# Review of: "Current Trends in the Use of Machine Learning for Error Correction in Ukrainian Texts"

Juan José Casal[1]

1 University of Buenos Aires

Review of the Paper: "Current Trends in the Use of Machine Learning for Error Correction in Ukrainian Texts"

General Comments:

The paper titled "Current Trends in the Use of Machine Learning for Error Correction in Ukrainian Texts" provides a comprehensive overview of current methodologies, tools, and models used for grammatical error correction (GEC) in Ukrainian texts. The authors have done a commendable job in summarizing the state-of-the-art techniques and their applicability to the Ukrainian language. However, there are several areas where the paper could be improved to provide a clearer and more impactful contribution to the field.

Detailed Analysis:

Clarity and Structure: The paper is generally well-organized, but certain sections could benefit from a more concise presentation. For instance, the sections on various tools and pre-trained models contain repetitive information that could be streamlined to enhance readability. Additionally, a more detailed abstract summarizing key findings and implications would provide better context for the reader.

Methodology: The methodology section outlines the comparative analysis of different approaches to GEC. While the discussion is thorough, the paper would benefit from a more explicit description of the evaluation criteria used to compare the models. For example, detailing the metrics and datasets employed in the comparative studies would strengthen the validity of the conclusions drawn.Comparative Analysis: The comparative analysis of rule-based, syntax-based, statistical, and machine learning methods is informative. However, it lacks a critical discussion of the inherent limitations and potential biases associated with each approach. Including a more nuanced analysis of why certain methods perform better in specific contexts would add depth to the discussion.Data and Corpora: The section on available corpora for Ukrainian is well-researched, highlighting the scarcity of resources compared to English. To enhance this section, it would be beneficial to include a table summarizing the key characteristics of each corpus mentioned. Additionally, the discussion on the need for more annotated data could be expanded to include specific recommendations for future corpus development.

Pre-trained Models: The review of pre-trained models and their application to Ukrainian GEC is thorough. However, the paper could be improved by providing more concrete examples of how these models have been fine-tuned for Ukrainian, including any challenges encountered during this process. Moreover, a comparison of model performance on standardized

benchmarks would provide a clearer picture of their relative effectiveness.

Case Studies and Examples:

Example 1: Neural Machine Translation (NMT) Approach In the study by Cho et al. (2014), the NMT model demonstrated superior performance in grammatical error correction tasks compared to rule-based and statistical methods. This improvement is attributed to NMT's ability to consider the broader context of sentences, making it particularly effective for languages with complex morphologies, such as Ukrainian (Cho et al., 2014). Cho, K., van Merrienboer, B., Gülçehre, Ç., Bougares, F., Schwenk, H., & Bengio, Y. (2014). Learning phrase representations using RNN encoder-decoder for statistical machine translation. arXiv preprint arXiv:1406.1078. https://arxiv.org/abs/1406.1078

Example 2: Seq2Seq Transformer Model The QC-NLP team (2023) reported that their seq2seq Transformer model, fine-tuned on the UA-GEC dataset, outperformed traditional methods. The key innovation was the two-stage fine-tuning process, first on synthetic data and then on "gold" standard data, which significantly enhanced the model's accuracy (Palma Gomez et al., 2023).

Example 3: Rule-Based vs. Machine Learning Models LanguageTool, a rule-based system, was less effective in handling contextual errors compared to the NMT-based approach used by WebSpellChecker (2023). The latter's RedPenNet model leveraged a transformer architecture, allowing it to implement source-to-target transformations efficiently, thus providing more accurate corrections (Didenko & Sameliuk, 2023). Didenko, B., & Sameliuk, A. (2023). RedPenNet for grammatical error correction: Outputs to tokens, attentions to spans. In Proceedings of the Second Ukrainian Natural Language Processing Workshop, Dubrovnik, Croatia. Association for Computational Linguistics. https://www.aclweb.org/anthology/2023.unlp-1.pdf

Results and Evaluation: The results section presents the performance of various models but lacks a detailed explanation of the experimental setup. For instance, the paper should specify the datasets used for training and testing, as well as the evaluation metrics. Including this information would allow for a more transparent assessment of the models' performance.

Future Directions: The conclusion emphasizes the need for further research and development of specialized models for Ukrainian GEC. To strengthen this section, the authors should provide more concrete recommendations for future work. This could include suggestions for specific research directions, potential collaborations, or the development of new tools and resources. For future corpus development, the following recommendations are made to enhance the effectiveness of grammatical error correction (GEC) systems for the Ukrainian language:

- Collect texts from various domains such as formal writing (e.g., academic papers, legal documents), informal writing (e.g., social media posts, blogs), and conversational text (e.g., chat logs, transcripts).
- Include texts from different dialects and regions to ensure the corpus captures the full linguistic diversity of Ukrainian.
- Morphological Annotations: Tags for parts of speech, grammatical number, case, gender, and verb conjugations.
- Syntactic Annotations: Dependency parsing and phrase structure trees to capture sentence syntax.
- Semantic Annotations: Entity recognition, coreference resolution, and semantic role labeling to understand text meaning.

- Error Annotations: Detailed tags for different types of errors (spelling, grammatical, punctuation, and stylistic errors) and their corrections.

- Contextual Annotations: Annotate errors in context, noting not just the error and correction, but also the surrounding sentence to help models understand context-based corrections. Include metadata such as the source of the text, the author's background, and the intended audience, which can influence language use and errors.

- Multimodal Data: Incorporate multimodal data where available (e.g., texts accompanied by images or audio) to support advanced NLP tasks like multimodal learning.

- Quality Control: Implement rigorous quality control procedures, including multiple rounds of annotation by different annotators and consensus mechanisms to resolve disagreements.

- Use professional linguists and native speakers to ensure high-quality annotations.

- Expansion of Existing Corpora: Continuously update and expand existing corpora with new data, ensuring that the datasets remain relevant and comprehensive.

- Create specialized subsets focusing on challenging aspects of the language, such as idiomatic expressions and highly inflected forms.

- Open Access and Collaboration: Make the corpus available as an open-access resource to encourage widespread use and collaboration. Foster partnerships with educational institutions, government agencies, and technology companies to contribute to and utilize the corpus.

Conclusion: The paper makes a significant contribution to the understanding of current trends in machine learning for error correction in Ukrainian texts. The comprehensive literature review and the detailed analysis of existing tools and models are particularly valuable. However, there are areas that require improvement, particularly in the clarity of methodology, comparative analysis, and the presentation of results.

With revisions addressing these points, the paper would be a strong candidate for publication. Its detailed exploration of a relatively under-researched area is commendable, and with some enhancements, it can provide even greater insights and practical recommendations for advancing the field of grammatical error correction in Ukrainian texts.