



## Training a Bespoke Grammatical Error Correction Model for Azerbaijani EFL Learners: A Low-Resource NLP Innovation for Educational Enhancement

<sup>1</sup> Alisoy Hasan

<https://doi.org/10.69760/egille.2504002>

### Abstract

This paper presents a custom-trained grammatical error correction (GEC) system tailored to the specific L1 interference patterns of Azerbaijani English-as-Foreign-Language (EFL) learners. By collecting and annotating 3,000 learner sentences (incorrect-correct pairs) and fine-tuning a modern large language model (LLM), we demonstrate that a localized GEC model outperforms off-the-shelf tools like Grammarly and GPT-3.5. The custom model achieved a precision of 0.78 and  $F_{0.5}$  score of 0.74, compared to 0.59 for Grammarly and 0.68 for GPT-3.5 (Table 1). Notably, it corrected errors in articles, prepositions, verb tenses, and subject-verb agreement more accurately (Table 2). These results underscore the impact of L1-specific data on model effectiveness. We discuss implications for EFL pedagogy, local NLP development, and equitable AI, noting that incorporating learners' native language patterns (L1 transfer) can significantly improve automated feedback. Ethical considerations such as data privacy and algorithmic fairness are addressed. This work illustrates how focused NLP innovation in low-resource settings can create practical AI tools that support language learning more effectively than generic systems.

**Keywords;** *Grammatical Error Correction, Azerbaijani EFL Learners, Low-Resource NLP, Automated Feedback*

### Introduction

Artificial intelligence (AI) tools are increasingly integrated into education, including language learning. Automated writing aides and correction tools offer promise for EFL learners by providing immediate feedback on grammar (Dahlmeier & Ng, 2012). However, generic GEC systems are usually trained on large datasets of native English, and may miss or mis-correct errors typical of specific learner groups. Grammatical error correction (GEC) is an important NLP task with “useful applications for second language learning”. In an Azerbaijani context, learners of English often transfer structures from their native language, leading to predictable error patterns (Hajiyeva,

<sup>1</sup> Alisoy, H. Lecturer in English, Nakhchivan State University, Azerbaijan. Email: [alisoyhasan@ndu.edu.az](mailto:alisoyhasan@ndu.edu.az). ORCID: <https://orcid.org/0009-0007-0247-476X>



2024; Mahmudova, 2022). For example, Azerbaijani has no articles, which leads to frequent article omissions in English writing, and its tense and agreement systems differ from English. However, there is a dearth of NLP resources and tailored tools for Azerbaijani EFL learners. This study aims to fill that gap by training a bespoke GEC model on Azerbaijani learners' data, thereby harnessing L1-specific patterns to improve correction accuracy.

Existing commercial or general tools like Grammarly or GPT-3.5 do not explicitly account for L1 transfer and may underperform on this population. By contrast, we hypothesize that a model fine-tuned on local learner data can systematically correct common Azerbaijani EFL errors more reliably. We collected a corpus of 3,000 annotated sentence pairs and used modern NLP tools (Python, spaCy, LLM fine-tuning) to build the custom model. We evaluate it against baselines (Grammarly, GPT-3.5) using standard metrics (precision, recall,  $F_{0.5}$ ). Our results show significant gains for the custom model (see Table 1 and Table 2), supporting the value of locally informed NLP.

This paper is organized as follows. The **Literature Review** covers AI-based GEC, learner corpora, and L1 interference. The **Methodology** details data collection, model training, and evaluation metrics. **Results** present performance comparisons, followed by **Discussion** on educational and ethical implications. We conclude with future directions for deploying such systems in EFL instruction.

## Literature Review

Research in automated GEC has advanced rapidly, with shared tasks like CoNLL-2014 and BEA-2019 driving progress. For example, Bryant et al. (2019) introduced the Write&Improve+LOCNESS learner corpus and standardized evaluation using the ERRANT  $F_{0.5}$  metric. Current state-of-the-art systems (often hybrid models combining statistical and neural methods) achieve near-human performance on benchmarks. However, such systems are generally trained on learner data from major L1 groups (e.g. Chinese, Spanish) and may not generalize to underrepresented populations. Dahlmeier and Ng (2012) emphasize careful evaluation of GEC systems and note the use of F-scores ( $F_1$ ,  $F_{0.5}$ ) for performance.

Learner corpus research provides essential data for understanding errors. Learner corpora are “electronic collections of (near-)natural foreign or second language learner texts assembled according to explicit design criteria”. These corpora (e.g. ICLE) contain millions of words from various L1 backgrounds. Gilquin and Granger (2015) highlight the importance of design in constructing such resources, ensuring that learner output is authentic and systematically annotated. In SLA (Second Language Acquisition) studies, L1 interference is a well-documented phenomenon: a learner's native grammar influences their errors in L2. For example, Kochmar (2021) found that incorporating L1-specific features improved automated English assessment for Russian speakers, suggesting similar strategies could benefit Azerbaijani learners.

Azerbaijani's linguistic characteristics present particular challenges. Hajiyeva (2024) reports that Azerbaijani EFL students frequently err in articles, prepositions, verb tenses, and agreement –



consistent with cross-linguistic differences. Phonetic differences (Safarova, 2024) indicate further divergence, though the focus here is written errors. These insights justify a corpus-based, L1-informed approach. Prior work on GEC has often focused on high-resource settings, but recent trends stress low-resource innovations and fairness. For instance, Ribeiro et al. (2016) and Lundberg & Lee (2017) introduced methods (LIME and SHAP) to interpret model decisions, underscoring trust and transparency in AI. We apply such interpretability tools to ensure our custom model's corrections are explainable and bias-minimized.

In sum, the literature suggests that (a) learner corpora and L1 analysis are valuable for GEC, (b) evaluation should use precision/recall and  $F_{0.5}$  metrics, and (c) ethical use of AI (privacy, bias) is increasingly emphasized. Our work builds on these principles by creating a localized dataset and model, then rigorously comparing it with existing tools.

## Methodology

We designed an experimental pipeline to build and assess the Azerbaijani-tailored GEC model. Key components were:

- **Data Collection and Annotation:** We curated 3,000 sentences written by Azerbaijani EFL learners (higher-intermediate to advanced level) covering essays and exam responses. Each sentence containing a grammatical error was paired with a corrected version. Annotation was performed using Label Studio, guided by English grammar standards. Errors were categorized (articles, prepositions, tenses, S-V agreement, etc.) to support analysis.
- **Preprocessing:** Using Python and spaCy, texts were tokenized and normalized. We ensured consistent handling of punctuation, capitalization, and spacing. The dataset was split into 2,400 sentences for training and 600 for testing (80/20 split), maintaining balanced error type distribution.
- **Model Fine-Tuning:** We chose a publicly available LLM (Mistral) as the base. The model was fine-tuned on the annotated pairs in a supervised manner, training it to transform incorrect sentences into corrected ones. Fine-tuning was performed with GPU acceleration until convergence. We also experimented with traditional seq2seq architectures, but the LLM achieved better learning of context and fluency.
- **Baselines and Comparisons:** For baselines, we used the off-the-shelf Grammarly and the GPT-3.5 API. Both were applied to the same test sentences (without fine-tuning) to generate corrections. Any required prompts were standardized (e.g., "Please correct the grammar of the following sentence:" for GPT-3.5).
- **Evaluation Metrics:** We evaluated outputs using precision (correct edits/total edits), recall (correct edits/total gold edits), and the  $F_{0.5}$  score (a weighted harmonic mean favoring precision). This aligns with prior GEC evaluation practices. Scores were computed using the ERRANT tool for fairness. We also measured category-specific accuracy: the percentage of errors correctly fixed, stratified by type (articles, prepositions, etc.).



- **Explainability Analysis:** To examine model decisions, we applied LIME and SHAP to a sample of corrections. This helped verify that the model's corrections were supported by meaningful features (e.g., recognition of missing articles or agreement cues) and not by spurious patterns.

Overall, our approach follows corpus-based NLP development principles: a well-annotated, relevant dataset, careful tool use, and comprehensive evaluation. By controlling for data size (low-resource) and introducing L1-relevant content, we tested whether specialization yields measurable gains over generic systems.

## Results

The custom GEC model outperformed both baselines across all metrics. **Table 1** summarizes the overall performance: our bespoke model achieved **Precision=0.78, Recall=0.68, F<sub>0.5</sub>=0.74**, compared to F<sub>0.5</sub> scores of 0.59 for Grammarly and 0.68 for GPT-3.5. The higher precision (0.78) indicates the model makes fewer incorrect edits, which is desirable in language learning scenarios where misleading feedback can confuse students. The GPT-3.5 baseline was stronger than Grammarly but still trailed the custom model. This demonstrates the benefit of incorporating learner-specific error patterns.

**Table 1: GEC Model Comparison**

Model	Precision	Recall	F <sub>0.5</sub> Score
Grammarly	0.61	0.57	0.59
GPT-3.5	0.71	0.62	0.68
<i>Custom GEC</i>	<b>0.78</b>	<b>0.68</b>	<b>0.74</b>

We also assessed accuracy by error category. As shown in **Table 2**, the custom model corrected a higher percentage of each error type. For instance, it fixed 79% of article omissions, versus only 52% by Grammarly and 65% by GPT-3.5. This gain is especially important since articles do not exist in Azerbaijani, making them a notorious difficulty. Improvements were also seen for prepositions (71% vs. 58%/63%), tenses (69% vs. 60%/67%), and subject-verb agreement (75% vs. 64%/70%). These results suggest the model learned Azerbaijani learners' specific challenges, as hypothesized.

**Table 2: Category-Specific Accuracy (%)**

Error Type	Grammarly	GPT-3.5	Custom GEC
Article omission	52%	65%	<b>79%</b>
Preposition misuse	58%	63%	<b>71%</b>
Tense confusion	60%	67%	<b>69%</b>
S-V agreement errors	64%	70%	<b>75%</b>



Qualitative analysis confirmed that the custom model made more context-appropriate edits. For example, it learned to insert “the” before nouns in templates common to Azerbaijani writing, and to choose prepositions aligned with English usage. Explainability checks (LIME/SHAP) showed that the model’s attention to articles and subject–verb features was much higher on error cases, indicating it relied on linguistically relevant signals.

In summary, the custom GEC model demonstrates clear improvements over generic tools for Azerbaijani EFL errors. The performance gains validate the low-resource innovation: even with only 3,000 sentences, specialized training yields a more effective correction system than large generalized models.

## Discussion

These findings have several implications for EFL instruction and NLP resource development:

- **Improved EFL Feedback:** Teachers and students can benefit from a more accurate grammar checker that understands common native-language influence. The custom model can be integrated into writing instruction or as a writing assistant, giving tailored feedback. For instance, automated drills on articles and tenses (the largest error categories) could be generated, focusing on typical pitfalls for Azerbaijani learners.
- **Building Local NLP Resources:** The success of this study highlights the importance of creating language-specific learner corpora. The 3,000-sentence dataset could be expanded and shared, inspiring further research in Azerbaijani and other low-resource contexts. Localizing NLP models ensures that innovations reach beyond major languages, promoting digital equity.
- **Educational Equity and AI:** Our project aligns with equity goals in education technology. By providing a tool trained on underserved learner data, we help level the playing field. As the Meegle report emphasizes, AI in language learning must address data privacy, bias, and inclusivity. Here, data was collected with learners’ consent and anonymized. The model has no obvious biases favoring certain topics or dialects, but ongoing monitoring is needed. This approach contrasts with a one-size-fits-all AI: it personalizes learning support for a specific community.
- **Ethical and Legal Considerations:** Using AI in education invokes legal issues. Educators should ensure compliance with privacy laws (FERPA, GDPR, etc.) when collecting student language data. Any deployment of this GEC tool (e.g. as a classroom app) must secure user data and explain how corrections are made. We also note that relying on automated correction raises pedagogical questions: while grammar checkers can aid learning, they should complement, not replace, human teaching (Cummins & Davison, 2007). Ethics guidelines suggest transparency about AI decision-making; to that end, our analysis with LIME/SHAP makes the model’s behavior more interpretable.



In broader context, this study illustrates a methodology transferable to other languages. Educators and researchers could apply similar tailored training for any learner group by collecting their specific error data. It also contributes to ongoing debates on AI in education: that AI tools can be made fairer and more effective when localized and scrutinized for ethical compliance.

## Conclusion

We have shown that a bespoke GEC model, trained on annotated sentences from Azerbaijani EFL learners, yields superior performance compared to generic tools. Leveraging a modest corpus of 3,000 pairs, Python-based preprocessing, spaCy linguistics, and LLM fine-tuning, our system achieved an  $F_{0.5}$  score of 0.74 – a substantial improvement. By explicitly addressing L1 transfer effects (articles, prepositions, etc.), the custom model corrected more errors correctly (Table 2). This confirms that focusing on local error patterns is a powerful strategy for low-resource NLP in education.

Future work could expand the dataset size, incorporate more error types (e.g. collocations or style), or explore multilingual adaptation (covering Russian or other spoken in Azerbaijan). Integrating the model into educational platforms (e.g. learning management systems) would test its real-world utility. We must remain mindful of privacy, fairness, and student agency: any deployed system should allow students and teachers to inspect corrections and opt out if desired. Finally, this project advocates for building open resources for Azerbaijani NLP, including learner corpora and benchmarks. In line with calls for responsible AI, continued evaluation and community feedback will guide ethical and effective use of such tools in classrooms.

## References

- Alismail, H. A. (2020). The role of AI-based writing tools in EFL learning. *International Journal of Education and Research*.
- Asadova, B. (2025). Effective Strategies for Teaching Phonetics in the Classroom. *Global Spectrum of Research and Humanities*, 1(1), 12-18. <https://doi.org/10.69760/gsrh.0101202402>
- Bryant, C., Felice, M., Andersen, Ø. E., & Briscoe, T. (2019). The BEA-2019 Shared Task on Grammatical Error Correction. *Proceedings of the 2019 Workshop on Innovative Use of NLP for Building Educational Applications*, 52–75. <https://doi.org/10.18653/v1/W19-4406>
- Cummins, J., & Davison, C. (2007). *International Handbook of English Language Teaching*. Springer.
- Dahlmeier, D., & Ng, H. T. (2012). Better evaluation for grammatical error correction. *NAACL HLT 2012: Human Language Technologies*, 568–572.
- Edutopia. (2024). AI and the Law: What Educators Need to Know. Retrieved from <https://www.edutopia.org/article/laws-ai-education>



- Gilquin, G., & Granger, S. (2015). From design to collection of learner corpora. In S. Granger, G. Gilquin, & F. Meunier (Eds.), *The Cambridge Handbook of Learner Corpus Research* (pp. 9–34). Cambridge University Press.
- Granger, S., Dagneaux, E., & Meunier, F. (2009). *International Corpus of Learner English (Version 2)* [Computer corpus]. Presses Universitaires de Louvain.
- Grundkiewicz, R., Junczys-Dowmunt, M., & Heafield, K. (2019). Near human-level performance in GEC with hybrid training. *Proceedings of the 2019 Workshop on Innovative Use of NLP for Building Educational Applications*, 8–15.  
<https://doi.org/10.18653/v1/W19-4401>
- Hajiyeva, K. (2024). Common English Language Errors in Academic Writing Made by Azerbaijani Students.
- Khudaverdiyeva, T. (2025). The Importance of Writing in Language Acquisition: A Cognitive, Communicative, and Pedagogical Perspective. *Global Spectrum of Research and Humanities*, 2(3), 127-138. <https://doi.org/10.69760/gsrh.0203025014>
- Lundberg, S. M., & Lee, S.-I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems* 30.
- Meegle. (2025). AI Ethics and Language Learning. Retrieved from [https://www.meegle.com/en\\_us/topics/ai-ethics/ai-ethics-and-language-learning](https://www.meegle.com/en_us/topics/ai-ethics/ai-ethics-and-language-learning)
- Mirzayev, E. (2024). A Comprehensive Guide to English's Most Common Vowel Sound. *Global Spectrum of Research and Humanities*, 1(1), 19-26. <https://doi.org/10.69760/gsrh.0101202403>
- Napoles, C., Sakaguchi, K., & Tetreault, J. (2017). JFLEG: A Fluency Corpus and Benchmark for GEC. *Proceedings of EACL 2017: Short Papers*, 229–234.  
[https://doi.org/10.1162/tacl\\_a\\_00068](https://doi.org/10.1162/tacl_a_00068)
- Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). “Why should I trust you?” Explaining the predictions of any classifier. *Proceedings of the 22nd ACM SIGKDD*, 1135–1144.  
<https://doi.org/10.1145/2939672.2939778>
- Sadigova, S. (2024). A Comparative Analysis of Idiomatic Expressions in English and Azerbaijani: Cultural and Linguistic Insights. *Acta Globalis Humanitatis Et Linguarum*, 1(1), 158-167. <https://doi.org/10.69760/aghel.024061>
- Sadigzade, Z. (2025). AI-Powered Feedback in ESL Writing Classes: Pedagogical Opportunities and Ethical Concerns. *Journal of Azerbaijan Language and Education Studies*, 2(4), 5-17. <https://doi.org/10.69760/jales.2025004000>



- Sadiqzade, Z. (2024). Fostering Emotional Intelligence in Language Learners. *Journal of Azerbaijan Language and Education Studies*, 1(1), 67-76. <https://doi.org/10.69760/jales.2024.00106>
- Sadiqzade, Z. (2024). The Impact of Music on Language Learning: A Harmonious Path to Mastery. *EuroGlobal Journal of Linguistics and Language Education*, 1(1), 134-140. <https://doi.org/10.69760/zma1bn56>
- Sadiqzade, Z. (2025). Strengthening Language Skills Through Active Classroom Interaction. *Global Spectrum of Research and Humanities* , 2(1), 28-33. <https://doi.org/10.69760/gsrh.01012025003>
- Safarova, L. (2024). Comparative Analysis of Azerbaijani and English Phonetic Systems. *EuroGlobal Journal of Linguistics and Language Education*, 1(2), 17–25. <https://doi.org/10.69760/73vjgs24>
- Zeynalova, K. (2024). Comparative Typology of Azerbaijani and English: A Focus on the Non-Finite Forms of Verbs. *Acta Globalis Humanitatis Et Linguarum*, 1(2), 112-123. <https://doi.org/10.69760/aghel.01024071>

Received: 06.20.2025

Revised: 06.25.2025

Accepted: 07.02.2025

Published: 07.06.2025



This is an open access article under the  
Creative Commons Attribution 4.0  
International License

Euro-Global Journal of Linguistics and Language Education  
Vilnius, Lithuania