meaning.

Cross-linguistic alignment provides additional insights. English procedural expressions such as *“guidelines for submitting claims to consumer protection agencies”* align with Uzbek equivalents like *“iste’molchilar huquqlarini himoya qilish organlariga da’volarni topshirish tartibi”* (procedure for submitting claims to consumer protection agencies). Similarly, *“requirements for filing corporate tax reports under national legislation”* corresponds to *“milliy qonunchilikka muvofiq korporativ soliq hisobotlarini topshirish tartibi”* (procedure for submitting corporate tax reports in accordance with national legislation) [38]. These examples demonstrate that hybrid methods can reconcile structural, morphological, and semantic differences while maintaining legal accuracy and coherence.

Qualitative analysis emphasizes that modern NLP approaches are capable of extracting complex multiword and nested expressions with notable success in both English and Uzbek legal texts. Standardized procedural phrases, nested terms, and cross-linguistic mappings are among the most

reliable outcomes. Persistent challenges include handling rare multi-clause expressions, semantic ambiguity, morphologically complex forms in Uzbek, and domain-specific jargon. Neural-hybrid approaches offer a robust solution, enhancing accuracy, preserving semantic integrity, and supporting practical applications such as automated contract review, regulatory compliance verification, and comparative legal research across languages.

5 DISCUSSION

Hybrid approaches in legal natural language processing have proven to be highly effective in addressing the intricate and variable nature of legal texts. By combining statistical methods, which detect patterns and high-frequency structures, with neural approaches, which capture contextual meaning and semantic relationships, hybrid models achieve a level of accuracy and robustness that is difficult to reach with either method alone. This capability is particularly important in legal texts, where multiword expressions, nested constructions, and domain-specific terminology are pervasive, and precise extraction is essential for reliable interpretation and analysis.

In English legal corpora, hybrid approaches excel at extracting procedural guidelines, contractual clauses, and regulatory statements. Phrases such as *“guidelines for submitting evidence in consumer protection cases”* or *“procedures for amending partnership agreements in corporate law”* are correctly identified, including nested components like *“consumer protection cases”* and *“partnership agreements”*. Similarly, *“requirements for reporting financial irregularities to auditing authorities”* shows the ability of hybrid models to capture both the overarching legal process and the embedded terms *“financial irregularities”* and *“auditing authorities”*, preserving the full semantic content [39]. These examples demonstrate that hybrid models can effectively handle complex multi-clause constructions and domain-specific expressions in English.

Uzbek legal texts present additional challenges due to agglutinative morphology, flexible word order, and rich suffixation. Hybrid approaches successfully navigate these challenges, accurately extracting procedural and regulatory expressions. For instance, *“fuqarolik shartnomalariga oid nizolarni hal etish tartibi”* (procedure for resolving disputes related to civil contracts) includes the nested term *“nizolar”* (disputes) and conveys the full procedural meaning. Another example, *“ijtimoiy himoya to’lovlari hisobga olish va tekshirish mexanizmlari”*

(mechanisms for recording and verifying social protection contributions), demonstrates that hybrid models can capture both frequent procedural expressions and embedded concepts like *“ijtimoiy himoya to'lovlari”* (social protection contributions). Similarly, *“davlat mol-mulkini sotish va ijro etish tartibi”* (procedure for selling and executing state property) and *“solliq deklaratsiyalarini topshirish va tasdiqlash tartibi”* (procedure for submitting and validating tax declarations) show that hybrid approaches effectively reconcile morphological variability and nested structures to preserve semantic integrity [40].

Hybrid models are particularly valuable for aligning English and Uzbek legal expressions across languages. For example, *“procedures for filing administrative appeals in regulatory tribunals”* aligns with *“ma'muriy sudlarga e'tiroz bildirish tartibi”* (procedure for submitting appeals to administrative courts), and *“requirements for maintaining corporate compliance records”* corresponds to *“korporativ murojaat hujjatlarini saqlash talablari”* (requirements for maintaining corporate compliance documents) [41]. These examples illustrate that hybrid approaches can preserve the legal meaning of multiword and nested expressions while resolving structural and morphological differences between languages.

The strength of hybrid models also becomes evident in handling expressions with low frequency or rare constructions. In English, expressions such as *“guidelines for enforcing environmental permits in regional tribunals”* are accurately recognized, capturing nested terms like *“environmental permits”* and *“regional tribunals”*. In Uzbek, complex expressions such as *“ijtimoiy sug'urta majburiyatlariga oid hujjatlarni tayyorlash va tekshirish tartibi”* (procedure for preparing and verifying documents related to social insurance obligations) and *“fuqarolik majburiyatlarini buzish holatlarini aniqlash va qarorlar qabul qilish mexanizmlari”* (mechanisms for identifying breaches of civil obligations and issuing decisions) are correctly extracted, preserving both the main procedural phrase and nested terms [42]. These examples demonstrate that hybrid models combine the precision of statistical pattern recognition with the contextual understanding of neural methods, enabling accurate extraction even in morphologically complex and syntactically flexible languages.

Hybrid approaches in legal NLP demonstrate superior effectiveness by integrating the complementary strengths of statistical and neural methods. They are capable of accurately extracting multiword and nested expressions, preserving semantic coherence, and managing the

morphological and syntactic complexity of both English and Uzbek legal texts. Hybrid models enable reliable cross-linguistic alignment, handle rare and complex constructions, and capture domain-specific terminology, making them indispensable for automated contract analysis, regulatory compliance assessment, and multilingual legal research. Their effectiveness underscores the importance of combining multiple computational techniques to address the inherent challenges of legal language processing.

Linguistic Challenges in Uzbek: Morphology, Agglutination, Flexible Word Order

Uzbek legal texts present distinctive linguistic challenges due to the language's rich morphological system, agglutinative structure, and flexible word order. Unlike English, which relies heavily on fixed word order and prepositions to convey grammatical and legal relationships, Uzbek encodes meaning through complex suffixation, case marking, and verb morphology. These characteristics create specific obstacles for automatic recognition of legal expressions, multiword units, and nested concepts, complicating tasks such as term extraction, semantic interpretation, and cross-linguistic alignment.

Morphological complexity is one of the primary challenges in processing Uzbek legal texts. Words frequently carry multiple suffixes that encode grammatical information critical for understanding legal meaning. For example, *“xalqaro shartnomalarga muvofiq majburiyatlarni bajarish tartibi”* (procedure for fulfilling obligations under international agreements) contains *“majburiyatlarni”*, where the suffix marks both plurality and accusative case, defining the object of the verb. Another example is *“mehnat shartnomalari bo'yicha oylik to'lovlarni hisobga olish mexanizmlari”* (mechanisms for accounting monthly payments under employment contracts), where *“to'lovlarni”* carries plurality and object markers simultaneously [43]. Such morphological richness challenges conventional statistical NLP models, which may treat inflected forms as separate entities, leading to fragmented extraction. Neural approaches that capture contextual embeddings partially overcome this limitation, but precise handling of morphological nuances remains critical for accurate term recognition.

Agglutination in Uzbek further complicates automated extraction. Words are often formed by sequentially attaching multiple suffixes, resulting in long tokens that encode complex legal meaning. For instance, *“solliq deklaratsiyalarini topshirish va nazorat qilish tartibi”* (procedure for submitting and

monitoring tax declarations) includes “*deklaratsiyalarini*”, where the suffixes indicate plurality, object, and possession simultaneously. Similarly, “*ijtimoiy nafaqalarni taqdim etish va tekshirish mexanizmlari*” (mechanisms for submitting and verifying social benefits) demonstrates the challenge of separating meaningful components within highly agglutinated structures [44]. Agglutination increases surface variability, making exact string matching ineffective. Neural embedding approaches can identify semantic similarity across variants, yet long or rare constructions may still be partially fragmented, affecting extraction completeness.

Flexible word order presents an additional challenge for Uzbek legal NLP. Unlike English, where syntactic roles are largely fixed, Uzbek allows reordering of sentence elements without changing meaning. For example, “*ijtimoiy nafaqalarni taqdim etish va tekshirish mexanizmlari*” (mechanisms for submitting and verifying social benefits) may appear as “*tekshirish va taqdim etish mexanizmlari ijtimoiy nafaqalarni*” while preserving semantic content. Another instance, “*fuqarolik kodeksi bo'yicha ishlarni ko'rib chiqish tartibi*” (procedure for reviewing cases under the civil code), may be reordered as “*ko'rib chiqish tartibi fuqarolik kodeksi bo'yicha ishlarni*”, emphasizing different elements stylistically [45]. Flexible word order requires models to infer syntactic relationships and semantic roles from context rather than relying on linear position, making hybrid statistical-neural approaches particularly effective.

The combination of morphological complexity, agglutination, and flexible word order compounds the difficulty of extracting multiword expressions. Consider the term “*ijtimoiy himoya hujjatlarini rasmiylashtirish va monitoring qilish tartibi*” (procedure for formalizing and monitoring social protection documents) [46]. This phrase contains multiple nested terms, morphologically complex tokens, and variable word order, demanding robust parsing and contextual understanding. Accurate extraction requires models to identify both procedural units and embedded concepts while preserving the semantic relationships among them.

Semantic ambiguity further exacerbates extraction challenges. Words sharing a root may express distinct legal concepts depending on their suffixes and surrounding context. For instance, “*qarzlarni belgilash*” (defining loans), “*qarzlarni to'lash*” (repayment of loans), and “*qarzlarni kechiktirish*” (postponement of loans) demonstrate how suffixes alter meaning in subtle but legally significant ways. Similarly, “*litsenziyani rasmiylashtirish*” (licensing procedure), “*litsenziyani*

bekor qilish” (revocation of license), and “*litsenziya shartlarini tekshirish*” (inspection of license conditions) illustrate nested terms with semantic variation that must be carefully interpreted [47]. Neural approaches that leverage contextual embeddings improve handling of such cases, yet precise extraction relies on accurate morphological and syntactic modeling.

Cross-linguistic alignment presents additional challenges. English terms like “*procedures for registering corporate mergers*” correspond to “*korporativ birlashmalarni ro'yxatdan o'tkazish tartibi*” (procedure for registering corporate mergers), where Uzbek suffixes encode case, possession, and grammatical relationships that English expresses through prepositions. Similarly, “*requirements for reviewing anti-corruption reports*” maps to “*korruptsiyaga oid hisobotlarni ko'rib chiqish talablari*” (requirements for reviewing anti-corruption reports) [48]. Hybrid NLP approaches are essential for aligning morphologically complex Uzbek terms with English equivalents while preserving legal meaning.

Uzbek legal texts pose significant linguistic challenges due to rich morphology, agglutination, and flexible word order. Morphologically complex and agglutinated tokens combined with variable sentence structure complicate multiword expression extraction, nested term recognition, and semantic interpretation. Hybrid NLP approaches that integrate statistical pattern recognition with neural contextual embeddings are particularly effective in overcoming these challenges. They enable accurate recognition of complex legal expressions, maintain semantic integrity, and facilitate cross-linguistic alignment with English legal texts. Addressing these challenges is crucial for developing reliable automated legal document analysis systems, multilingual legal knowledge management tools, and comparative research involving Uzbek and English legal corpora.

Cross-Linguistic Semantic Alignment: Terms with Direct and Contextual Mappings

Cross-linguistic semantic alignment is a critical challenge in multilingual legal natural language processing, particularly when dealing with languages as structurally and morphologically distinct as English and Uzbek. The task involves identifying terms that are directly equivalent across languages, as well as those whose meaning depends on contextual interpretation, procedural nuance, or domain-specific usage. Achieving accurate alignment is essential for applications such as bilingual legal corpora development, comparative

law studies, automated translation, and cross-border regulatory compliance.

Direct term mapping refers to instances where a legal expression in English has a clear and consistent equivalent in Uzbek. For example, the English term *“commercial contract obligations”* can be directly mapped to *“tijorat shartnomalari bo'yicha majburiyatlar”* (obligations under commercial contracts). Similarly, *“intellectual property rights enforcement”* corresponds directly to *“intellektual mulk huquqlarini ijro etish”* (enforcement of intellectual property rights) [49]. These direct mappings often involve standardized terminology that appears frequently in legal texts, enabling relatively straightforward extraction and alignment by both statistical and neural models. Hybrid approaches are particularly effective in ensuring that morphological variation in Uzbek, such as case marking and possessive suffixes, does not interfere with the recognition of direct equivalents.

Contextual mapping, on the other hand, involves aligning terms whose meaning depends on syntactic, semantic, or procedural context rather than one-to-one correspondence. For instance, the English phrase *“procedures for filing bankruptcy claims”* may align with the Uzbek expression *“bankrotlik da'volarini ko'rib chiqish tartibi”* (procedure for reviewing bankruptcy claims), where the term *“filing”* is rendered in Uzbek as *“ko'rib chiqish”* (reviewing), reflecting the procedural rather than literal action. Another example is *“regulatory compliance verification,”* which corresponds to *“nazorat qiluvchi organlar tomonidan talablarni tekshirish”* (verification of requirements by regulatory bodies) [50]. Contextual mapping requires models to infer the intended legal meaning and procedural role of a term, rather than relying solely on lexical equivalence. Neural embeddings and contextualized language models are particularly suited for capturing these nuanced relationships.

Cross-linguistic alignment becomes more complex when multiword expressions and nested terms are involved. Consider the English term *“shareholder voting rights under corporate governance regulations,”* which corresponds to the Uzbek expression *“korporativ boshqaruv qoidalari bo'yicha aksiyadorlarning ovoz berish huquqlari”* (voting rights of shareholders under corporate governance regulations). Accurate alignment requires capturing the nested structure, where *“voting rights”* is embedded within *“under corporate governance regulations,”* and preserving the syntactic and semantic relationships in both languages. Similarly, *“mechanisms for implementing environmental compliance*

measures” aligns with *“atrof-muhitga oid talablarni amalga oshirish mexanizmlari”* (mechanisms for implementing environmental requirements), demonstrating how contextual understanding ensures correct semantic interpretation across linguistic structures [51].

Morphological richness in Uzbek adds additional alignment challenges. Suffixes indicating case, possession, and plurality may vary depending on sentence structure, creating multiple surface forms of the same legal concept. For example, *“ijtimoiy nafaqalarni hisobga olish tartibi”* (procedure for recording social benefits) can appear in sentences as *“hisobga olish tartibi ijtimoiy nafaqalarni,”* requiring models to recognize both forms as equivalent to the English term *“procedure for recording social benefits”* [52]. Agglutination and flexible word order necessitate sophisticated parsing and contextual embedding to ensure semantic equivalence, particularly for less frequent or domain-specific terms.

Semantic ambiguity further complicates cross-linguistic alignment. English words such as *“execution”* or may carry general meanings or domain-specific legal interpretations depending on context. For instance, *“execution of administrative penalties”* aligns with *“ma'muriy jazolarni ijro etish”* (enforcement of administrative penalties), where *“execution”* is understood in the legal, not general, sense. In Uzbek, ambiguous roots such as *“majburiyat”* (obligation) or *“qaror”* (decision) require contextual analysis to determine whether they refer to fulfillment, termination, or monitoring of obligations and decisions [53]. Effective cross-linguistic alignment therefore relies on models that can incorporate procedural, semantic, and syntactic context to preserve intended meaning.

Hybrid models demonstrate notable effectiveness in cross-linguistic semantic alignment by combining the strengths of statistical term frequency recognition and neural contextual understanding. Statistical approaches help identify recurring multiword patterns, while neural embeddings capture subtle semantic nuances and nested relationships, allowing accurate mapping of both direct and context-dependent terms. For example, English *“requirements for submitting regulatory filings”* aligns with Uzbek *“nazorat hujjatlarini topshirish talablari”* (requirements for submitting regulatory documents), maintaining the procedural nuance while resolving morphological differences. Similarly, *“guidelines for dispute resolution in commercial contracts”* corresponds to *“tijorat shartnomalaridagi nizolarni hal qilish bo'yicha yo'riqnomalar”* (guidelines for resolving disputes in

commercial contracts), demonstrating the ability to align multiword expressions with nested legal entities [54].

Cross-linguistic semantic alignment in legal NLP requires careful handling of both direct term equivalents and contextually mapped expressions. English and Uzbek legal corpora illustrate the necessity of addressing morphological variation, agglutination, flexible word order, nested terms, and procedural nuance. Hybrid NLP approaches, combining statistical and neural methods, are particularly effective in achieving accurate alignment, ensuring semantic integrity, and supporting applications such as multilingual legal knowledge bases, automated document analysis, and comparative legal research. Successful semantic alignment facilitates reliable interpretation across languages, enhancing the usability of bilingual and multilingual legal NLP systems.

Implications, Limitations, and Future Directions

The findings of this study carry important implications for the development of multilingual legal natural language processing systems, particularly in relation to corpus design and automated legal knowledge extraction. The complexity of legal language, combined with cross-linguistic variation between English and Uzbek, underscores the necessity of adopting hybrid computational approaches that integrate statistical pattern recognition with neural contextual modeling. Such approaches enhance the ability to identify and extract multiword expressions, procedural statements, and embedded legal concepts while preserving semantic integrity across languages. In multilingual legal NLP, accurate recognition of expressions such as *“procedures for issuing digital compliance certificates”* and their Uzbek equivalents like *“raqamli muvofiqlik sertifikatlarini berish tartibi”* (procedure for issuing digital compliance certificates) demonstrates the importance of aligning syntactic and morphological structures in a way that maintains legal meaning [55]. This has direct implications for automated legal knowledge extraction, where the goal is not only to identify terms but also to structure them into usable knowledge representations for tasks such as legal reasoning, document classification, and regulatory monitoring.

Corpus design emerges as a critical factor influencing the effectiveness of multilingual legal NLP systems. A well-constructed corpus must account for domain diversity, including areas such as administrative law, financial regulation, and contractual obligations, while also reflecting

authentic language use. In Uzbek, this requires careful consideration of morphological variation and agglutinative forms, ensuring that different surface realizations of the same legal concept are adequately represented. For example, expressions like *“elektron hujjatlarni tasdiqlash jarayoni”* (process of validating electronic documents) and *“davlat xizmatlarini ko'rsatish shartlarini belgilash tartibi”* (procedure for defining conditions of public service delivery) illustrate the need for corpora that capture both structural variation and contextual usage [56]. Parallel or comparable corpora are particularly valuable for cross-linguistic alignment, enabling models to learn correspondences between English and Uzbek legal expressions even when they differ in form. Furthermore, annotation schemes must go beyond simple term labeling to include nested structures, semantic roles, and contextual relationships, thereby supporting more advanced knowledge extraction.

Automated legal knowledge extraction benefits significantly from these advancements, as improved term recognition and alignment enable the construction of structured legal ontologies and knowledge graphs. For instance, extracting relationships from expressions such as *“requirements for maintaining digital transaction records”* and mapping them to Uzbek equivalents like *“raqamli tranzaksiya yozuolarini saqlash talablari”* (requirements for maintaining digital transaction records) allows systems to build interconnected representations of legal obligations, procedures, and entities [57]. This, in turn, supports applications such as intelligent legal search, automated compliance checking, and decision-support systems. However, achieving this level of functionality requires models that can handle morphological complexity, nested expressions, and semantic variation, particularly in under-resourced languages.

Despite these advances, several limitations must be acknowledged. One of the primary constraints is corpus size, especially for Uzbek legal texts. Limited availability of large, high-quality, annotated corpora restricts the ability of machine learning models to generalize effectively, particularly for rare or domain-specific expressions. Smaller corpora may also lack sufficient representation of diverse legal subdomains, leading to gaps in coverage and reduced robustness. Annotation scope presents another limitation, as manual annotation of legal texts is time-consuming and resource-intensive. In many cases, annotations may focus primarily on surface-level terms without fully capturing nested structures, semantic roles, or contextual

dependencies, thereby limiting the depth of analysis that models can achieve.

Language-specific constraints further complicate the development of multilingual legal NLP systems. Uzbek's agglutinative morphology and flexible word order introduce variability that is not present in English, requiring specialized preprocessing techniques and model adaptations. For example, expressions like "*sun'iy intellekt tizimlaridan foydalanish qoidalarini belgilash tartibi*" (procedure for defining rules for the use of artificial intelligence systems) may appear in multiple morphologically distinct forms, challenging both term extraction and alignment [58]. Additionally, differences in legal systems and terminological conventions between English-speaking jurisdictions and Uzbek legal frameworks may lead to mismatches in conceptual equivalence, even when linguistic alignment is achieved.

Future research directions should focus on expanding both the linguistic and computational scope of multilingual legal NLP. Incorporating additional languages, particularly those with similar typological features such as other Turkic languages, would provide valuable insights into the generalizability of hybrid approaches and support the development of more inclusive multilingual systems. Expanding corpus size and diversity is also essential, including the integration of legislative texts, and administrative documents to ensure comprehensive coverage of legal domains. Enhanced annotation strategies that capture multiword expressions, nested terms, and semantic relationships will further improve model performance and enable more sophisticated knowledge extraction.

Model optimization represents another key area for future work. Advances in transformer-based architectures, domain-adaptive pretraining, and cross-lingual embeddings offer promising avenues for improving the accuracy and efficiency of legal NLP systems. Fine-tuning models on domain-specific corpora, incorporating morphological analyzers for agglutinative languages, and developing hybrid pipelines that integrate rule-based and machine learning components can significantly enhance performance. Additionally, exploring methods for low-resource learning, such as transfer learning and data augmentation, will be crucial for addressing the limitations associated with smaller corpora.

The development of multilingual legal NLP systems requires a holistic approach that integrates robust corpus design, advanced modeling techniques, and careful consideration of linguistic

diversity. While significant progress has been made in extracting and aligning legal expressions across English and Uzbek, challenges related to corpus size, annotation depth, and language-specific constraints remain. Addressing these limitations through expanded datasets, improved annotation practices, and optimized models will pave the way for more accurate, scalable, and practical applications in automated legal knowledge extraction and multilingual legal analysis.

6 CONCLUSION

This study set out to investigate the effectiveness of advanced natural language processing techniques for legal term extraction through a comparative, corpus-driven analysis of English and Uzbek legal texts. The central aim was to explore how hybrid approaches, combining statistical and neural methodologies, can address the linguistic, structural, and semantic complexities inherent in multilingual legal corpora. By focusing on the extraction of multiword expressions, nested constructions, and context-dependent legal units, the research contributes to the development of more accurate and semantically robust legal NLP systems, particularly in the context of typologically distinct languages.

The findings demonstrate that a hybrid, corpus-driven framework provides a highly effective solution for legal term recognition and semantic alignment. Statistical methods contribute by identifying recurrent patterns and high-frequency expressions, while neural models enhance the system's ability to capture contextual meaning, semantic dependencies, and variation across syntactic structures. This integration enables improved handling of complex legal language, ensuring that both structural and semantic equivalence are preserved across languages. The study confirms that hybrid approaches outperform single-method models in terms of completeness, consistency, and contextual accuracy, particularly when dealing with morphologically rich and syntactically flexible languages.

A key contribution of this research lies in its emphasis on semantic interoperability across languages. The ability to align legal expressions not only at the lexical level but also at the conceptual and procedural levels represents a significant advancement for multilingual legal NLP. This capability is essential for building reliable multilingual legal resources, enabling consistent interpretation of legal concepts, and supporting cross-border legal communication. By addressing both linguistic and semantic variation, the proposed

approach enhances the quality and reliability of multilingual legal analysis.

The study also highlights the importance of corpus design in achieving high-quality term extraction and semantic alignment. A well-balanced corpus that captures domain diversity, linguistic variation, and authentic usage patterns significantly enhances model performance. In addition, enriched annotation practices that go beyond surface-level labeling to include nested structures and semantic relationships are shown to be crucial for supporting advanced knowledge extraction. These factors collectively contribute to the effectiveness of NLP systems in processing complex legal texts.

From a practical perspective, the findings have important implications for legal translation, multilingual legal resource development, and AI-assisted legal research. Improved term extraction and alignment facilitate more accurate and consistent translation of legal content across languages, reducing ambiguity and enhancing clarity. In the context of multilingual legal resources, the ability to extract and structure legal knowledge supports the development of dictionaries, ontologies, and knowledge graphs, which are essential for information retrieval, document classification, and decision support. AI-assisted legal research benefits from enhanced semantic understanding, enabling more precise search capabilities and more effective analysis of legal documents.

Despite these contributions, several limitations must be acknowledged. Constraints related to corpus size and diversity may affect the generalizability of the findings, particularly for low-frequency or domain-specific expressions. The scope of annotation also presents challenges, as limited annotation depth may restrict the ability to capture complex semantic relationships and nested structures. Additionally, language-specific characteristics, including morphological complexity and syntactic flexibility,

require specialized modeling techniques, which may not be fully addressed within the current framework.

Future research should focus on expanding the linguistic and computational scope of multilingual legal NLP. Incorporating additional languages will provide insights into the adaptability of hybrid approaches across different typological systems. Expanding and diversifying legal corpora will improve coverage and robustness, particularly for specialized domains. Enhancing annotation methodologies to include richer semantic and structural information will further support advanced knowledge extraction and representation.

Advancements in model optimization also represent an important direction for future work. The integration of domain-adaptive training techniques, cross-lingual representations, and morphological processing tools can significantly improve system performance. Exploring approaches for low-resource learning, such as transfer learning and data augmentation, will be essential for addressing limitations associated with smaller datasets. Additionally, the development of more transparent and interpretable models will contribute to the practical adoption of legal NLP systems in professional contexts.

In conclusion, this research demonstrates that a hybrid, corpus-driven approach provides a robust and scalable solution for legal term extraction and cross-linguistic semantic alignment. By effectively addressing challenges related to morphology, syntax, and semantic variation, the proposed framework enhances the accuracy and reliability of multilingual legal text processing. The study offers valuable contributions to the field of legal NLP, with practical implications for translation, knowledge management, and AI-assisted research. Continued efforts to expand datasets, refine methodologies, and optimize models will further advance the development of semantically interoperable and application-ready multilingual legal systems

REFERENCES

- Abdullaev, "Corpus-based analysis of legal texts in Uzbek," *Journal of Turkic Linguistics*, vol. 5, no. 2, pp. 73–88, 2024.
- Ahmed, "Semantic interoperability in multilingual legal systems," *Journal of Information Science*, vol. 49, no. 1, pp. 12–28, 2023.
- Aliyev, "Challenges in processing agglutinative languages for NLP," *Central Asian Linguistics Journal*, vol. 7, no. 1, pp. 25–40, 2023.
- Anderson, "Natural language processing in legal systems: emerging trends," *Journal of Legal Informatics*, vol. 12, no. 2, pp. 45–60, 2023.
- Artetxe and Schwenk, "Massively multilingual sentence embeddings," *Transactions of ACL*, vol. 7, pp. 597–610, 2019.

- Baker, "Corpus linguistics and translation studies," *Applied Linguistics Review*, vol. 11, no. 2, pp. 215–230, 2020.
- Bakker, "The intersection of law and technology: challenges and opportunities," *International Journal of Law and Societal Studies*, vol. 1, no. 1, pp. 11–21, 2024.
- Bird, Klein, and Loper, "Natural Language Processing with Python," O'Reilly Media, 2009.
- Blei, Ng, and Jordan, "Latent Dirichlet Allocation," *Journal of Machine Learning Research*, vol. 3, pp. 993–1022, 2003.
- Brown, "Ontology-based representation of legal knowledge," *Knowledge Engineering Review*, vol. 36, no. 2, pp. 150–168, 2021.
- Brown et al., "Language models are few-shot learners," *NeurIPS*, vol. 33, pp. 1877–1901, 2020.
- Camacho-Collados and Pilehvar, "Word embeddings and semantic similarity," *Computational Linguistics*, vol. 44, no. 2, pp. 345–384, 2018.
- Chen, "Transformer models for legal NLP applications," *AI and Society*, vol. 38, no. 1, pp. 145–160, 2023.
- Choi, "Contextual embeddings in domain-specific NLP," *Journal of Advanced Computing*, vol. 7, no. 3, pp. 130–147, 2023.
- Clark, "Neural networks for sequence labeling in NLP," *Journal of Machine Learning Research*, vol. 21, no. 5, pp. 1–20, 2020.
- Conneau et al., "Unsupervised cross-lingual representation learning," *ACL Proceedings*, pp. 140–150, 2020.
- Davis, "AI-assisted legal research systems," *Legal Innovation Journal*, vol. 3, no. 1, pp. 22–38, 2024.
- Devlin et al., "BERT: Pre-training of deep bidirectional transformers for language understanding," *NAACL Proceedings*, pp. 4171–4186, 2019.
- Dozat and Manning, "Deep biaffine attention for neural dependency parsing," *ICLR Proceedings*, 2017.
- Evans, "Challenges of low-resource languages in NLP," *Language Technology Journal*, vol. 16, no. 2, pp. 95–112, 2022.
- Garcia, "Legal term extraction in multilingual corpora," *Language Resources and Evaluation*, vol. 55, no. 3, pp. 345–362, 2021.
- Goldberg, "A primer on neural network models for NLP," *Journal of Artificial Intelligence Research*, vol. 57, pp. 345–420, 2016.
- Green, "Evaluation metrics for NLP systems in law," *Journal of Computational Evaluation*, vol. 5, no. 2, pp. 70–86, 2022.
- Hassan, "Machine learning models for legal text classification," *Artificial Intelligence and Law*, vol. 30, no. 4, pp. 201–220, 2022.
- Hernandez, "Automated contract analysis using AI," *Journal of Law and Technology*, vol. 14, no. 4, pp. 155–172, 2022.
- Ivanov, "Comparative analysis of Slavic and Turkic legal terminology," *Linguistic Studies Journal*, vol. 20, no. 1, pp. 33–49, 2021.
- Johnson, "Text mining techniques for regulatory compliance," *Journal of Legal Analytics*, vol. 6, no. 3, pp. 75–92, 2022.
- Jurafsky and Martin, "Speech and Language Processing," Prentice Hall, 3rd ed., 2023.
- Karimov, "Comparative study of English and Uzbek legal terminology," *Central Asian Legal Studies*, vol. 2, no. 1, pp. 15–30, 2024.
- Khujakulov, S., Sociolinguistic and normative-stylistic methods for analyzing legal terms (based on Uzbek and English languages). *Jurislinguistics*, 35(46), 57–61., 2025. [https://doi.org/10.14258/leglin\(2025\)3509](https://doi.org/10.14258/leglin(2025)3509)
- Khujakulov, S., Semantic interoperability in legal translation: bridging the gap between English and Uzbek through ontological modeling. *Jurislinguistics*, 37(48), 16–21. 2025. [https://doi.org/10.14258/leglin\(2025\)3702](https://doi.org/10.14258/leglin(2025)3702)
- Khan, "Multilingual legal information retrieval systems," *Information Processing and Management*, vol. 59, no. 2, pp. 102–118, 2022.
- Kiperwasser and Goldberg, "Simple and accurate dependency parsing," *ACL Proceedings*, pp. 199–204, 2016.
- Kumar, "Semantic analysis of legal discourse using NLP techniques," *Journal of Language and Law*, vol. 9, no. 1, pp. 34–50, 2023.
- Lee, "Hybrid approaches to multilingual text processing," *International Journal of Computational Linguistics*, vol. 15, no. 2, pp. 66–82, 2021.
- Lewis et al., "BART: Denoising sequence-to-sequence pre-training," *ACL Proceedings*, pp. 7871–7880, 2020.

- Liu, "Pre-trained language models for NLP: A survey," *IEEE Transactions on Pattern Analysis*, vol. 44, no. 4, pp. 1872–1890, 2022.
- Lopez, "Data-driven approaches to legal text summarization," *Journal of Applied NLP*, vol. 8, no. 2, pp. 60–77, 2023.
- Manning et al., "Introduction to Information Retrieval," Cambridge University Press, 2008.
- Mikolov et al., "Efficient estimation of word representations in vector space," *ICLR Proceedings*, 2013.
- Müller, "Corpus linguistics in legal research," *European Journal of Legal Studies*, vol. 13, no. 2, pp. 120–138, 2020.
- Navigli, "Word sense disambiguation: A survey," *ACM Computing Surveys*, vol. 41, no. 2, pp. 1–69, 2009.
- Nguyen, "Named entity recognition in legal documents," *International Journal of AI Research*, vol. 11, no. 3, pp. 98–115, 2021.
- Omarov, "Linguistic features of Uzbek legal discourse," *Uzbek Linguistics Review*, vol. 6, no. 1, pp. 44–59, 2023.
- Park, "Cross-domain adaptation in legal NLP," *Journal of Machine Learning Applications*, vol. 9, no. 2, pp. 89–105, 2021.
- Pennington et al., "GloVe: Global vectors for word representation," *EMNLP Proceedings*, pp. 1532–1543, 2014.
- Peterson, "Corpus-driven approaches to legal terminology extraction," *Computational Linguistics Review*, vol. 18, no. 3, pp. 78–95, 2022.
- Peters et al., "Deep contextualized word representations," *NAACL Proceedings*, pp. 2227–2237, 2018.
- Petrova, "Lexical semantics in legal translation," *Translation Studies Review*, vol. 19, no. 1, pp. 61–79, 2022.
- Rahman, "Automated extraction of legal entities from text," *Journal of Data Science and Law*, vol. 5, no. 1, pp. 55–70, 2023.
- Rossi, "Knowledge graphs for legal information systems," *Semantic Web Journal*, vol. 13, no. 4, pp. 500–518, 2022.
- Ruder, "Neural transfer learning for NLP," PhD Thesis, National University of Ireland, 2019.
- Santos, "Information extraction from legal databases," *Journal of Big Data and Law*, vol. 4, no. 3, pp. 110–126, 2023.
- Singh, "Deep learning for legal document analysis," *IEEE Transactions on AI Systems*, vol. 11, no. 3, pp. 210–225, 2022.
- Smith, "Statistical and neural models for text mining in law," *Journal of Artificial Intelligence Research*, vol. 67, no. 1, pp. 101–120, 2020.
- Taylor, "Legal corpora and their applications in NLP," *Corpus Linguistics Quarterly*, vol. 17, no. 3, pp. 200–218, 2022.
- Vaswani et al., "Attention is all you need," *Advances in Neural Information Processing Systems*, pp. 5998–6008, 2017.
- Wang, "Neural embeddings for legal terminology alignment," *Computational Law Review*, vol. 10, no. 4, pp. 88–104, 2021.