

# I’m Sorry, but I Can’t Help with Braille: Revealing Accessibility Failures in State-of-the-Art LLMs

Abdullah Abdullah

orinu Inc.

Hwaseong, Republic of Korea

abdullah.flickdone@gmail.com

## Abstract

Large Language Models (LLMs) perform strongly on many language tasks, but their capability in structurally constrained, accessibility-critical modalities such as Braille remains unclear. We evaluate state-of-the-art LLMs on bidirectional Korean–Braille translation using a human-annotated dataset. Despite expectations that multilingual, instruction-tuned models can generalize to Braille via text representations, we find consistently poor, unstable outputs and substantial disagreement with human judgments. These results point to missing Braille-aware tokenization and weak alignment between Korean and Braille patterns. In contrast, supervised fine-tuning of a small model (T5-small) on the same data yields large and stable gains over zero-shot and prompted LLM baselines across standard metrics (SacreBLEU, ChrF++, CER, BLEU, ROUGE-L, METEOR, CIDEr). Our findings reveal a systematic limitation of current LLMs and demonstrate the effectiveness of modest task-specific supervision.

## 1 Introduction

LLMs (Achiam et al., 2023; Team et al., 2023; Yoo et al., 2024; Anthropic, 2025) have demonstrated strong performance in a wide range of natural language generation and understanding tasks, including machine translation, summarization, and reasoning (Brown et al., 2020; Chowdhery et al., 2023; Touvron et al., 2023). As these models scale, they are often assumed to generalize broadly across languages, scripts, and modalities. However, recent work has shown that such generalization remains uneven, particularly for low-resource languages, non-standard scripts, and accessibility-related representations (Joshi et al., 2020; Blasi et al., 2022).

Braille is a critical written modality for blind and visually impaired users, yet it remains largely overlooked in NLP research. Unlike standard text

translation, Braille conversion is highly character-sensitive, rule-governed, and language-specific, with strict conventions governing contractions, numerals, symbols, and spacing. These properties make it challenging for general-purpose LLMs, which are rarely exposed to Braille during pretraining. Although recent work has begun exploring Braille modeling in other languages (Huang et al., 2025), differences in linguistic structure and Braille conventions limit direct transfer to Korean Grade 2 Braille.

We investigate whether state-of-the-art LLMs meaningfully support Korean–Braille translation under the official Korean-Braille regulations. Using a large human-annotated parallel corpus, we evaluate both Korean-to–Braille and Braille-to–Korean directions. LLMs frequently produce refusals, hallucinations, or invalid outputs, revealing a systematic blind spot in accessibility-critical settings.

To address this gap, we introduce **BT5**, a lightweight Braille-aware model based on T5 (Raffel et al., 2020). With straightforward supervised fine-tuning on expert-annotated data, BT5 substantially outperforms zero-shot and prompted LLM baselines across character-level and generation-based metrics.

Our contributions are threefold: (1) the first systematic evaluation of LLMs on Korean–Braille translation, (2) evidence that small task-specific models can surpass much larger general-purpose LLMs with proper supervision, and (3) identification of Braille processing as an essential yet under-explored direction for inclusive NLP.

## 2 Methods

### 2.1 Dataset

We evaluate both Braille-to–Korean and Korean-to–Braille translation using the NIKL Korean Print–Braille Parallel Corpus 2023 (v1.0) (National Institute of Korean Language, 2024) as a







## Limitations

This study focuses on Korean–Braille translation, and the findings may not directly generalize to other languages or Braille systems with different linguistic structures, contraction rules, and encoding conventions. While recent work has explored Braille modeling in other languages, cross-lingual transfer is non-trivial and was not investigated in this work.

Our evaluation includes a set of representative state-of-the-art LLMs and a Korean LLM available at the time of experimentation; however, it does not exhaustively cover all possible models, architectures, or prompting strategies. In particular, many proprietary LLMs exhibited refusal behaviors, empty outputs, or malformed generations when prompted for Braille, which prevented consistent large-scale quantitative evaluation. As a result, comparisons with these systems are based on controlled samples and qualitative analysis rather than full test-set benchmarking.

Additionally, differences in tokenizer design and pretraining data introduce inherent disparities between BT5 and general-purpose LLMs. Although BT5 benefits from explicit exposure to Braille through supervised fine-tuning and a dedicated tokenizer, most LLMs lack Braille-aware tokenization and aligned training data, making the direct comparison imperfect. Our results should therefore be interpreted as highlighting capability gaps rather than as strictly controlled architectural comparisons.

Although we employ standard automatic metrics (e.g., BLEU, ChrF++, CER, ROUGE), these metrics primarily capture surface-level similarity and may not fully reflect functional usability, readability, or correctness under official Braille standards. Human-centered evaluation with Braille users was beyond the scope of this work but is essential for real-world validation.

Furthermore, our approach relies on supervised fine-tuning with human-annotated parallel data, which may be costly or unavailable in other low-resource settings. We do not explore data augmentation, semi-supervised learning, or cross-lingual transfer, which could improve scalability.

Finally, errors in Braille translation can have significant real-world consequences in accessibility-critical contexts. Consequently, we do not claim that any evaluated model is suitable for direct deployment without rigorous validation, robustness

testing, and adherence to official Braille standards.

## Acknowledgments

This work was supported by the Starting Growth Technological R&D Program (RS-2025-25465816) funded by the Ministry of SMEs and Startups (MSS, Korea). The author also thanks the team members at orinu Inc. for their support with the research environment, project coordination, and data infrastructure during this work.

## Ethical Considerations

This study uses the publicly available NIKL Korean Print–Braille Parallel Corpus 2023 (v1.0), released by the National Institute of Korean Language. The dataset contains written Korean text and corresponding Korean–Braille transcriptions constructed according to the Korean–Braille Regulations (2024) and does not include personally identifiable information. Experiments evaluate existing language models using this dataset and standard evaluation metrics. We report dataset sources, preprocessing steps, data splits, model configurations, fine-tuning procedures, and evaluation protocols to support independent replication, although minor variability may arise from differences in hardware, software versions, or model access. Multiple state-of-the-art LLMs were evaluated as experimental subjects for Korean-to–Braille and Braille-to–Korean translation using publicly available APIs and were treated as black-box systems.

## Use of Generative AI Tools

Generative AI tools were used only for language editing and clarity improvements in this manuscript. All experimental design, data preparation, model training, evaluation, and analysis were conducted by the authors, and no generative AI system was used to generate datasets, annotations, or experimental results.

## References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Anthropic. 2025. *Claude opus 4.5*. Model release. Accessed January 16, 2025.

- Satanjeev Banerjee and Alon Lavie. 2005. **METEOR: An automatic metric for MT evaluation with improved correlation with human judgments**. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.
- Damian Blasi, Antonios Anastasopoulos, and Graham Neubig. 2022. Systematic inequalities in language technology performance across the world’s languages. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5486–5505.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, and 1 others. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113.
- Tianyuan Huang, Zepeng Zhu, Hangdi Xing, Zirui Shao, Zhi Yu, Chaoxiong Yang, Jiaxian He, Xiaozhong Liu, and Jiajun Bu. 2025. Braillem: Braille instruction tuning with large language models for braille domain tasks. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 28589–28600.
- Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. 2020. The state and fate of linguistic diversity and inclusion in the nlp world. *arXiv preprint arXiv:2004.09095*.
- Liblouis Developers. 2024. Liblouis: Open-source braille translation software. <https://liblouis.io/>. Version used in this work.
- Chin-Yew Lin. 2004. **ROUGE: A package for automatic evaluation of summaries**. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- National Institute of Korean Language. 2024. Nikl korean–korean braille parallel corpus 2023 (v1.0). <https://kli.korean.go.kr/corpus>.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.
- Maja Popović. 2017. chrf++: words helping character n-grams. In *Proceedings of the second conference on machine translation*, pages 612–618.
- Matt Post. 2018. **A call for clarity in reporting BLEU scores**. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, and 1 others. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, and 1 others. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. 2015. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4566–4575.
- Kang Min Yoo, Jaegeun Han, Sookyo In, Heewon Jeon, Jisu Jeong, Jaewook Kang, Hyunwook Kim, Kyung-Min Kim, Munhyong Kim, Sungju Kim, and 1 others. 2024. Hyperclova x technical report. *arXiv preprint arXiv:2404.01954*.



### Task: Korean → Braille

#### Korean:

양교는 미래혁신위원회 활동을 통해 △교원 교류 △학생 및 학점 교류 △온라인 공동교과목 개설 △비교과 공동프로그램 개설 △해외교류 프로그램 참가 △인공지능 (AI), 데이터과학 (DS), 의생명분야 공동연구 수행 △해외대학과의 공동연구 등에 상호 협력하게 된다.

#### Liblouis:

양교는 미래혁신위원회 활동을 통해 △교원 교류 △학생 및 학점 교류 △온라인  
공동교과목 개설 △비교과 공동프로그램 개설 △해외교류 프로그램 참가 △인공지능  
(AI), 데이터과학 (DS), 의생명분야 공동연구 수행 △해외대학과의 공동연구 등에  
상호 협력하게 된다.

#### Ours:

양교는 미래혁신위원회 활동을 통해 △교원 교류 △학생 및 학점 교류 △온라인  
공동교과목 개설 △비교과 공동프로그램 개설 △해외교류 프로그램 참가 △인공지능  
(AI), 데이터과학 (DS), 의생명분야 공동연구 수행 △해외대학과의 공동연구 등에  
상호 협력하게 된다.

#### REFERENCE:

양교는 미래혁신위원회 활동을 통해 △교원 교류 △학생 및 학점 교류 △온라인  
공동교과목 개설 △비교과 공동프로그램 개설 △해외교류 프로그램 참가 △인공지능  
(AI), 데이터과학 (DS), 의생명분야 공동연구 수행 △해외대학과의 공동연구 등에  
상호 협력하게 된다.

Figure 5: Example outputs for the Korean-to-Braille task for the rule-based Liblouis system and the BT5 model. Incorrect or mismatched segments are highlighted in red. Liblouis produces structurally inconsistent Braille sequences due to rigid rule application, whereas BT5 closely matches the reference output.