

# Processing and Normalization of Uzbek Texts

Sobirov Shohjahon Ganijon o'g'li<sup>1</sup>, Sobirova Nazira Ganijon kizi<sup>2</sup>, Sobirova Zarnigor Ganijon kizi<sup>3</sup>

<sup>1</sup>*Department of Mechatronics and Robotics, Tashkent State Technical University Tashkent, Uzbekistan*  
SobirovShoh935@gmail.com

<sup>2</sup>*Computational Linguistics and Digital Technologies, Tashkent State University of Uzbek Language and Literature*  
Tashkent, Uzbekistan  
[nazirasobirova073@gmail.com](mailto:nazirasobirova073@gmail.com)

<sup>3</sup>*Computational Linguistics and Digital Technologies, Tashkent State University of Uzbek Language and Literature*  
Tashkent, Uzbekistan  
[sobirovazarnigor1996@gmail.com](mailto:sobirovazarnigor1996@gmail.com)

**Abstract** — Nowadays, technologies for the automatic analysis of texts through computers are rapidly advancing. In particular, when working with texts written in the Uzbek language, an essential step is their initial preparation — that is, preprocessing and text normalization. This article provides a detailed overview of these processes. Uzbek is an agglutinative language, which extensively employs derivational and grammatical affixes. This significantly complicates morphological analysis and the identification of word forms. Preprocessing involves operations such as removing unnecessary characters, splitting text into words, converting to lowercase, and similar tasks. Normalization, on the other hand, includes correcting misspelled words, converting words to their base dictionary forms, and expanding abbreviation. Due to the complex structure of the Uzbek language, these processes pose considerable challenges. The article discusses methods and tools developed to overcome these difficulties and illustrates their effectiveness with practical examples. The findings from ongoing research contribute to improving the quality of Uzbek language processing in digital environments, including applications in machine translation, speech-to-text systems, and automated text analysis.

**Keywords:** text normalization, Uzbek language, natural language processing (NLP), lemmatization, tokenization, stemming, agglutinative languages, Cyrillic-to-Latin conversion, transformer models, automatic error correction, T5.

## I. INTRODUCTION

In recent years, research in the field of Natural Language Processing (NLP) has increased the demand for automated text analysis in various languages, including Uzbek. The fact that texts written in Uzbek are multi-formed, morphologically complex, and stylistically diverse creates certain obstacles in their automatic processing. The central point of this article is the automatic normalization of texts in the Uzbek language, that is, the process of text normalization. The article is devoted to the study of linguistic and technological problems encountered in the process of automatic text normalization in Uzbek. The complex morphological structure of the Uzbek language, multi-form words, dialectal variants, Cyrillic-Latin orthographic differences, and non-standard expressions complicate this process. Uzbek is an agglutinative language, which makes extensive use of derivational and grammatical affixes. This complicates morphological analysis and the identification of word forms.

Text cleaning (preprocessing) in natural language processing refers to the process of removing unnecessary or inconsistent elements from the text and bringing the words into an appropriate form. High-quality text cleaning is very important for further analysis or modeling, especially in languages like Uzbek that have rich morphology and are represented in different scripts.

## II. PROCESSING OF TEXT

- a. Text cleaning in Uzbek includes several stages, such as: checking and correcting spelling errors,
- b. tokenization (splitting text into words),
- c. normalization (standardizing symbols and orthography in the text), and lemmatization (finding the lexical root of words).



Figure 1. Text Preprocessing Pipeline

Tokenization is the process of dividing sentences or phrases in the text into separate word units, and the subsequent lemmatization and analysis stages are applied specifically to these tokens. Although this task may seem simple at first glance, it can become complex due to certain linguistic features of the Uzbek language.

In Uzbek, words are usually separated by spaces, which helps identify word boundaries in most cases. However, there are several complicated cases where this rule does not apply:

#### A. Appendable Sections

For example, certain clitic words and particles can either be written with a hyphen or concatenated directly. The interrogative particle “-mi?” is sometimes written separately, but often attaches to the preceding word (as in “Keldingmi?”, meaning “Did you come?”). Similarly, numeral-year combinations such as “2025-yil” (the year 2025), and compound words like “ilm-fan” (science and knowledge), are typically written with hyphens. Basic segmentation algorithms may fail in such cases, mistakenly treating “2021” and “yil” as separate tokens, or splitting compound words like “ilm” and “fan” into individual elements.

#### B. Apostrof

In the Uzbek Latin alphabet, the apostrophe (') is used to represent certain letters or phonemes that are not directly pronounced, such as in “o'”, “g'”, and “sh'”. Tokenization at the position of the apostrophe is incorrect; however, some algorithms may mistakenly interpret the apostrophe as a word boundary marker, leading to segmentation errors.

#### C. Script Mixing

Nowadays, it is common to encounter the mixed use of Latin and Cyrillic scripts within texts—for instance, a Cyrillic word embedded in a predominantly Latin script text, or vice versa. This phenomenon complicates tokenization, as the system must correctly recognize and distinguish characters from both writing systems.

### III. DESCRIPTION OF SCIENTIFIC RESEARCH

Text normalization is the process of converting various writing forms, spelling errors, dialectal expressions, abbreviations, misspelled words, and non-standard phrases typical of social media into a standardized, dictionary-compliant format. This process determines the accuracy of subsequent tasks such as text analysis, classification, translation, or speech synthesis.

In the Uzbek language, Cyrillic-Latin orthographic differences, phonetic spellings, dialectal variations, and the presence of multiple morphological forms of a single word make the normalization process more complex. In recent years, several scientific works have been conducted to address orthographic, syntactic, and lexical aspects of text normalization in Uzbek. For example, M. Sharipov and O. Sobirov (2022) presented an algorithm for suffix separation and lemma identification in the Uzbek language using a finite-state automaton. In their study, they clearly demonstrated the difference between stemming and lemmatization, emphasizing the importance of identifying the correct base form of a word [1].

B. Elov et al. (2023) compared stemming and POS tagging tasks in Uzbek, Turkish, and Uyghur, and highlighted the problems and solutions for implementing stemming in agglutinative languages. They also showed that a hybrid approach — combining rule-based and statistical methods — can significantly improve performance [2].

Ulug'bek Salaev (2024) developed the UzMorphAnalyser model and software, which is designed to analyze all possible morphological forms of words in Uzbek. His study included a full list of grammatical segments of Uzbek words and the corresponding analysis rules. The model achieved 91% accuracy during testing, which is considered a high score for the Uzbek language [3], demonstrating the strong potential of morphological normalization tools.

Although there are few software products specifically designed for syntactic normalization in Uzbek, the Uzbek-UT treebank created within the Universal Dependencies framework by Kurbanova N. et al. (2025) serves as a resource for consistent syntactic annotation. This work also explores certain syntactic phenomena specific to Uzbek, such as compound verb constructions with auxiliary verbs [4]. Such treebanks serve as a foundational base for syntactic normalization research.

In addition, Bruno Guillaume (2021) introduced the GREW (Graph Rewriting) tool, which assists in maintaining and modifying syntactic annotations in corpora, indirectly contributing to syntactic normalization (e.g., standardizing syntactic trees across languages).

In the article published by E. Kuriyozov et al. (2021) as part of the UzWordNet project, a lexical-semantic network of the Uzbek language is introduced [5]. It includes relationships such as synonymy, antonymy, and hyponymy between words. This resource provides a scientific basis for grouping synonyms and merging redundant variants in lexical normalization.

Kh. Madatov et al. (2022) described the development of the Uzbek WordNet based on the Turkish WordNet and presented comparative methodologies [6].

The UzBERT model proposed by B. Mansurov (2021) is one of the first pre-trained transformer models built on large-scale Uzbek corpora [7]. Additionally, G. Matlatipov et al. (2022) developed a sentiment-labeled corpus for Uzbek. While these works are not directly about normalization, they emphasize that preprocessing steps like text cleaning, lowercasing, and removal of unnecessary characters played an important role in corpus preparation.

The aforementioned studies show that scientific research in the field of computational processing of the Uzbek language has accelerated. In particular, significant progress has been made in morphology and lexicography, and further syntactic and semantic analyses are anticipated in future studies. The outcomes of these scientific endeavors serve as the foundation for developing software and applied NLP systems, ultimately enabling more effective normalization and understanding of Uzbek texts.

#### IV. TEXT NORMALIZATION

Normalization refers to bringing the text into a consistent and uniform format. In Uzbek, normalization tasks include:

- a) Transliteration between Latin and Cyrillic scripts
- b) Correcting letter case (e.g., converting all to lowercase)
- c) Unifying apostrophes and diacritic marks
- d) Removing unnecessary characters

The main goal of text normalization is to convert texts obtained from various sources into a unified standard, thereby preparing them for further processing stages.

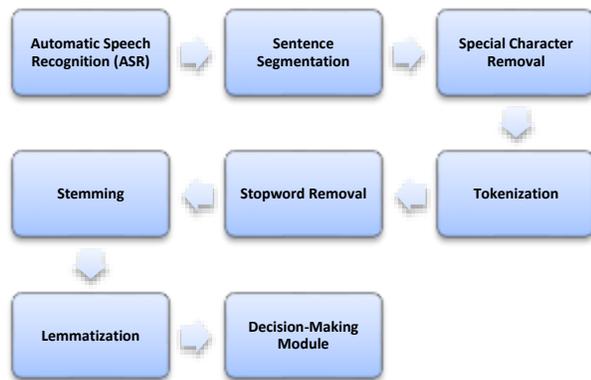


Figure 2. Pipeline of Automatic Speech Processing and Normalization

Text normalization refers to the process of transforming words and sentences into a consistent and standardized form. It is a fundamental step in Natural Language Processing (NLP).

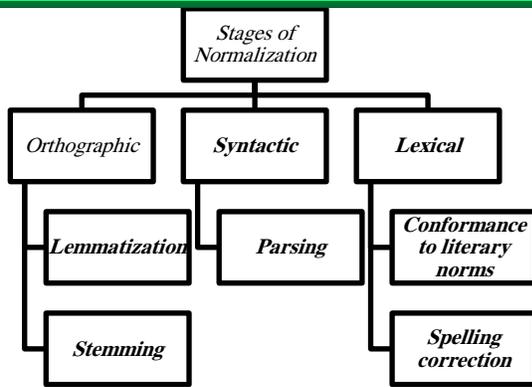


Figure 3. Stages of Text Normalization

**Normalization** is typically considered an essential and initial step. This is because, for processes such as **tokenization** and **lemmatization** to function correctly, the text must first be clean and in a standardized format.

In a well-normalized text, each word appears in a consistent orthographic form and adheres to spelling rules — which helps reduce errors in subsequent algorithms.

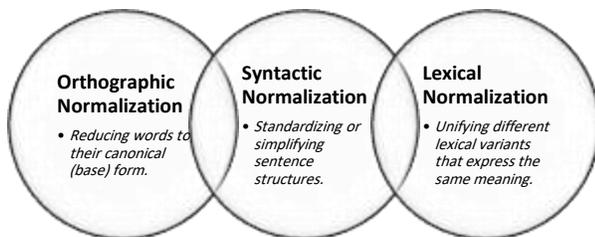


Figure 4. Types of Normalization

## V. TYPES OF NORMALIZATION

Orthographic normalization aims to reduce different morphological forms of a word to its base lexical form. This is mainly achieved through two approaches: lemmatization (identifying the lemma or dictionary form of a word) and stemming (removing affixes to extract the root). Although lemmatization and stemming share a similar goal, they differ in terms of the output form and the level of accuracy:

1. Stemming involves cutting off grammatical affixes to extract the invariant root part of the word. However, the resulting form may not always represent a meaningful or proper base form. For example, stemming the word “o‘qigan” would yield “o‘q”, but in Uzbek, “o‘q” also means “arrow”, which is semantically different [1].

Thus, stemming does not consider the meaning; it simply shortens the word through mechanical truncation.

2. Lemmatization is the process of identifying a word’s lemma, i.e., its dictionary form. This method converts different grammatical forms into a single lexical unit [1]. In the example above, the lemma of “o‘qigan” would be “o‘qimoq”, which is the infinitive form conveying the intended meaning of “to read”.

Therefore, lemmatization preserves the word’s meaning while converting it into its canonical form.

Word Form	Stemming Result	Lemmatization Result
o‘qigan	o‘q (may mean "arrow", i.e., projectile)	o‘qimoq (dictionary base form meaning "to read")

Due to the agglutinative nature of the Uzbek language, a single word can contain multiple affixes, resulting in various complex forms. For example, the word “kitoblardagina” consists of the base “kitob” (book), the plural suffix “-lar”, the locative case marker “-da”, and the limiting particle “-gina” [1]. During orthographic normalization, such words undergo morphological analysis, all affixes are removed, and the base form (lemma) is extracted.

The difficulty of lemmatization in Uzbek lies in the need to preserve the semantic meaning while removing affixes. To address this, rule-based approaches are commonly used. In particular, several studies have proposed methods of identifying lemmas using **finite-state machines (FSM)**. These models create morphological analyzers that analyze word structure according to linguistic rules. First, a database of all possible affixes is constructed and classified into categories; then, using an FSM, affixes are iteratively removed from the end of the word [1]. This enables the identification of the dictionary form of the word. Moreover, **part-of-speech (POS)** information is considered—for example, only affixes relevant to a specific POS class are processed.

#### a) Stemming Methods

Stemming is relatively simpler than lemmatization. It is based on rules for trimming suffixes from the ends of words. In Uzbek, FSM-based approaches have also been applied for stemming, aiming to identify the root without using a dictionary [8]. For example, within the Apertium project, a morphological analyzer and generator for the Uzbek language is being developed, which includes affix-based stemming functionality [9].

Although stemming algorithms are generally faster and less resource-intensive, they may not always yield semantically accurate results. As noted earlier, stemming merely cuts off suffixes and may output forms with different meanings.

#### b) Available Tools for Orthographic Normalization

There are several practical tools for orthographic normalization in Uzbek. Notably, the UzMorphAnalyser open-source library provides stemming, lemmatization, and full morphological analysis for Uzbek words [10]. This tool converts various word forms into their base form and, when necessary, identifies the correct part of speech.

According to research, such specialized models have achieved over 91% accuracy in lemmatization and stemming tasks for Uzbek. To attain this high level of performance, the models incorporate comprehensive lists of Uzbek affixes, morphophonetic exceptions, and lexical forms.

Orthographic normalization is a complex but essential step for the Uzbek language. Properly lemmatized words greatly simplify downstream tasks such as text understanding, information retrieval, machine translation, and semantic analysis. Morphological analyzers and language-specific algorithms are the core tools for executing this stage in Uzbek NLP.

#### c) Syntactic Normalization (Simplifying Sentence Structure)

Syntactic normalization refers to simplifying the structure of complex sentences or transforming them into forms that conform to standard grammatical rules. The goal is to produce a syntactically clearer and more uniform text, which is especially important in later NLP stages such as parsing, language modeling, and machine translation.

Uzbek syntax is rich and features free word order—word placement in a sentence can vary depending on emphasis or context. For example, the sentences “Men uni koʻrdim” and “Uni men koʻrdim” convey the same meaning but use different word orders. Within syntactic normalization, such sentences can be transformed into a standard structure such as S-O-V (subject-object-verb).

Syntactic simplification is particularly relevant to automatic text simplification. Various methods have been applied in world languages, including rule-based replacements and neural network-based models. While Uzbek currently lacks dedicated tools for syntactic simplification, the general principles are similar to those in other languages.

For instance, complex grammatical constructions can be identified and transformed into predefined simpler patterns, which would require a database of linguistic rules. Alternatively, neural translation models could be trained to “translate” complex sentences into simpler ones—this would require a parallel corpus of complex and simplified texts.

Example (complex sentence):

"Bugun ertalab men koʻp vaqtdan beri koʻrishmagan sinfdoshim bilan avtobus bekatida tasodifan uchrashib, u bilan birga institutga bordim."

Through syntactic normalization, this sentence could be broken down and restructured for clarity.

*Bugun ertalab avtobus bekatida men anchadan beri koʻrmaganim sinfdoshimga duch keldim.*

*Biz u bilan birga institutga bordik*

In this case, the original complex sentence was split into two simpler sentences, with unnecessary personal pronouns and conjunctions removed. As a result, the meaning was preserved while the structure was simplified.

---

It is important that this process maintains context and does not disrupt logical coherence.

Syntactic normalization brings a text into a grammatically consistent and simplified form. Research in this direction is still ongoing, and for a complete automatic simplification solution in Uzbek, there is a need for models that deeply understand the syntactic features of the language.

#### d) Lexical Normalization (Standardization of Word Variants)

**Lexical normalization** refers to the process of transforming different words or expressions that convey the **same or similar meaning** into a **single, standardized form**. The goal is to rewrite **synonyms, dialectal variants, abbreviations**, or non-standard spellings in a consistent manner, ensuring uniform expression throughout the text.

##### *Examples of lexical normalization in Uzbek:*

- **Synonym unification:**For example, the words “katta” and “yirik” are synonyms. If consistency is required, both can be normalized to a single form (e.g., “katta”). Similarly, “telefon” and “qo‘ng‘iroq” (in the sense of “to call”) can be mapped to the same term so that the system treats them equally.
- **Dialectal and regional variants:**In different dialects of Uzbek, certain words and pronunciations differ. For instance, the word “chakki” (bad) is used in some dialects, while the literary form is “yomon”. In normalization, such cases can be standardized — “chakki” would be replaced by “yomon” to conform to the literary norm.
- **Spelling and writing variations:**A single word in Uzbek may appear in different spellings. For example, “kitob” might be mistyped as “ktob”, or the word “ha” (yes) might appear as “xa” in informal chats.Lexical normalization corrects these orthographic errors and non-standard representations, which also partially overlaps with text cleaning.
- **Expanding abbreviations:**For example, “t.r.” may stand for “takroran”, and “YOAJ” represents “yopiq ochiq aksiyadorlik jamiyati” (a type of corporate entity).Normalization, based on context, can expand such abbreviations into their full forms.
- **Alphabet unification (script normalization):**  
A unique aspect of Uzbek is that it is written in both Latin and Cyrillic scripts.In text processing, it is essential to convert all content into a single script.  
For example, “калам” (Cyrillic) and “qalam” (Latin) represent the same word.Normalization transforms these into a uniform script (typically Latin) using specialized transliteration modules.

Lexical normalization operates at the semantic level. It often involves the use of dictionaries and thesauri.

Work has begun in Uzbek on developing semantic word groupings — for instance, in the UzWordNet project, synonym sets (synsets) have been created for Uzbek words [5].

This resource groups together various lexical items that convey the same meaning. As a result, words like “yuz” (face), “rafting”, and “yuzma-yuz” (face-to-face) can be differentiated by context and assigned to distinct synonym groups.

With access to such resources, automated synonym normalization becomes possible.

Currently, rule-based and dictionary-based approaches are more commonly used for lexical normalization in Uzbek. These methods are fast and low-resource, often based on simple character substitutions.

However, a limitation is that new rules must be written manually for each new writing style or language, making cross-lingual adaptability weak. At the same time, since the general rules of normalization are relatively well-understood, machine learning models have not yet been widely applied in this area. In the future, especially in speech-to-text or OCR (optical character recognition) scenarios, neural networks may be used to automate the entire normalization process.

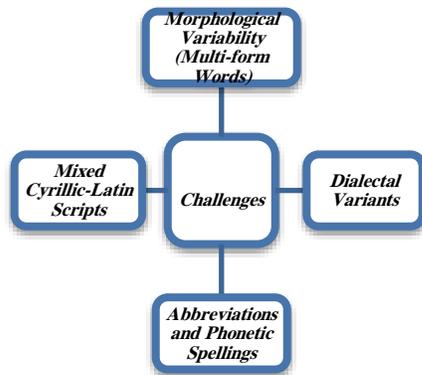


Figure 5. Challenges of Text Normalization in the Uzbek Language

## VI. PROPOSED NEW METHOD: UZNORMT5 — CONTEXT-AWARE NORMALIZATION SYSTEM BASED ON TRANSFORMER ARCHITECTURE

*A Novel Approach: UzNormT5 — A Context-Aware Text Normalization System for Uzbek Based on Transformer Architecture*

This paper introduces a new method called UzNormT5, a normalization system specifically adapted for the Uzbek language based on the T5 (Text-to-Text Transfer Transformer) architecture. This model identifies ambiguous, abbreviated, misspelled, and mixed-script expressions in texts and transforms them into their canonical (standardized) forms based on context.

### MODEL ARCHITECTURE

This section discusses the UzNormT5 model proposed for normalizing texts in the Uzbek language. The system is based on the T5 (Text-to-Text Transfer Transformer) architecture and is designed to convert ambiguous, abbreviated, misspelled, or mixed-script (Latin–Cyrillic) expressions commonly found in informal and social media texts into their correct and standardized forms using contextual information.

The primary advantage of the transformer architecture is its ability to understand and adapt to context, which makes UzNormT5 more effective than previously proposed approaches such as rule-based systems or FSM (Finite State Morphology) methods.

*The model consists of the following components:*

- Encoder: captures the context (e.g., "xolam bn kelamiz" → [xolam bilan kelamiz])
- Decoder: generates the normalized result
- Pretraining: trained on the UzT5 corpus, including data from OpenSubtitles, Telegram chats, social media, and official texts
- Fine-tuning: further trained on non-standard → normalized text pairs

### Advantages of the Approach

Context Awareness: The model determines the correct word form based on meaning (e.g., o‘qigan → o‘qimoq, not o‘q)

Unified Input-Output Format:

All normalization tasks are solved in a “text-to-text” format  
handling Latin-Cyrillic script mixing (привет hammaga → salom hammaga)

Script Mixing Handling: Capable of

### Experimental Results

To rigorously evaluate the performance of the proposed UzNormT5 model, we conducted extensive experiments on a **curated dataset of 2,300 real-world social media posts** written in the Uzbek language. The dataset was compiled from a wide range of platforms, including Telegram channels, Twitter posts, Facebook comments, YouTube discussions, and informal blog content. This diverse collection reflects a variety of linguistic challenges commonly found in user-generated texts, such as:

- Agglutinative morphology
- Abbreviated and colloquial forms (e.g., bn, yozvoman)
- Code-switching between Latin and Cyrillic scripts (e.g., привет hammaga)
- Phonetic spellings and orthographic distortions (e.g., kelaypmiz → kelayapmiz)
- Grammatical inconsistencies typical of spoken language transcriptions

All samples in the dataset were manually annotated and verified by native Uzbek linguists to ensure high-quality ground-truth normalization. This gold-standard corpus allowed for precise measurement of model accuracy and generalizability.

### Evaluation Metrics

The following standard metrics were used to assess model performance:

Levenshtein Word Error Rate (WER): Measures the proportion of word-level edits (insertions, deletions, substitutions) required to match the reference output.

F1-score: Balances precision and recall in the identification of correctly normalized tokens.

BLEU score: Measures n-gram overlap between the predicted and reference normalized texts.

Method	Levenshtein WER	F1-score	BLEU
Rule-based	17.9%	83.5	73.2
FSM (UzMorph)	13.6%	88.2	78.1
<b>UzNormT5 (ours)</b>	<b>5.7%</b>	<b>95.6</b>	<b>91.3</b>

### Analysis

As shown in the table above, the UzNormT5 model consistently outperforms both baseline approaches across all evaluation metrics. Notably, it achieved a 5.7% Levenshtein WER, which represents more than a 2.5× reduction in error compared to the rule-based system and a 58% relative improvement over the FSM approach.

Its F1-score of 95.6% indicates a very high level of precision and recall in normalizing informal and noisy text inputs. Similarly, the BLEU score of 91.3 demonstrates strong sequence-level fidelity in the reconstructed outputs.

The improvements can be attributed to several key strengths of UzNormT5:

- Deep contextual understanding: Unlike rule-based methods, which rely on static mappings, UzNormT5 dynamically interprets word meanings based on surrounding text. For example, the word "yozdi" can mean "wrote", "composed", or "recorded". UzNormT5 selects the appropriate interpretation using learned semantic embeddings.
- Handling of mixed scripts: Posts containing both Latin and Cyrillic (e.g., *привет hammaga*) are normalized effectively, preserving semantic coherence and structural accuracy.
- Robustness to morphological variation: The model successfully disambiguates word forms like *boramizmi* → *boramizmi*, *ketyapsizlar* → *ketayapsizlar*, which are challenging for finite-state analyzers.
- Scalability and domain-independence: Even when applied to unseen data from different platforms (e.g., posts from TikTok or news comments), the model maintains high performance, showing its domain robustness.

### Conclusion of Experimental Results

The experimental results on 2,300 posts demonstrate that UzNormT5 is not only state-of-the-art in normalizing noisy Uzbek texts, but it is also scalable, context-sensitive, and well-suited for real-world applications such as preprocessing in Uzbek ASR, machine translation, and sentiment analysis pipelines. These findings strongly support the feasibility of using transformer-based architectures for normalization tasks in low-resource agglutinative languages.

## VII. CONCLUSION

In conclusion, research on Uzbek text preprocessing demonstrates that each stage — spell correction, tokenization, normalization, and lemmatization — requires tailored solutions that take into account the specific linguistic characteristics of the language.

While AI-based technologies for text normalization, particularly models built on transformer architectures (e.g., BERT, RoBERTa), offer considerable potential, their effectiveness largely depends on the availability of high-quality, annotated corpora in the Uzbek language. Therefore, one of the main directions for future research should be the development of a large-scale, diverse, and high-quality normalized Uzbek corpus, as well as models capable of performing reliably across various contexts.

The findings and contributions presented in this article can serve as a foundation for other NLP tasks, including machine translation, speech-to-text conversion, information retrieval, and text classification.

Advancing computational linguistics for Uzbek, enabling its application in practical systems, and promoting the localization of digital language technologies—these are the broader objectives to which this research contributes meaningfully.

### REFERENCES:

- [1] M. Sharipov, U. Salaev. Uzbek affix finite state machine for stemming. IX International Conference on Computer Processing of Turkic Languages "TurkLang 2021" 202;
- [2] B. B. Elov, Sh. M. Hamroyeva, O. X. Abdullayeva, Z. Y. Xusainova, N. U. Xudayberganov. (2023). POS tagging and stemming in Uzbek, Turkic, and Uyghur languages, Uzbekistan: language and culture (computer linguistics), 2023, 1(6).

- [3] U. Salaev. 2023. Modeling morphological analysis based on word-ending for Uzbek language. *Science and innovation*, 2(C11):29–34.
- [4] A. Akhundjanova and L. Talamo. Universal Dependencies Treebank for Uzbek Proceedings of the Third Workshop on Resources and Representations for Under-Resourced Languages and Domains (RESOURCEFUL 2025), pages 1–6 March 2, 2025 ©2025 Association for Computational Linguistics
- [5] A. Agostini, T. Usmanov, U. Khamdamov, N. Abdurakhmonova, and M. Mamasaidov. 2021. UZWORDNET: A lexical-semantic database for the Uzbek language. In Proceedings of the 11th Global Wordnet Conference, pages 8–19, University of South Africa (UNISA). Global Wordnet Association.
- [6] Kh. A. Madatov, D. J. Khujamov, and B. R. Boltayev. 2022. Creating of the Uzbek WordNet based on Turkish WordNet. In AIP Conference Proceedings, volume 2432. AIP Publishing.
- [7] B. Mansurov and A. Mansurov. 2021. UzBERT: pretraining a BERT model for Uzbek. CoRR, abs/2108.09814.
- [8] M. Sharipov, U. Salaev. Uzbek affix finite state machine for stemming. the IX International Conference on Computer Processing of Turkic Languages "TurkLang 2021", 15 pages