

# Automatic Text Normalization in Uzbek: Problems, Tools, And Solutions

Sobirova Nazira G'anijon qizi

PhD candidate at Alisher Navoiy Tashkent State University of Uzbek Language and Literature, Uzbekistan

**Received:** 23 April 2025; **Accepted:** 19 May 2025; **Published:** 21 June 2025

**Abstract:** In recent years, research in the field of Natural Language Processing (NLP) has increased the demand for automated text analysis across multiple languages, including Uzbek. The multi-form, morphologically complex, and stylistically diverse nature of texts written in Uzbek poses certain challenges for automatic analysis. The central focus of this article is the automatic normalization of Uzbek texts—that is, the process of text normalization. It is dedicated to studying the linguistic and technological issues that arise during automatic text normalization in the Uzbek language. Complex morphological structures, polyform words, dialectal variants, Cyrillic-Latin script differences, and non-standard expressions complicate this process. The results of this research contribute to the deeper digital processing of the Uzbek language and to improving the quality of systems for machine translation, speech-to-text conversion, and text analysis.

**Keywords:** Uzbek language, text normalization, natural language processing, artificial intelligence, neural networks, rule-based approach, morphological analysis, BERT, writing systems, linguistic issues.

**Introduction:** Text normalization is the process of converting various writing styles, spelling errors, dialectal expressions, abbreviations, incorrectly written words, and non-standard phrases commonly found on social media into a standardized, dictionary-compliant form. This process is crucial for the accuracy of subsequent tasks such as text analysis, classification, translation, or speech synthesis. In Uzbek, the differences between Cyrillic and Latin scripts, phonetic writing practices, dialectal variations, and the existence of multiple morphological forms of a single word add further complexity.

In recent years, several scientific studies have focused on the normalization of Uzbek texts, particularly in terms of orthography, syntax, and lexicon.

M. Sharipov and O. Sobirov (2022), in their article, presented an algorithm for affix separation and lemma identification in the Uzbek language using a finite-state automaton. They demonstrated the difference between stemming and lemmatization with a concrete example, emphasizing the importance of identifying the correct root form of a word [1].

B. Elov et al. (2023) compared stemming and POS

tagging across Uzbek, Turkish, and Uyghur languages, discussing the challenges and solutions of implementing stemming in agglutinative languages. They also showed that a hybrid approach—combining rule-based and statistical methods—improves effectiveness [2].

The UzMorphAnalyser model and software, developed by Ulug'bek Salaev (2024), analyzes all possible forms of words in Uzbek. The study compiled a list of all grammatical affixes in Uzbek and developed analysis rules for each. When tested, the model achieved 91% accuracy [3], which is a high result for the Uzbek language. This indicates the strong potential of morphological normalization tools.

While there are few dedicated tools for syntactic normalization in Uzbek, the Uzbek-UT treebank developed within the Universal Dependencies framework (Kurbanova N. et al., 2025) serves as a resource for consistent syntactic annotation. This work highlights challenges in annotating specific syntactic features of Uzbek—for example, compound verb constructions formed with auxiliary verbs [4]. Such treebanks provide a foundational basis for syntactic normalization research.

Additionally, the GREW (graph rewriting) tool proposed by Bruno Guillaume (2021) helps consistently preserve and transform syntactic annotations in corpora, which indirectly supports syntactic normalization (e.g., harmonizing parse trees across languages).

The article by E. Kuriyozov et al. (2021), within the UzWordNet (Uzbek WordNet) project, introduces a lexical-semantic network for the Uzbek language [5]. The resource includes synonymy, antonymy, and hyponymy relations between words. This database provides a scientific basis for lexical normalization tasks such as grouping synonyms and consolidating redundant variants.

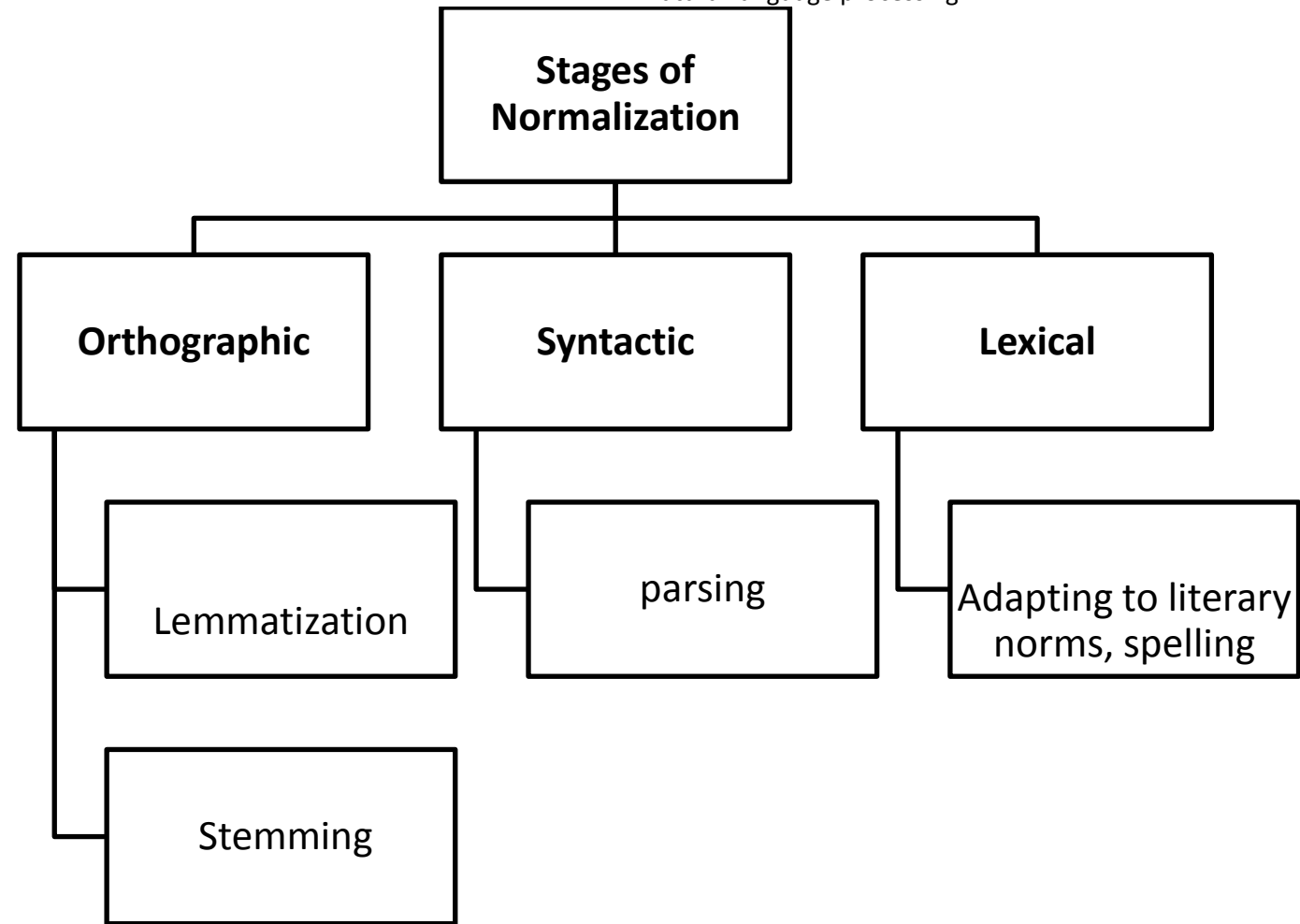
Kh. Madatov et al. (2022) described the process of creating the Uzbek WordNet based on the Turkish WordNet and presented comparative approaches [6].

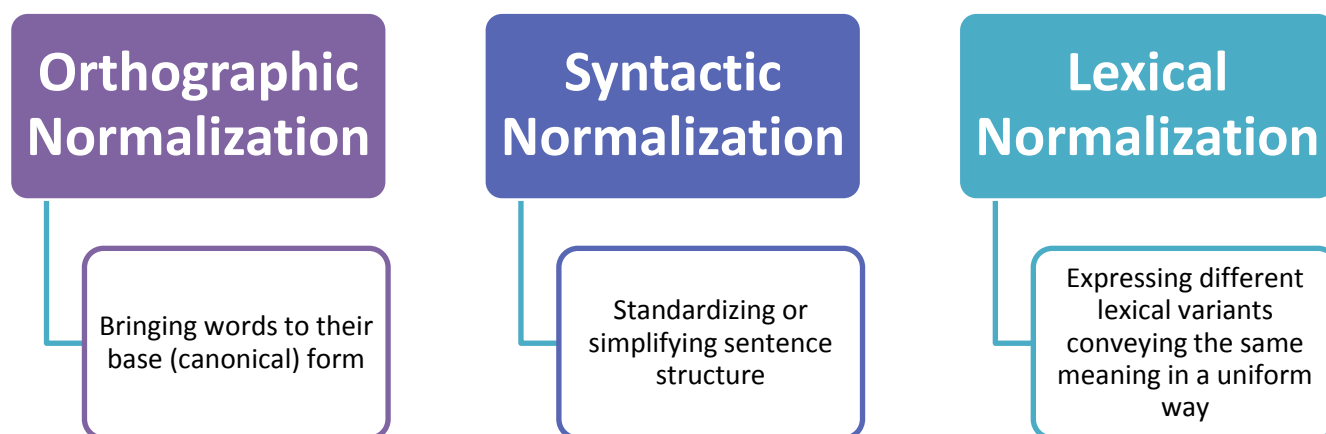
The UzBERT model introduced by B. Mansurov (2021) is one of the first transformer-based models pre-trained on large Uzbek corpora [7].

Additionally, G'. Matlatipov et al. (2022) created a labeled corpus for sentiment analysis in Uzbek. Although these studies are not directly about normalization, they note that during corpus preparation, normalization (e.g., text cleaning, case folding, removing unnecessary characters) was a key step.

These referenced studies demonstrate a growing scientific interest in computational processing of the Uzbek language. In particular, significant progress has been made in morphology and lexicography. Future research is expected to delve deeper into syntactic and semantic analysis. The results of these studies form the foundation for software and practical NLP systems, enabling more efficient normalization and understanding of Uzbek texts.

Text normalization is the process of converting the words and sentences in a text into a uniform, standardized form. Normalization plays a critical role in natural language processing.





### Orthographic Normalization (Lemmatization and Stemming)

Orthographic normalization aims to bring various morphological forms of a word into its main lexical form. This is mainly carried out through the methods of lemmatization (finding the lemma — lexical root of the word) and stemming (removing affixes to extract the root). While lemmatization and stemming share similar goals, they differ in the resulting form and accuracy:

**1. Stemming** – cutting off grammatical affixes from the word to extract an unchanging root part. In this case, the full main form may not be obtained in terms of meaning.

For example, if the word “o‘qigan” is stemmed, the

word form	Stemming	Lemmatization
<b>o‘qigan</b>	<b>o‘q</b> (ma’nosi “yoy o ‘qi”, boshqa ma’no)	<b>o‘qimoq</b> (lug‘aviy asos, “o‘qimoq” – ‘to read’)

### Orthographic Normalization (Lemmatization and Stemming) in the Uzbek Language

Due to the agglutinative nature of the Uzbek language, a single word can include numerous suffixes and appear in various forms. For example, the word “kitoblardagina” consists of the morphemes kitob (root), -lar (plural), -da (locative case), and -gina (restrictive particle) [1]. In the process of orthographic normalization, such words undergo morphological analysis: all affixes are separated, and the root form (lemma) is identified.

The complexity of lemmatization in Uzbek arises from the need to preserve the lexical meaning of a word while stripping away its affixes. For this reason, rule-based approaches are often used. In particular, research has proposed the use of a finite state machine (FSM) to remove affixes and find the lemma. This involves constructing a morphological analyzer based

result is “o‘q”, but in Uzbek, “o‘q” also means “yoy o‘qi” (arrow), which refers to a different meaning.

Thus, stemming does not take semantics into account, it only shortens the form through a technical cutting process.

**2. Lemmatization** – identifying the lemma, i.e., the lexical root of a word. This method standardizes the grammatical forms of a word into a single lexical unit.

In the example above, the lemma of “o‘qigan” is determined as “o‘qimoq”, which is the infinitive (base) form of the word, meaning “o‘qish amali” (to read).

Therefore, lemmatization preserves the meaning of the word while converting it into its base form.

on regular rules: a database of all affixes is created and categorized into groups. Then, starting from the end of the word, the FSM sequentially removes matching affixes [1]. As a result, the lexical base of the word can be determined.

In this process, the part of speech (POS)—such as noun, verb, adjective, etc.—is also taken into account. For example, only affixes relevant to a specific part of speech are considered.

#### Stemming Methods

Stemming is a simpler process compared to lemmatization and is based only on rules for trimming affixes from the end of a word. For stemming in Uzbek, the FSM approach has also been suggested—this aims to identify the root without using a dictionary, by merely removing affixes [8].

For example, under the Apertium project, a

morphological analyzer and generator for the Uzbek language is being developed, which also performs affix-based stemming [9]. Stemming algorithms are typically faster and require fewer resources, but as mentioned earlier, the results may not always preserve the correct lexical meaning.

### Available Tools for Uzbek

Several practical tools exist for orthographic normalization in Uzbek. One example is the open-source library called UzMorphAnalyser, which provides functionalities for stemming, lemmatization, and full morphological analysis of Uzbek words [10]. This tool can normalize various word forms to their base and, if needed, also specify the part of speech.

According to research, using such specialized models has achieved over 91% accuracy in Uzbek lemmatization and stemming tasks. To achieve this level of performance, the model includes a comprehensive list of Uzbek particles and affixes, their morphophonetic exceptions, and a database of lexical forms.

Orthographic normalization in the Uzbek language is complex but essential. Correctly lemmatized words significantly simplify downstream tasks such as text understanding, search, translation, and analysis. Morphological analyzers and algorithms developed specifically for Uzbek are key components in enabling this step.

### Syntactic Normalization (Simplification of Sentence Structure)

Syntactic normalization refers to the simplification of complex sentence structures or transforming them into forms that align with standard grammatical rules. The goal is to make the syntactic composition of the text more understandable and consistent, which is essential for subsequent NLP tasks (e.g., parsing, language modeling, or machine translation).

The syntax of the Uzbek language is rich and allows for flexible word order – words in a sentence can change positions based on emphasis and context. For example, “Men uni ko’rdim” and “Uni men ko’rdim” have the same meaning, but different word orders. Within syntactic normalization, such sentences can be transformed into a unified standard order (e.g., S-O-V – Subject-Object-Verb order).

Syntactic simplification is especially important in tasks like automatic text simplification. Various methods have been applied in world languages in this area, including rule-based transformations and neural network-based models. Although dedicated syntactic simplification tools for Uzbek are still under development, the general principles are similar to

those used in other languages.

For example, complex grammatical constructions can be identified and replaced with predefined simpler patterns (which requires a base of linguistic rules). Alternatively, using neural translation techniques, a model can be trained to “translate” complex sentences into simpler ones – this requires a parallel corpus of simplified text.

### Example:

“Bugun ertalab men ko’p vaqtdan beri ko’rishmagan sinfdoshim bilan avtobus bekatida tasodifan uchrashib, u bilan birga institutga bordim.”

This complex sentence can be normalized and reconstructed as:

- “Bugun ertalab avtobus bekatida men anchadan beri ko’rmaganim sinfdoshimga duch keldim.”
- “Biz u bilan birga institutga bordik.”

Here, the original complex sentence is split into two simple sentences; unnecessary pronouns and conjunctions are removed, and the structure is simplified while preserving the meaning. Of course, it is important during this process to maintain context and logical coherence.

Syntactic normalization brings text into a grammatically consistent and simplified form. Research in this area is ongoing, and fully automated simplification solutions for the Uzbek language require models that deeply understand its syntactic features.

### Lexical Normalization (Unifying Word Variants)

Lexical normalization is the process of bringing different words and expressions that convey the same or similar meaning into a unified, standard form. The goal is to rewrite synonyms, dialectal variants, abbreviations, or non-standard forms uniformly, ensuring consistency throughout the text.

### Examples of lexical normalization in Uzbek:

-  **Unifying synonyms:**

For instance, the words “katta” and “yirik” are synonyms. If consistency is needed in a text, they can be standardized to one form (e.g., using only “katta”). Similarly, “telefon” and “qo’ng’iroq” (in the sense of making a call) – these can be normalized to a single form so that the system interprets them as the same.

-  **Dialect and regional variants:**

Some words and pronunciations differ in various Uzbek dialects. For example, “chakki” (bad) is used in some dialects, while the literary equivalent is “yomon”. During normalization, “chakki” can be replaced with “yomon” to align with standard Uzbek.

- ✓ **Spelling and writing variants:**

The same word may appear in different spellings in Uzbek – for instance, “kitob” might be mistakenly written as “ktob”, or “ha” (yes) may be written as “xa” in chats. Lexical normalization includes correcting such misspellings and informal writing (this overlaps with text cleaning).

- ✓ **Expanding abbreviations:**

For example, “t.r.” stands for “takroran”, or “YOAJ” – “yopiq ochiq aksiyadorlik jamiyati”. The normalization system can replace such abbreviations with their full forms based on context.

- ✓ **Standardizing script (alphabet) usage:**

One characteristic of the Uzbek language is that it is written in two scripts (Latin and Cyrillic). During text processing, all words need to be converted into a single script. For instance, “qalam” (Cyrillic) and “qalam” (Latin) – are essentially the same word. Normalization converts these into one script (e.g., Latin) using special transliteration modules.

Lexical normalization is a form of semantic-level unification. Sometimes, lexical databases and thesauri are used in this process. Work has begun on creating word groupings by meaning in Uzbek – for example, within the UzWordNet project (Uzbek WordNet), synonym sets (synsets) are being developed [5]. This database groups different lexical items expressing the same meaning. As a result, words like “yuz”, “rafting”, and “yuzma-yuz” can be distinguished by meaning and grouped appropriately. If lexical normalization relies on such resources, standardizing synonyms in text can be automated.

## CONCLUSION

In summary, although artificial intelligence technologies for text normalization, especially models based on transformer architectures (such as BERT, RoBERTa, and others), possess great potential, their effectiveness is directly dependent on the availability of a high-quality and sufficiently large annotated corpus in the Uzbek language. Therefore, one of the main directions for future research should be the creation of a large, diverse, and high-quality normalized corpus in Uzbek, as well as the development of models capable of flexible performance in various contexts. The results of this study also serve as a foundation for other NLP systems such as automatic text translation, speech-to-text conversion, information retrieval, and text classification. This research represents an important step in advancing computational linguistics research in Uzbek, applying it in practical systems, and localizing digital language technologies.

## REFERENCES

Sharipov M, Salaev U. Uzbek affix finite state machine for stemming. IX International Conference on Computer Processing of Turkic Languages “TurkLang 2021” 202;

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).

Ulugbek Salaev. 2023. Modeling morphological analysis based on word-ending for Uzbek language. Science and innovation, 2(C11):29–34.

Arofat Akhundjanova and Luigi 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

Alessandro Agostini, Timur Usmanov, Ulugbek Khamdamov, Nilufar Abdurakhmonova, and Mukhammadsaid 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.

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.

B. Mansurov and A. Mansurov. 2021. UzBERT: pretraining a BERT model for Uzbek. CoRR, abs/2108.09814.

Maksud Sharipov, Ulugbek Salaev. Uzbek affix finite state machine for stemming. the IX International Conference on Computer Processing of Turkic Languages “TurkLang 2021”, 15 pages

wiki.apertium.org.

pypi.org