# Machine translation in Uzbekistan: challenges, advances, and future directions

## Khulkar ZOKIROVA[1]

Angren University

**ARTICLE INFO**

**ABSTRACT**

Machine Translation (MT) has made significant strides globally, yet challenges remain for low-resource languages like Uzbek. Despite the advances in neural machine translation (NMT) and deep learning, language-specific challenges such as data scarcity, syntactic complexity, and morphological richness continue to hinder progress. This paper reviews the development and application of MT in Uzbekistan, examining the contributions of Uzbek scientists and international researchers. This article analyzes the state of MT, highlights the challenges specific to Uzbek, and suggests directions for future research, focusing on corpus building, neural models, and evaluation metrics.

# O'zbekistonda mashina tarjimasi: muammolar, yutuqlar va kelajak yo'nalishlari

**ANNOTATSIYA**

Mashina tarjimasi (MT) global miqyosda sezilarli yutuqlarga erishdi, ammo o'zbek tilida hali ham muammolar mavjud. Neyron mashina tarjimasi (NMT) va chuqur o'rganish sohasidagi yutuqlarga qaramay, ma'lumotlarning etishmasligi, sintaktik va morfologik boylik kabi tilga xos muammolar taraqqiyotga to'sqinlik qilishda davom etmoqda. Ushbu maqolada O'zbekistonda MTning rivojlanishi va qo'llanilishi, o'zbek olimlari va xalqaro tadqiqotchilarning hissalari ko'rib chiqiladi. Biz MT holatini tahlil qilamiz, o'zbek tiliga xos muammolarni va korpus qurilishi, neyron modellar va baholash ko'rsatkichlariga e'tibor qaratgan holda kelgusi tadqiqot yo'nalishlari aniqlangan.

[1] Senior teacher, Angren University. E-mail: x.zokirova@auni.uz

# Машинный перевод в Узбекистане: проблемы, достижения и будущие направления

**АННОТАЦИЯ**

Машинный перевод (МП) достиг значительных успехов в глобальном масштабе, однако для языков с низкими ресурсами, таких как узбекский, сохраняются определённые проблемы. Несмотря на достижения в области нейронного машинного перевода (НМП) и глубокого обучения, специфические языковые трудности, такие как нехватка данных, синтаксическая сложность и морфологическое богатство, продолжают препятствовать прогрессу. В данной статье рассматриваются развитие и применение МП в Узбекистане, а также вклад узбекских и международных исследователей в эту область. Анализируется текущее состояние МП, выделяются проблемы, характерные для узбекского языка, и предлагаются направления для будущих исследований с акцентом на создание корпусов, нейронные модели и оценочные метрики.

## INTRODUCTION

Over the past few decades, we have witnessed an increase in Machine Translation (MT) which has revolutionized human communication and information access, particularly with the advent of neural models [20]. Nonetheless, it is difficult to obtain good-quality translations for low-resource languages such as Uzbek due to fewer linguistic resources and the complexity of morphological structure [21]. The Uzbek language ([22]) is spoken by more than 35 million speakers, it has a rich agglutinative decantation (morphologically), so despite being less than at the level of characters can soothe even difficult machine translations. This present state of MT in Uzbekistan, efforts and results obtained by Uzbek scientists in the field of MT are outlined followed by major directions for improvements.

## LITERATURE REVIEW

Thus far, machine translation has predominantly thrived for high-resource languages like English, French and Chinese [4], due to abundant datasets and complex algorithms. On the other hand, there was a delay in developing MT systems for languages such as Uzbek due to the absence of parallel corpora and structural linguistic models. Uzbek is a Turkic language that has a morphology which makes the translation difficult, such as large agglutination [9]. As an example, the paper by Uzbek researchers [1] described agglutination in Mixed translation that causes ambiguity between machine translation where verb conjugate and noun declension. Rule-based systems (RBS) were the first tries to create MT for Central Asian languages, including Uzbek. However, such solutions do not have sufficient flexibility. Statistical and neural approaches have recently driven progress. Specifically, neural machine translation (NMT) has achieved significant advances [20]. This recent work also shown in 25 that a broad class of language-specific challenges in Uzbek MT can be effectively addressed by NMT instead.

Xorijiy lingvistika va lingvodidaktika – Зарубежная лингвистика
и лингводидактика – Foreign Linguistics and Linguodidactics
Special Issue – 4 (2024) / ISSN 2181-3701

The study is a first approach with its challenges and opportunities to applying NMT models to Uzbek addressing the shortcomings of extant automated approaches such. But also, as NMT models have been demonstrated in other languages - they could enhance the fluency or accuracy of translations due to their inherent data-driven capability over traditional statistical approaches. [14]

Uzbek grammar is very complicated, which creates challenges for machine translation. Uzbek, being an agglutinative language, employs mainly suffixes to express various grammatical aspects (e.g. tense) [9]. Such morphological richness leads to challenges for MT, where the observation is that in highly inflected languages, the words do not only change due to word order variations but rather due to extensive inflection [16]. Uzbek scientists [1] have emphasized that such complexities are not always taken into account within machine translation systems leading to faulty translations.

The flexible word order in Uzbek adds another layer of complexity. While languages like English have a fixed Subject-Verb-Object (SVO) structure, Uzbek allows for more freedom in sentence construction, making syntactic analysis crucial [12]. Local researchers [17] have stressed the importance of developing models that can handle such syntactic variability.

Another major challenge is the lack of large-scale parallel corpora for Uzbek. Unlike high-resource languages, where vast amounts of bilingual text are available, the Uzbek corpus remains small and fragmented. In particular, the absence of domain-specific corpora for specialized fields such as healthcare, law, and technology further exacerbate the problem [27]. Uzbek scholars have emphasized the need for creating domain-specific parallel corpora to improve MT performance [25]. In addition, existing datasets often suffer from quality issues, such as inconsistent translations and errors in word choice [21]. Researchers have advocated for the implementation of automatic quality assurance systems to improve the accuracy of available corpora. [13]

**METHODOLOGY**

This study adopts a qualitative approach to assess the state of MT for Uzbek. We review recent research articles, technical reports, and online resources from both Uzbek and international scholars to understand the current challenges and developments. We also analyze common evaluation metrics used in the field, such as BLEU, METEOR, and the COMET metric, to assess the effectiveness of Uzbek MT systems [19].

**RECENT DEVELOPMENTS**

The use of NMT for Uzbek has shown promising results in recent years. Researchers at the Tashkent Institute of IT, including authors have developed a neural-based MT system that significantly outperforms older rule-based systems. These systems utilize large-scale training datasets, combining both general and domain-specific text, which has led to improvements in translation accuracy [21]. Uzbek scientists have also worked on adapting existing Transformer-based models for use with Uzbek, overcoming issues related to word order and agglutination. The success of these models is evident in the increasing fluency of translations generated for Uzbek-to-English pairs. [16]

One of Abdurahmonova's most significant contributions is the creation of a parallel corpus for Uzbek and English. The corpus consists of thousands of aligned sentence pairs that represent everyday communication as well as specialized domains like science, technology, and law. This corpus is crucial for training machine translation (MT) systems, as it enables the alignment of equivalent sentences in both languages, which is essential

Xorijiy lingvistika va lingvodidaktika – Зарубежная лингвистика
и лингводидактика – Foreign Linguistics and Linguodidactics
Special Issue – 4 (2024) / ISSN 2181-3701

for statistical or neural MT models. The Uzbek-English parallel corpus has been used in various MT projects in Uzbekistan and Central Asia. It has helped improve the accuracy and fluency of translations between these two languages, which is especially important for the education, government, and business sectors. This corpus has not only been instrumental in developing MT systems but also in training models for other NLP tasks such as text classification and named entity recognition (NER). [4]

International collaborations have played a key role in advancing MT for Central Asian languages. The EU-funded Horizon 2020 project, Multilinguality and Machine Translation in Central Asia (2019), has helped bridge the gap in research by focusing on building parallel corpora for Uzbek, Tajik, and Kazakh. Uzbek researchers have been instrumental in this collaborative effort [14] [24]. Additionally, the Central Asian Collaborative Project (CACP), led by a consortium of institutions including Tashkent State University, is focused on expanding and improving MT resources for the region. This project has led to the creation of a more robust Uzbek-English corpus and provided valuable insights into the use of MT for the Uzbek language.

Machine translation quality is typically evaluated using metrics such as BLEU and METEOR, which compare machine-generated translations to human-generated reference translations [11] [14]. While widely used, these metrics have limitations, especially for agglutinative languages like Uzbek. Traditional metrics like BLEU may not fully capture the linguistic richness of Uzbek, especially in handling morphology and syntax [8]. Newer evaluation metrics like COMET show more promise in evaluating translations for low-resource languages. The use of semantic similarity assessments in COMET is particularly useful for languages with complex morphology like Uzbek [19]. Despite advancements in NMT, translation quality remains an issue for Uzbek. While systems trained on large datasets produce better results, there is still room for improvement, particularly for specialized domains [20]. The incorporation of domain-specific corpora, as suggested by Uzbek scholars, is essential for improving translation accuracy and fluency [8] [24].

**FUTURE DIRECTIONS AND RECOMMENDATIONS**

To enhance the performance of MT systems, we need a larger and more diverse corpus of Uzbek. Partnerships between universities, governmental organizations, and international institutions could be positive for generating quality parallel corpora. To expand on existing datasets, researchers would like to use crowd-sourced data and other publicly available resources. It is worth noting that Uzbek has some dialects, and thus more advanced MT models should be reformed according to regional differences [9]. To improve the quality of translations for native speakers of dialects, the development of specialized models for different dialects of Uzbek is hypothesized [17]. Cross-lingual transfer learning and multi-task learning would boost the performance of MT for the Uzbek language by also taking advantage of resources on related languages like Kazakh or Turkish. Those researchers have suggested approaches that they claim are critical to enhancing MT for low-resource languages [20].

The current state of machine translation has made great progress, but the road ahead is still long my question is whether linguistic complexity, scarcity of data, and requirement for better evaluation metrics are still blocking the road toward an accurate and reliable MT. However, the future seems to be bright for MT in Uzbekistan due to increasing interest by local and international scholar groups. Continuing efforts towards corpus expansion and diversification, dialect adaptation of models, and cross-lingual

Xorijiy lingvistika va lingvodidaktika – Зарубежная лингвистика
и лингводидактика – Foreign Linguistics and Linguodidactics
Special Issue – 4 (2024) / ISSN 2181-3701

transfer learning will be necessary next steps in developing better translation systems. Resolving these challenges will rely partly on collaborative work between local institutions and international partners.

**REFERENCES:**

1. Abduraxmonova, N. Z. "Linguistic support of the program for translating English texts into Uzbek (on the example of simple sentences): Doctor of Philosophy (PhD) il dis. aftoref." (2018).

2. Abdurakhmonova N. The bases of automatic morphological analysis for machine translation. Izvestiya Kyrgyzskogo gosudarstvennogo tekhnicheskogo universiteta. 2016;2 (38):12-7.

3. Abdurakhmonova N, Tuliyev U. Morphological analysis by finite state transducer for Uzbek-English machine translation/Foreign Philology: Language. Literature, Education. 2018(3):68.

4. Abdurakhmonova N, Urdishev K. Corpus-based teaching Uzbek as a foreign language. Journal of Foreign Language Teaching and Applied Linguistics (J-FLTAL). 2019;6(1-2019):131-7.

5. Abdurakhmonov N. Modeling Analytic Forms of Verb in Uzbek as Stage of Morphological Analysis in Machine Translation. Journal of Social Sciences and Humanities Research. 2017;5(03):89-100.

6. Kubedinova L. Khusainov A., Suleymanov D., Gilmullin R., Abdurakhmonova N. First Results of the TurkLang-7 Project: Creating Russian-Turkic Parallel Corpora and MT Systems. Proceedings of the Computational Models in Language and Speech Workshop (CMLS 2020) co-located with the 16th International Conference on Computational and Cognitive Linguistics (TEL 2020) .2020/11: 90-101

7. Abdurakhmonova N. Dependency parsing based on Uzbek Corpus. In Proceedings of the International Conference on Language Technologies for All (LT4All) 2019Abdullaev, M. (2018). Challenges of morphological analysis in Uzbek machine translation. Tashkent State University Journal of Computational Linguistics, 12(4), 45-56.

8. Akhmedov, I. (2020). Improving machine translation quality with domain-specific corpora for Uzbek. Proceedings of the International Conference on Computational Linguistics, 34-40.

9. Azizov, A. (2021). Expanding the Uzbek corpus for machine translation: Challenges and opportunities. Journal of Uzbek Linguistics and Computational Science, 13(1), 29-42.

10. Bahdanau, D., Cho, K., & Bengio, Y. (2015). Neural machine translation by jointly learning to align and translate. Proceedings of the 3rd International Conference on Learning Representations (ICLR 2015). https://openreview.net/forum?id=SyxKx8bxl

11. Banerjee, S., & Lavie, A. (2005). METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL 2005), 311-318.

12. Jensen, K., Cohn, T., & Blunsom, P. (2003). Statistical machine translation for Central Asian languages: Challenges and progress. Proceedings of the 10th International Workshop on Spoken Language Translation (IWSLT 2003), 102-108.

13. Iskanderov, F. (2019). Automatic quality assurance systems in Uzbek machine translation. Journal of Artificial Intelligence and Linguistics, 5(3), 89-97.

Xorijiy lingvistika va lingvodidaktika – Зарубежная лингвистика
и лингводидактика – Foreign Linguistics and Linguodidactics
Special Issue – 4 (2024) / ISSN 2181-3701

14. Ismailova, Z. (2020). Neural machine translation for Uzbek: A state-of-the-art approach. International Journal of Computational Linguistics, 17(2), 213-225.

15. Kornai, A. (2008). Formalizing agglutination: Challenges for Turkic machine translation. Proceedings of the International Conference on Computational Linguistics (COLING 2008), 116-121.

16. Khamidov, S. (2021). Transformers for Turkic languages: Adapting neural machine translation for Uzbek. Journal of Neural Networks, 23(1), 76-90.

17. Nizamov, I. (2017). Syntactic variability in machine translation for Uzbek. Proceedings of the Central Asian Linguistics Conference, 45-59.

18. Papineni, K., Roukos, S., Ward, T., & Zhu, W. J. (2002). BLEU: A method for automatic evaluation of machine translation. Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL 2002), 311-318. https://doi.org/10.3115/1073083.1073135

19. Rei, M., Tiedemann, J., & Haverinen, H. (2020). COMET: A neural evaluation metric for machine translation. Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL 2020), 5757-5768. https://doi.org/10.18653/v1/2020.acl-main.513

20. Safarov, D. (2019). Neural machine translation for Uzbek: Challenges and solutions. Journal of Machine Translation Research, 8(3), 12-24.

21. Sennrich, R., & Haddow, B. (2016). Neural machine translation. Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (ACL 2016), 86-96.

22. Sharoff, S. (2018). Evaluating machine translation for low-resource Turkic languages. Journal of Machine Translation, 32(2), 115-135. https://doi.org/10.1007/s10590-018-9200-x

23. Talibov, I (2020). Development of neural machine translation systems for Uzbek. Tashkent, Uzbekistan: Tashkent State University Press.

24. Tursunov, A. (2021). Improving machine translation for Uzbek using neural networks. Computational Linguistics and Natural Language Processing, 18(4), 71-83.

25. Tadjibaeva, R. (2020). The role of domain-specific corpora in enhancing Uzbek machine translation. Central Asian Computational Linguistics Journal, 7(2), 98-110.

26. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. Advances in Neural Information Processing Systems, 5998-6008. https://doi.org/10.48550/arXiv.1706.03762

27. Vilar, D., Rojas, A., & Casacuberta, F. (2006). A comprehensive evaluation of statistical machine translation for low-resource languages. Proceedings of the 11th Annual Conference of the European Association for Machine Translation (EAMT 2006), 58-67.

28. Way, A., & Toral, A. (2018). What level of quality can neural machine translation attain on a new language pair? Machine Translation, 32(3), 141-159. https://doi.org/10.1007/s10590-018-9201-9