

EVALUATION OF MACHINE TRANSLATION IN CENTRAL ASIA: CURRENT CHALLENGES AND EMERGING SOLUTIONS

Khulkar Izzatillayevna Zokirova
Senior teacher of Angren University
x.zokirova@auni.uz

Abstract

Machine Translation (MT) has rapidly advanced in recent years, driven by developments in neural networks and large-scale datasets. However, these advancements have been unevenly applied across languages, with limited focus on the complex linguistic diversity of Central Asia. This study evaluates MT performance for Central Asian languages—Kazakh, Kyrgyz, Uzbek, Tajik, and Turkmen—and identifies key obstacles to achieving accurate and contextually appropriate translations.

Keywords: Machine Translation, Central Asia, BLEU Score, Neural Machine Translation, Linguistic Diversity, Evaluation Metrics, Kazakh, Uzbek.

Machine Translation (MT) systems have become increasingly sophisticated, with substantial improvements in accuracy and fluency, particularly in resource-rich languages like English, Chinese, and German (Vaswani et al., 2017). However, for low-resource languages—such as those spoken in Central Asia—the effectiveness of these technologies remains limited. The Central Asian region, which includes languages like Kazakh, Kyrgyz, Uzbek, Tajik, and Turkmen, presents unique challenges for MT development due to the scarcity of linguistic data, significant morphological complexity, and regional dialectal variations (Akhtyamova et al., 2020; Sennrich & Haddow, 2016).

A core question in MT research for Central Asia concerns the quality of translation output. While widely used evaluation metrics, such as BLEU (Bilingual Evaluation Understudy) and METEOR (Metric for Evaluation of Translation with Explicit ORdering), offer quantitative assessments of translation accuracy, they are often inadequate for languages with complex morphology and flexible syntax, like Kazakh and Uzbek (Papineni et al., 2002; Banerjee & Lavie, 2005). This paper provides an in-depth analysis of MT evaluation in the context of Central Asian languages and explores how MT systems can be adapted to meet the region's linguistic needs more effectively.

Evaluating machine translation typically involves metrics like BLEU and METEOR, which measure the overlap between machine-generated translations and human-generated reference translations. BLEU, introduced by Papineni et al. (2002), is one of the most common metrics used in MT evaluation and relies on matching n-grams between the reference and candidate translations. However, BLEU scores often fail to account for contextual or idiomatic language

use, leading to misleadingly high scores for low-quality translations in morphologically complex languages (Toral & Way, 2018).

METEOR, introduced by Banerjee and Lavie (2005), offers improvements by accounting for synonyms, stemming, and other linguistic features, making it somewhat more effective for languages with varied morphological patterns. However, these traditional metrics still fall short when applied to Central Asian languages, which have high degrees of inflection and syntactic variability (Koehn, 2009).

Recent studies have introduced neural network-based evaluation metrics, such as COMET, which uses cross-lingual embeddings to capture semantic similarities beyond mere lexical overlap (Rei et al., 2020). These metrics hold promise for capturing the nuanced translation needs of Central Asian languages, but their implementation is limited by the lack of sufficient training data for these languages (Conneau et al., 2020).

Central Asian languages are marked by morphological richness, complex syntactic structures, and unique vocabularies, all of which pose challenges for MT. For example, Kazakh and Uzbek, both Turkic languages, display complex case and agreement systems that complicate direct translations into languages with simpler structures (Sharoff, 2018). Additionally, regional dialects within each language add layers of complexity that current MT systems struggle to handle effectively (Ahmad et al., 2021).

Data scarcity exacerbates these issues, as Central Asian languages lack the extensive corpora available for languages like English or Spanish. Akhtyamova et al. (2020) emphasize that while some projects, like the Tatar-Kazakh Parallel Corpus, have aimed to address this data gap, progress remains slow. Without sufficient linguistic data, even advanced MT models like NMT (Neural Machine Translation) and Transformer-based architectures struggle to produce coherent translations for Central Asian languages (Wu et al., 2016).

This study conducted a systematic literature review, drawing from scholarly databases, including IEEE Xplore, ACL Anthology, and Google Scholar, to identify relevant research on MT evaluation metrics, NMT advancements, and the specific challenges facing Central Asian languages. Key terms included “machine translation evaluation,” “Central Asian languages,” “BLEU score,” and “low-resource languages.” Articles were selected based on their relevance to MT evaluation, linguistic challenges, and resource constraints specific to Central Asia.

The reliance on BLEU and METEOR for MT evaluation in Central Asia raises several concerns. BLEU’s emphasis on exact n-gram matching often fails to account for the syntactic flexibility and morphological variations in languages like Kazakh and Uzbek (Toral & Way, 2018). METEOR’s integration of linguistic features offers improvements but remains insufficient for adequately capturing Central Asian linguistic nuances.

Recent studies advocate for neural-based metrics, such as COMET, which rely on cross-lingual embeddings to evaluate MT output more effectively by capturing semantic equivalence rather than lexical similarity (Rei et al., 2020). Although promising, these metrics require robust

training data, which is lacking for most Central Asian languages. Sennrich and Haddow (2016) suggest that back-translation—using MT to generate synthetic parallel corpora from monolingual data—may help mitigate the data shortage, especially for low-resource languages.

Data scarcity presents a fundamental challenge for MT in Central Asia. While efforts like the Kazakh-English corpus and Tatar-Kazakh Parallel Corpus have been instrumental in expanding linguistic resources, they cover only a fraction of what is needed for high-quality MT (Akhtyamova et al., 2020). Techniques such as transfer learning, which allows models to leverage knowledge from related languages, could provide a viable path forward (Conneau et al., 2020). Another significant issue is dialectal variation within Central Asian languages. For instance, Uzbek has distinct dialects that complicate standardized translations (Sharoff, 2018). Dialect-aware models, which can adapt to regional linguistic differences, are necessary for ensuring accurate MT across different Central Asian communities (Ahmad et al., 2021).

The evaluation of machine translation in Central Asia remains an underexplored yet crucial area of research. Traditional metrics like BLEU and METEOR provide a foundation for assessing MT quality, but they fail to capture the full range of linguistic complexities present in Central Asian languages. Neural-based metrics, combined with data augmentation techniques like back-translation and transfer learning, offer promising solutions. However, without continued investment in data resources, dialectal support, and culturally relevant evaluation metrics, MT systems will remain inadequate for meeting the region's translation needs. Future work should focus on developing open-access corpora, enhancing MT models with dialectal sensitivity, and designing metrics tailored to Central Asian linguistic features.

References

1. Ahmad, W., Zhang, P., Ma, X., Peng, N., & Chang, K.-W. (2021). Cross-Lingual and Multilingual Pre-training for Neural Machine Translation. *Transactions of the Association for Computational Linguistics*, 9, 333-345.
2. Akhtyamova, L., Tyers, F., & Washington, J. N. (2020). Building an open-source morphological transducer for Tatar. *Proceedings of the 12th Language Resources and Evaluation Conference*, 1221-1228.
3. Banerjee, S., & Lavie, A. (2005). METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments. *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, 65-72.
4. Conneau, A., Khandelwal, K., Goyal, N., Chaudhary, V., Wenzek, G., Guzmán, F., ... & Stoyanov, V. (2020). Unsupervised Cross-lingual Representation Learning at Scale. *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 8440-8451.

5. Koehn, P. (2009). Statistical Machine Translation. Cambridge University Press.
6. 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, 311-318.
7. Rei, R., Farahani, M., & Thompson, B. (2020). COMET: A Neural Framework for MT Evaluation. Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, 2685-2702.
8. Sennrich, R., & Haddow, B. (2016). Improving Neural Machine Translation Models with Monolingual Data. Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, 86-96.
9. Sharoff, S. (2018). Evaluating MT for Low-Resource Turkic Languages. Journal of Machine Translation, 32(2), 115-135.
10. Toral, A., & Way, A. (2018). What Level of Quality Can Neural Machine Translation Attain on a New Language Pair? Machine Translation, 32(3), 141-159.
11. 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.
12. Wu, Y., Schuster, M., Chen, Z., Le, Q. V., Norouzi, M., Macherey, W., ... & Dean, J. (2016). Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation. arXiv preprint arXiv:1609.08144.