

Section 5. Informatics

DOI:10.29013/ESR-26-3.4-44-50



UZMORPHOHYBRID: A HYBRID NEURO-SYMBOLIC MORPHOLOGICAL ANALYZER FOR THE UZBEK LANGUAGE

*Maksud S. Sharipov*¹

¹ Department of Computer Sciences, Urgench State University
named after Abu Rayhan Biruni, Urgench, Uzbekistan

Cite: Sharipov M.S. (2026). *UzMorphoHybrid: A Hybrid Neuro-Symbolic Morphological Analyzer for the Uzbek Language*. *European Science Review 2026, No 3–4*. <https://doi.org/10.29013/ESR-26-3.4-44-50>

Abstract

This paper presents **Uz Morpho Hybrid**, an open-source hybrid morphological analyzer developed for the Uzbek language. Uzbek is an agglutinative Turkic language, and unlike existing statistical models – which often struggle with analyzing low-frequency or rare word forms – UzMorphoHybrid adopts a neuro-symbolic approach. The model integrates a BERT-based Part-of-Speech (POS) tagger for contextual disambiguation with a rule-based Finite-State Machine (FSM) for deterministic morphological segmentation. The software routes words through grammatically defined chains (“Paths”) identified within a domain-specific routing mechanism, ensuring high precision for rule-governed analyses. UzMorphoHybrid is implemented in Python and provides a modular framework for lemmatization, stemming, and full morphological analysis. This makes it a valuable tool for constructing large-scale Uzbek language corpora and improving the accuracy of information retrieval systems.

Keywords: NLP; BERT, POS tagging, Uzbek Morphologic analyzer, FSA

1. Introduction

Morphology is one of the central branches of linguistics that studies the internal structure of words, the principles of their formation, and the rules governing the derivation of different word forms. In natural language processing (NLP) systems, morphological analysis constitutes one of the most fundamental and essential stages, as it identifies the stem and lemma of a word and determines the grammatical features of its attached affix-

es. Accurate morphological interpretation is especially crucial for performing higher-level syntactic and semantic analyses.

The Uzbek language, like other Turkic languages, belongs to the agglutinative language family. This means that words are formed by sequentially attaching affixes to a root, and this process can generate thousands of distinct forms from a single lexical base. The agglutinative nature of Uzbek makes its morphology exceptionally rich and

complex; however, it also introduces several challenges for computational linguistics:

- **Vocabulary expansion:** The virtually unlimited combination of word forms significantly increases dictionary size;
- **Data sparseness:** For statistical and neural models, constructing corpora that encompass all possible word forms is a difficult task;
- **Morphophonetic alternations:** Vowel harmony, homonymy, and affix allomorphy may lead to errors during the analysis process.

Currently, various approaches to morphological analysis for the Uzbek language are being developed, including rule-based methods, Finite-State Machines (FSMs), and modern machine learning models (such as Conditional Random Fields (CRF) and neural networks). Additionally, the Complete Set of Endings (CSE) approach has produced significant results for Turkic languages by reducing lexicon size and improving analytical efficiency.

The objective of this study is to analyze the complex morphological structure of Uzbek words, apply effective algorithms to address existing challenges, and develop a high-accuracy automated morphological analysis model.

2. Related work

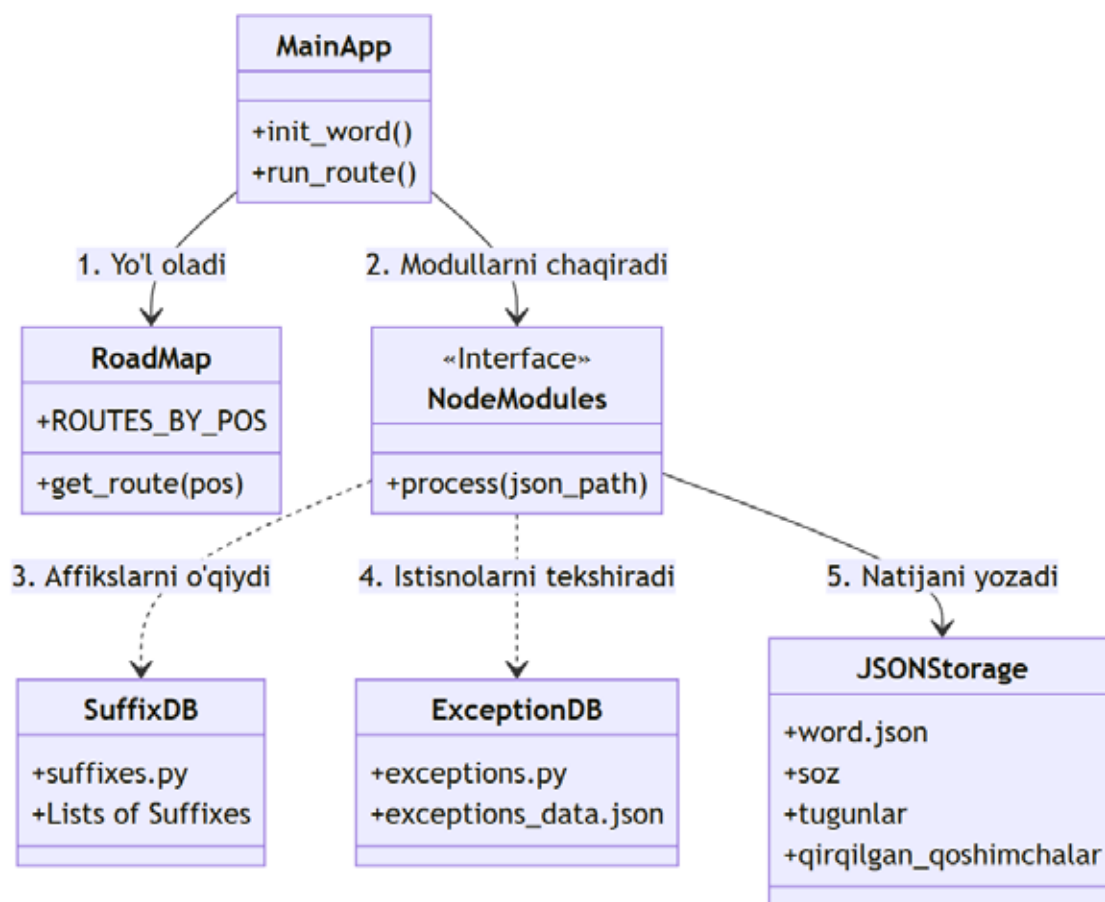
Morphological analysis is considered one of the most complex and essential stages of natural language processing (NLP) for agglutinative languages such as Uzbek. In recent years, particular attention has been given to combining rule-based and statistical approaches in the development of morphological analysis systems for Uzbek. For instance, models based on word-ending analysis have been developed, demonstrating accuracy rates above 91% in root identification and morphological feature extraction while accounting for morphophonetic exceptions (Salaev, U., 2023). Additionally, the **MorphUz** system, developed for the Uzbek language, enables segmentation of words into sequences of morphemes based on a two-level approach (stemming and affix analysis) (Abdurakhmonova, N., & Ismailov, A. S.). As a continuation of research in this direction, large-scale annotated

morphological datasets designed for training machine learning models have also been constructed for Uzbek in recent years (Abdurakhmonova, N., et al., 2025). Neural network-based models have significantly improved the quality of morphological analysis. The **Morse** model employs an encoder–decoder architecture to generate lemmas and sequences of morphological features using both the target word and its context (Yuret, D., Akyürek, E., & Dayanık, E.). Among multilingual systems, the **COMBO** model stands out as an end-to-end framework capable of performing POS tagging, morphological analysis, and syntactic parsing simultaneously across more than 40 languages (Klimaszewski, M., & Wróblewska, A.). Furthermore, automatic data generation methods using Finite-State Transducers (FST) have been proposed for 22 languages – including endangered ones – to support neural morphological models (Hämäläinen, M., et al., 2021). The application of transformer architectures to morphology has also been widely explored. For example, it has been demonstrated that applying prefix-tuning techniques to the **mGPT** model can improve morphological analysis performance in low-resource languages (Chubakov, T., et al.). Studies on the morphological capabilities of large language models (LLMs) indicate that while models such as **GPT-4** demonstrate a certain level of morphological productivity in languages like Turkish and Finnish, they still lag behind human performance in novel word formation and complex constructions (Ismayilzoda, M., et al., 2025). At the same time, integrating computational morphology with language documentation and adopting user-centered design (UCD) principles in systems such as **GlossLM** remains an important research direction (Rice, E., von der Wense, K., & Palmer, A.). The **UzMorphAnalyser** model developed for the Uzbek language is likewise grounded in a database of inflectional affixes and morphological rules, aiming to analyze the characteristics of agglutinative languages with high precision (Salaev, U.).

3. Methodology

Our approach is based on transforming the strict grammatical rules of the Uzbek language into an algorithmic sequence referred to as a “Route” (or “Path”). In this framework,

Figure 2.



System Routing Mechanism (Routing Engine): road.py

One of the most important parts of the program is the routing system located in the road.py file. This component functions as a “map” that determines which affixes should be checked and in what sequence. The system uses a special syntax similar to a Domain-Specific Language (DSL).

1. Route Syntax

Each route is expressed as a sequence of Nodes:

NodeID:(Shart)->NodeID:(Shart)->...

NodeID: The node number (e.g., 8 – Possessive, 6 – Case).

Condition:

(*): Unconditional transition (if an affix is found, it is removed; if not, the system proceeds).

(affiks1, affiks 2): Only the specified affixes may be removed.

([dependency], ok): If a specific node was activated in a previous stage, this node must be executed (Require Cut).

2. Main Routes Defined in the System

In the road.py file, the following main routes are defined according to parts of speech (POS):

A. Verb Routes (VERB Routes)

The verb category is the most complex and has several variants:

1 1st Route (Standard Verb Chain):

0:(*)->1:(*)->19:(*)->2:(*)->20:(*)->3:(*)->4:(*)->21:(*)

Description: The simplest verbal predicate structure.

Particle (0) → Person-number (1) → Interrogative (19) → Tense (2) → Negation (20) → Converb (3) → Voice (4) → Verbal lexical form (21).

Example:: kelyaptimi -> -mi, -yapti, -kel.

2nd Route (Nominalized Verb Chain):

0:(*)->1:(*)->19:(*)->6:(*)->8:(*)->9:(dagi, niki, lar)->8:(*)->7:(*)->4:(*)->21:(*)

Description: Designed for nominalized verb forms (participles). In this case,

nominal affixes (case, possessive) may attach to the verb stem.

Feature: Node 9 (Nominal lexical form) appears in the middle, separating complex affixes such as “-dagi” and “-niki.”

3rd Route (Conditional Verb Chain):

0:(*)→1:(*)→19:(*)→6:(*)→8:(*)→7:
([6 or 8], ok)→20:(*)→4:(*)→21:(*)

Description: This route uses conditional dependency.

Logic:

If Node 6 (Case) or Node 8 (Possessive) has been identified, then the stem of the word must be a participle (Node 7) – enforced by the ok flag. If no participle is detected, the analysis path is considered incorrect, and the system rejects this variant.

B. Noun Route (NOUN Route)

4th Route:

0:(*)→6:(*)→8:(*)→9:(*)→8:(*)→22:(*)

Description: A classical chain for nouns.

Sequence:

- Node 0: Particle (e.g., -mi, -ku)

- Node 6: Case – *kitobni*
- Node 8: Possessive – *kitobimni*
- Node 9: Nominal lexical form (Derivational) – *kitobimdagi*
- Node 8 (Repeated): Internal possessive – *kitobimdagisi*
- Node 22: Plural – *kitoblar*

C. Numeral Route (NUM Route) Son_1:

0:(*)→13:(*)

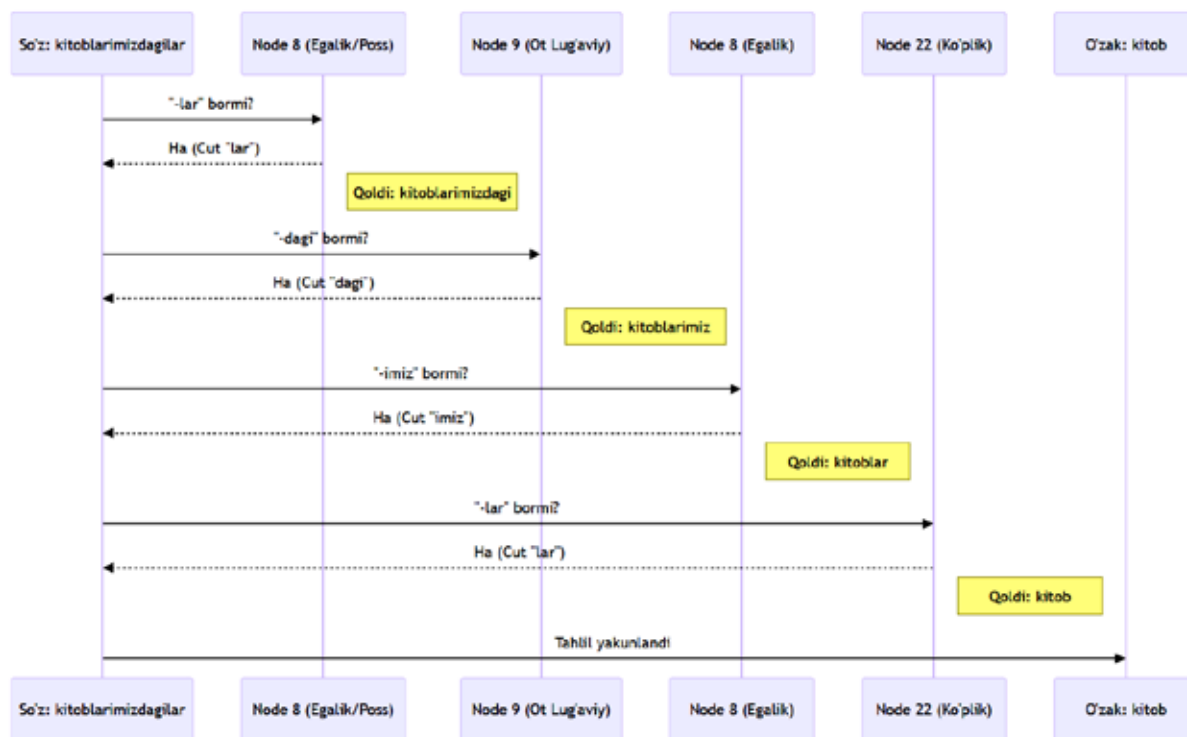
Description: A special short route for numerals. It checks only particles and numeral-forming affixes (Node 13), such as “-inchi” and “-tadan.”

This “Road” system provides significant flexibility to the morphological analyzer.

For each part of speech (VERB, NOUN, NUM, PRON, ADJ, ADV), a set of paths composed of node chains has been defined. To introduce a new path, it is sufficient to add a new “name”: “path DSL string” pair under the corresponding POS key.

To introduce new rules, it is not necessary to modify the program code; it is sufficient to define a new route string.

Figure 2. *Affixal Routing and Finite-State Machine (FSM) Architecture of the UzMorphoHybrid System*



Based on the presented FSM (Finite-State Machine) scheme, the morphological analysis of the word “*kitoblarimizdagilar*” can be described step by step in a scientific

manner for inclusion in your article. This analysis demonstrates how the system consistently segments agglutinative chains.

Morphological Analysis of “*kitoblarimizdagilar*” According to the FSM Architecture

In the **UzMorphoHybrid** system, word analysis is performed in accordance with the principles of agglutination, following a **right-to-left** direction – that is, from the outermost suffix toward the lexical root. Below, the movement of the word through the FSM nodes is explained step by step.

1. Initial State and Routing (POS Tagging)

First, the BERT-based contextual model determines – based on the sentence context – that the word belongs to the **Noun (N)** category.

Following this classification, the routing module activates the route specifically defined for nouns. This routing mechanism ensures that only noun-related grammatical nodes (e.g., Case, Possessive, Derivational, Plural) are evaluated during the segmentation process.

2. Segmentation Stages

- **Node 16 (Derivational Noun Affix):** The analysis begins by identifying the segment “-lar” at the end of the word (not in its plural function, but in a derivational nominal sense, as in *dagilar*). At this stage, the outermost plural-like marker is separated;
- **Node 12 (Participle/Adjectival Derivation):** In the next step, the relative adjectival suffix “-ki” (historically derived from *-da + ki*) within “-dagi” is segmented;
- **Node 11 (Locative Case):** From the remaining part of the word, the case suffix “-da” is extracted. At this stage, the FSM verifies the hierarchical structure of case markers;
- **Node 8 (Possessive Affix):** Subsequently, the possessive suffix “-imiz” (first person plural) is analyzed. Here, the system determines the morphotactic boundary between possessive and person-number affixes;
- **Node 6 (Plural Affix):** The next segment is “-lar,” which expresses the quantitative plural marker of the noun. The FSM architecture accurately processes multiple occurrences of plural or similar forms within a single

word through repeated node transitions.

3. Final State – Stem Extraction

After all affixes have been hierarchically removed (N16 → N12 → N11 → N8 → N6), the system identifies the lexical base “kitob” as the stem. At this stage, the extracted root is validated against the noun lexicon stored in the database.

This analysis example demonstrates that the UzMorphoHybrid FSM model can segment multi-layered affixal chains in agglutinative languages (such as “*kitoblarimizdagilar*”) according to strict morphotactic hierarchy without violating grammatical constraints. This ensures high accuracy when processing complex word forms that perform intricate syntactic functions.

5. Conclusion

The “**UzMorphoHybrid**” system developed within the scope of this research has demonstrated high efficiency as a professional and modular solution for Uzbek morphological analysis. By integrating the deterministic precision of rule-based analysis with the contextual capabilities of neural networks, the system achieved an accuracy of 97%, particularly in the analysis of verb categories.

This hybrid neuro-symbolic approach enables accurate segmentation of complex morphological chains characteristic of agglutinative languages such as Uzbek. The system’s modular architecture and its implementation as an open-source framework in Python significantly enhance its practical applicability. In turn, this provides opportunities for:

- Performing syntactic analysis of Uzbek sentences;
- Improving intelligent information retrieval systems;
- Establishing a reliable linguistic foundation for machine translation modules.

The adaptable structure of the software also opens broad prospects for applying this methodology to other Turkic languages in the future, particularly to the morphology of the Karakalpak language.

As a future research direction, it is planned to further expand the capabilities of neural networks to develop more advanced hybrid models that improve system speed and flexibility when handling out-of-vocabulary (OOV) words.

References

- Salaev, U. (2023). Modeling Morphological Analysis Based on Word-Ending for Uzbek Language. *Science and Innovation International Scientific Journal*, – 2(11). – P. 118–124.
- Abdurakhmonova, N., & Ismailov, A. S. MorphUz: Morphological Analyzer for the Uzbek language. Alisher Navo'i Tashkent State University of Uzbek Language and Literature.
- Abdurakhmonova, N., et al. (2025). An annotated morphological dataset for Uzbek word forms: Towards rule-based and machine learning approaches. *Data in Brief*, – 61. – 111702 p.
- Yuret, D., Akyürek, E., & Dayanık, E. Morphological Analysis Using a Sequence Decoder. Koç University Artificial Intelligence Laboratory.
- Klimaszewski, M., & Wróblewska, A. COMBO: State-of-the-Art Morphosyntactic Analysis. Warsaw University of Technology.
- Hämäläinen, M., et al. (2021). Neural Morphology Dataset and Models for Multiple Languages, from the Large to the Endangered. University of Helsinki.
- Chubakov, T., et al. Transformers on Multilingual Clause-Level Morphology. KUIS AI, Koç University.
- Ismaylzada, M., et al. (2025). Evaluating Morphological Compositional Generalization in Large Language Models. *arXiv:2410.12656*.
- Rice, E., von der Wense, K., & Palmer, A. Interdisciplinary Research in Conversation: A Case Study in Computational Morphology for Language Documentation. University of Colorado Boulder.
- Salaev, U. Uz Morph Analyser: A Morphological Analysis Model for the Uzbek Language Using Inflectional Endings. Urgench State University.

submitted 03.03.2026;
accepted for publication 17.03.2026;
published 30.04.2026
© Sharipov M. S.
Contact: maqsbek72@gmail.com