
Lost in Translation: Benchmarking Commercial Machine Translation Models for Dyslexic-Style Text

Gregory Price

School of Electrical Engineering and Computer Science
University of Ottawa
gpric024@uottawa.ca

Shaomei Wu

AImpower.org
Mountain View, CA 94043
shaomei@aimpower.org

Abstract

1 Dyslexia is a neurodivergence that impacts one’s ability to process and produce
2 textual information. While previous research has identified unique patterns in
3 the writings of people with dyslexia - such as letter swapping and homophone
4 confusion - that differ themselves from the text typically used in the training and
5 evaluation of common natural language processing (NLP) systems such as machine
6 translation (MT), it is unclear how current state-of-the-art NLP systems perform
7 for users with dyslexia. In this work, we explore this topic through a systematic
8 audit of the performance of commercial MT services using synthetic dyslexia data.
9 By injecting common dyslexia-style writing errors into popular benchmarking
10 datasets, we benchmark the performance of three commercial MT services and one
11 large language model (LLM) with various types and quantities of dyslexia-style
12 errors and show a substantial disparity in MT quality for dyslexic and non-dyslexic
13 text. While people with dyslexia often rely on modern NLP tools as assistive
14 technologies, our results shed light on the fairness challenges experienced by this
15 demographic with popular NLP services, highlighting the need to develop more
16 inclusive and equitable NLP models for users with diverse language use patterns.

17 **1 Introduction**

18 Dyslexia is one of the most common learning disabilities, estimated to affect 10% to 17% of English
19 speaking population [33, 3]. As a neuro-cognitive condition with no known cure, dyslexia impacts
20 one’s ability to process and produce textual information [29, 28], and can lead to long-term social,
21 emotional, and economic challenges such as less peer acceptance, poor self-image, lower educational
22 attainment, and reduced employment opportunities [11, 27].

23 The rapid development and adoption of neural language technologies - such as the ChatGPT - makes
24 them an important part of the information ecosystem and a promising assistive tool for people
25 with dyslexia [35, 8]. However, most of existing neural language models have been developed and
26 evaluated over typical text (e.g. WikiText [15], CommonCrawl¹), with little consideration of dyslexia
27 use case. The fairness and inclusivity of neural language technologies for users with dyslexia is thus
28 largely underexplored.

29 In this paper, we present an evaluation of the current state-of-the-art machine translation (MT)
30 models available via popular cloud services on dyslexia-style text. To evaluate potential biases

¹<https://commoncrawl.org/overview>

31 presented in machine translation against dyslexia text, we perturbed the source text in WMT14
32 (en2fr) dataset [2] with synthetic dyslexia style writing errors, and benchmark the performance of
33 four commercial machine translation systems using the perturbed data. Our results show all audited
34 models - including advanced LLMs - struggle with dyslexia-style input text, making substantially
35 more lexical and semantic mistakes. By varying the quantity and the types of dyslexia style errors
36 injected into the original text, we also observe a near linear relationship between the amount of
37 dyslexia errors and the decrease in performance for all services, especially for real-word errors such
38 as the confusion of homophones [23, 24]. Our contribution to language technology, AI fairness,
39 and accessibility research is two-fold: 1) Our findings uncover the disparities in the performance of
40 commercial machine translation systems to translate dyslexia style text; 2) Our systematical approach
41 in generating synthetic dyslexia datasets provides an useful instrument to further investigate the
42 potential sources and mechanism for such disparities in typically “black-boxed” systems when real
43 dyslexia datasets are scarce. As an early exploration in AI fairness and dyslexia, our work invites
44 further investment and urgent attention from NLP researchers and commercial companies to develop
45 inclusive and fair NLP models with people with dyslexia, a community deeply impacted by and
46 highly experienced with language technologies.

47 **2 Background and Related Work**

48 **2.1 Dyslexic Writing Style**

49 There are many spellcheckers available that try to correct spelling errors but most are not specifically
50 designed to address dyslexic-style writing. General use spell-checkers perform poorly when it comes
51 to real-word errors [18] (e.g. form v.s. from). [25] found that this comprises of 17% of the errors
52 made by English dyslexic people. There have been some efforts to create a dyslexia-style writing
53 support tool from [21], [18] and [25]. Unfortunately, these systems are designed in an academic
54 fashion and are not the most appropriate for a widespread writing style that is used in an everyday
55 life. Previous work from [35] utilized data and writing from social media to give more relevance to
56 everyday text. More recently, Goodman et al. [8] utilized a Large Language Model (LLM) to create
57 an email-writing interface tool for users with dyslexia. In this study, we do not try to “correct” any
58 dyslexic-style typographical errors, but to understand the capacity of commercial machine translation
59 systems at handling text that contains this style of writing.

60 Dyslexia-style text has been categorized in previous works from [22]. The typographical errors
61 presented were broken down into four categories. *Substitutions* were identified as letters that are
62 changed with one another (reelly v. really). *Insertions* were counted where a letter is inserted
63 (situation v. situatation) or where a word that was incorrectly split (sub marine v. submarine).
64 *Deletions* was when a letter is omitted (approch v. approach). *Transpositions* were considered as two
65 letters that were swapped and adjacent (articile v. article). Using these categories [23] found that the
66 substitutions were by far the most common type of dyslexic-style typographical error. It is to note
67 that this was found on a Spanish corpus of hand written text by students with dyslexia.

68 Work from [19] created a large confusion set of words. This set consists of real word errors collected
69 from dyslexic text [18] and also synthetically created samples that were used to test spell checkers.
70 The set is a mix of homophones, substitutions, insertions, omissions (deletions) and transpositions.
71 This is the most exhaustive set of dyslexic related errors that we were able to find. Previous work
72 utilized synthetic dyslexia writing for neural translation models and demonstrated success in creating
73 an assistive writing tool for dyslexic users on social media [35].

74 **2.2 Subgroup Performance Disparities in AI Systems**

75 Previous work from [4] and [5] has brought to light racial disparities in AI. Inequalities are often
76 caused by lack of awareness in training data, fairness in training [31, 16] and other inclusive consid-
77 erations. They found that people of the minority classes are the ones who suffer from shortcomings

78 of the machine learning models. Lack of data of different groups leads to the use of synthetic data
79 like in [12] for stutter data.

80 Object-recognition systems displayed disparities in terms of income levels and geographies [7,
81 9]. Smaller subgroups are more at risk for poorer performances. Work from [10] identified key
82 components (texture, occlusion and darker lighting) that lead to performance degradation of object-
83 recognition systems in lower income levels/geographical areas and show that it is possible to mitigate
84 these disparities.

85 Work from [30] spotlights the issue of the utilization of the English language when training models
86 creating disparities. These disparities range from people being unable to utilize the models due to a
87 language barrier to the models existing but not performing to par. Inequalities for resources, variation
88 and performances is seen as the industry norm when we apply NLP to underrepresented communities.

89 Unfortunately, there has not been much work researching the affects of artificial intelligence on
90 people with dyslexia. Researcher from [1] were able to gather an estimate of the amount of dyslexic
91 text documents on the Web (0.005%). They deemed their estimate much lower than the corresponding
92 number of dyslexic users (10-17%). If they considered spelling errors as dyslexic-style typographical
93 errors, the number would increase to 0.2%. Therefor, it is likely that models trained on data from the
94 Web is not reflective of the dyslexic population. This leads us to believe that the models are trained
95 on "perfect" data that has been filtered through spell checkers. This potential bias is what is being
96 studied in this paper for one NLP task.

97 **2.3 NLP Model Evaluation and Benchmarking**

98 Machine translation is a common NLP task where a source sentence is translated into a different
99 target language. The *Machine Translation Foundation*² provides a new dataset yearly during the
100 Conference on Machine Translation (WMT) to benchmark the performance of SOTA MT models on
101 various translation tasks. Many language pairs with parallel data are provided in the WMT datasets,
102 with public available source data and manually translated target references. We use WMT14 (en2fr)
103 dataset for this study. It contains news articles in English as source data, together with parallel manual
104 translation in French.

105 Following the breakthrough by Vaswani et al. [34], the transformer architecture has become increas-
106 ingly popular for MT models. We assume that widely used translation services from major cloud
107 service providers such as AWS, Google Cloud and Azure are utilizing this architecture. However,
108 the exact model structure is not public information nor the data that is used in the training for these
109 models. That means to understand and diagnose these systems, we have to rely on their APIs and
110 translation outputs of a wide range of source sentences to shed light on the black box models.

111 **3 Method**

112 For our scope of work, we leveraged and modified the WMT14 (en2fr) [2] dataset to evaluate a
113 machine translation task from English to French with injected synthetic dyslexic-style errors. We
114 select machine translation for our exploratory evaluation because the task is well-defined, with well-
115 established metrics and benchmarking datasets, as well as many popular consumer-facing applications
116 and services such as Google Translate³. We also limit our initial benchmarking to the translation
117 from English to French - two well-resourced languages for machine learning, to reduce potential
118 confounding factors due to languages. In this section, we review how we created the synthetic
119 dyslexic text corpora and the types of dyslexic writing errors injected. We then present and discuss
120 the commercial machine translation services we evaluated using the synthetic dyslexic text. Finally,
121 we describe the metrics and methods we utilized for benchmarking the performance of these services
122 in both lexical and semantic dimensions.

²<https://machinetranslate.org/about>

³<https://translate.google.com/>

123 3.1 Simulating Dyslexia

124 The lack of large scale and publicly available dyslexic text corpus has been a bottleneck for dyslexia-
125 related language technologies today [35, 8]. Direct collection of text written by people with dyslexia
126 faces both ethical and practical challenges. As an “invisible” disability that is highly stigmatized,
127 many people with dyslexia feel the pressure and need to conceal their dyslexia, spending extra
128 efforts to proofreading their writing or avoiding to write at all [26]. Even if people with dyslexia
129 consent to share their data, it is difficult to fully anonymize the data while preserving the unique
130 and personal writing styles of dyslexia. Encouraged by the success of using synthetic disability data
131 for data-intensive machine learning tasks [12, 35], we created a synthetic dataset of dyslexic text
132 by injecting typical dyslexic writing errors into a popular MT benchmarking dataset, namely, the
133 WMT14 (en2fr) test dataset [2]. Taking a similar approach proposed by Wu et al. [35], we perturbed
134 the English source sentences with the following three synthetic errors that are frequent in dyslexic
135 input text and less likely to be fixed by mainstream spellcheckers:

- 136 1. Letter confusion: substituting similar-looking or sounding letters (e.g. b v.s p). Letter
137 confusion is reported as the most frequently occurred errors in dyslexic writing [23].
- 138 2. Homophone: replacing a word with its homophones. Phonetically similar sounding words
139 are noted as another common but unique challenge for people with dyslexia [18], [23], and
140 can potentially create issues for NLP models as this type of error is relatively rare in typical
141 text used to train the models.
- 142 3. Confusion set: substituting a word with another word that are likely to be confused with
143 by people with dyslexia (e.g. “your” and “you”). Previous work found confusion sets
144 contribute a substantial percentage of dyslexic writing errors and are least likely to be caught
145 by conventional spellcheckers [18, 24, 35].

146 To simulate letter confusion, we constructed a letter substitution dictionary in which each letter is
147 associated with other letters people with dyslexia are often confused with [23]. The frequency of
148 letter confusion is controlled by a parameter p_l , which represents the probability for letter confusion
149 to occur in the original corpus. However, following empirical findings that letter confusion rarely
150 occur at the beginning of a word [36, 20, 18], the substitution of the first letter would ignored
151 95% of the time during error injection. Also, to be consistent with the observations that multiple
152 letter confusions are uncommon in dyslexic writing [23], we decreased the probability of another
153 substitution happening by 90% for that same word after one substitution is made.

154 To simulate homophone errors, we constructed a homophone dictionary in which each word is
155 associated with its phonetically similar sounding words. We leveraged free public resources such as
156 the Homophone Finder website⁴ to build the homophone dictionary. The frequency of homophone
157 error is again controlled by a parameter p_h , which represents the probability for us to swap the current
158 word with its homophone.

159 To simulate errors from confusion set, we constructed a dictionary using the confusion set identified
160 by Pedler and Mitton [19]. This set contains around 6000 pairs of words that are likely to be confused
161 with one another by people with dyslexia. The frequency of this type of error is controlled by p_s ,
162 representing the probability of a word being replaced by its paired word in the confusion set.

163 Examples of three types of injected errors are provided in Table 1. The original sentences are taken
164 from WMT14 (en2fr). Note that the perturbed sentences with homophone and confusion set errors
165 do not have misspellings but “real word errors” that are less likely to be detected and fixed by
166 spellcheckers before being sent for machine translations [24].

167 With this in mind, we are able to modify the WMT14 (en2fr) test dataset with different p values,
168 resulting different quantities of dyslexic errors injected into original source data. In this paper, we
169 focus on the percentage of words modified ranging from 10-20% as this follows findings from [23] in
170 real world dyslexic text error rate.

⁴<https://www.homophone.com>

Table 1: Example synthetic dyslexic sentences with injected dyslexic writing errors

| Error Injection | Original Sentence | Perturbed Sentence |
|------------------|--|---|
| Letter Confusion | In Nevada, where about 50 volunteers’ cars were equipped with the devices not long ago, drivers were uneasy about the government being able to monitor their every move. | In Nevada, where about 50 wolunteers ’ cars were equipped with thi devoces not iong ago, driverc were nneasy about the government being able to mohitor their every movov . |
| Homophone | New York City is looking into one. | New York City is looking into won . |
| Confusion Set | “The gas tax is just not sustainable,” said Lee Munnich, a transportation policy expert at the University of Minnesota. | “The gas tax is just knot sustainable,” said Lee Munnich, eye transportation policy export at the University of Minnesota. |

171 **3.2 Commercial Machine Translation Audit**

172 We chose to evaluate SOTA models that are deployed across major cloud computing platforms namely,
 173 AWS, Azure and Google Cloud. Based on a survey from *Public First* 51% of business utilize cloud
 174 services, majority of which are customers of AWS, Azure and Google Cloud ⁵. We also tested our
 175 dataset on GPT-3.5 (gpt-3.5-turbo-1106)⁶ a large language model (LLM). For each one of these
 176 services, we tested the performance of document translation, and for GPT we did a sentence-level
 177 translation (document translation was not available). For document translation, we submitted text
 178 files to the services for translation. For sentence-by-sentence translation, we were able to call the
 179 OpenAI API with Python scripts. All of these platforms require payment for the use of the translation
 180 services. For Google Cloud, we used the Cloud Translation API, for AWS, we used the Amazon
 181 Translate service and for Azure, we used the Translator in the Cognitive Services. Once the text was
 182 received we were able to evaluate the text.

183 **3.3 Evaluation Metrics**

184 We evaluate the performance of commercial MT services over synthetic dyslexic text with both
 185 lexical and semantic metrics. While the lexical metrics - such as BLEU [17] and WER [32] - allow
 186 us to benchmark against position our results in relation to a wide range of MT models and tasks, the
 187 semantic metrics - such as BERT and LaBSE - help illustrate how dyslexia might impact the user
 188 experience of these MT services.

189 **3.3.1 Lexical metrics**

190 Lexical based metrics have been commonly used in the evaluation of machine translation systems [13].
 191 One of the most popular lexical based metrics is Bilingual evaluation understudy (BLEU) [17], which
 192 is frequently used for in benchmarks and leaderboards. BLEU measures the n-gram similarity
 193 between MT output and the reference, and it is known for its simplicity, language-agnostics, and
 194 ability to measure both precision and fluency. BLEU score ranges from 0 to 1 where 1 indicates a
 195 perfect translation. State-of-the-art (SOTA) MT systems have reported BLEU score as high as 0.464
 196 for WMT14 (en2fr) task [14], which could be considered as generally “high quality translations”⁷.
 197 In contrast, BLEU scores lower than 0.2 would be considered “hard to understand” and “almost
 198 useless”.

199 The second lexical based metric we utilize is Word Error Rate (WER) [32], which measures the
 200 edit distance between MT output and the reference. As WER can be further broken down into the
 201 minimum number of word substitutions, insertions, and deletions required to convert the MT output

⁵<https://awsus.publicfirst.co/>

⁶<https://platform.openai.com/docs/models/gpt-3-5-turbo>

⁷BLEU Score Interpretations: <https://cloud.google.com/translate/automl/docs/evaluate>

202 to the reference sentence, this metric provides us additional insights into how the translation of
203 perturbed dyslexic sentences differ from the original sentences. While WER can range from zero to
204 infinity, a WER score higher than 0.5 generally suggests a poor performance.

205 3.3.2 Semantic Metrics

206 Since we are dealing with injected synthetic text, the lexical form of words are sometimes very similar
207 (for example in third row of Table 1 we have "knot" v. "not"). The edit distance between the two
208 samples is 1. However, the semantics of the words are completely different. This is where our lexical
209 metrics would likely fail. In order to fairly compare the sentences, we introduce semantic calculations.
210 The first method was using BERTScore [37] which computes a similarity score between 0 and 1
211 (where 1 is perfect) using contextual embeddings. The second evaluation metric we utilized was a
212 language independent method LaBSE [6] where we were able to use the source English sentences
213 from WMT directly for semantic comparison. We calculated the L2-norm of the sentence embeddings
214 from LaBSE to get the similarity between the source English sentences (without injections) to the
215 translations generated by the models. We called this the LaBSE score⁸. Same to the previous metric,
216 the score ranges between 0 and 1 where 1 indicates identical sentences and meaning. We must note
217 that a score of 1.0 requires the sentences to be syntactical identical. In other words, two sentences
218 with identical meanings but different writing would not score 1.0, but very close to 1.0.

219 4 Results

220 4.1 Lexical Divergence

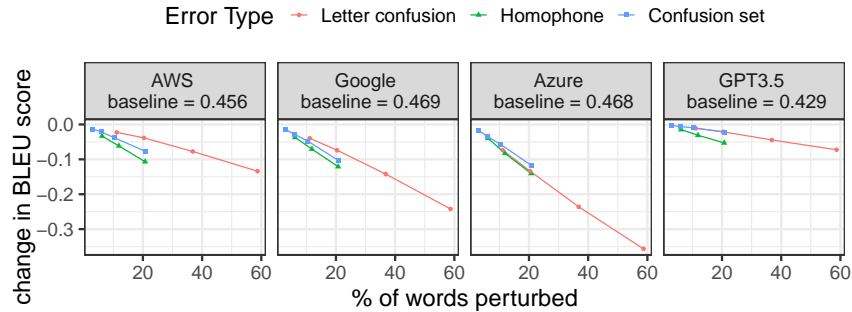
221 To measure how injected dyslexic errors influence translation results at a lexical level, we calculated
222 the BLEU and WER scores using the French translation from perturbed English sentences as
223 hypothesis and the original target sentences in French as references. We also calculated the BLEU
224 and WER scores for the translations generated by each MT service over the original, unperturbed
225 English data, as the baseline for our comparison.

226 We observed a SOTA level of performance in audited MT services at the baseline condition, with
227 BLEU score ranging from 0.429 (GPT3.5) to 0.469 (Google). However, the performance consistently
228 degrades as dyslexic style errors occur. Figure 1a shows a near linear drop in BLEU score, along with
229 the increase of words perturbed with dyslexic errors. While GPT3.5 has the lowest baseline BLEU
230 score, it is also least impacted by the increase of dyslexic errors. In contrast, the performance of
231 Azure MT drops most drastically when encountering more dyslexic errors. In terms of error types, we
232 notice that most services have more difficulties dealing with "real word errors" from homophone and
233 confusion set, rather than syntactic errors like letter confusion, with Azure being the only exception.
234 This observation is consistent with previous findings that real word errors in dyslexic writing pose
235 greater challenges for NLP models [19, 24].

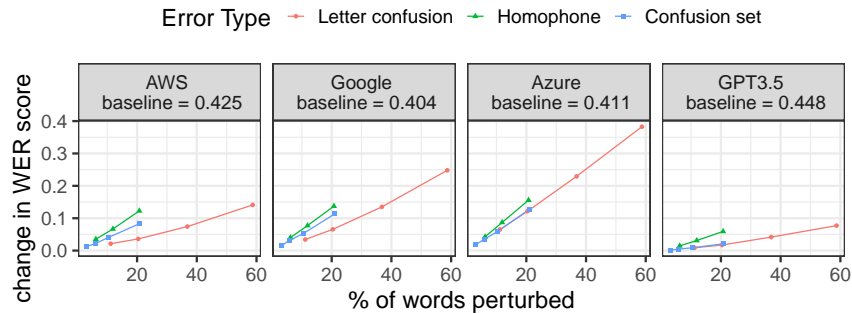
236 Similar trend is observed in WER scores. As shown in Figure 1b, for all audited services, their WER
237 scores increase steadily as more synthetic dyslexic errors are injected into the source data. The slope
238 of increase is greatest for homophone errors, and lowest for letter confusion. However, comparing to
239 AWS and GPT3.5, Google and Azure seem to be particularly challenged by letter confusion errors,
240 showing a degradation in translation quality almost as rapidly as when encountering synthetic real
241 word errors. Further inspection of their translation results in this condition suggests that the MT
242 services by Google and Azure are less likely to recover from a misspelled word, but tend to directly
243 copy it in the translation. For example, when the baseline sentence "*The American Civil Liberties
244 Union is deeply concerned*" is perturbed to become "*The American Cavil Liberties Union is deeply
245 concerned*", Google and Azure would translate the perturbed sentence to "*L'American Cavil Liberties
246 Union est profondément préoccupée*", with the misspelling "*Cavil*" preserved in the translation.

247 We also broke down the different types of edits used for calculating WER and inspect them separately.
248 Figure 2 shows the breakdown of substitutions, insertions, and deletions in the translation of 20%

⁸<https://huggingface.co/setu4993/LaBSE>



(a) BLEU scores drop as more dyslexic errors occur



(b) WER scores increase as more dyslexic errors occur

Figure 1: Change in lexical based metrics for all audited services. Baseline values indicate the metric score for unperturbed text, y-axis shows the change in corresponding metric compared to the baseline.

249 perturbed text from the reference. While the overall trends are similar for all MT services with three
 250 types of synthetic errors, we do observe some small difference in Azure and Google when handling
 251 letter confusion. These two services appear to make more deletions than insertions in their translation
 252 of text with letter confusion errors, suggesting potential loss of semantic information in the translation
 253 when source data contain significant amount of dyslexic misspellings. On the other hand, services
 254 like AWS and GPT3.5, despite more robust performance, tend to insert words in their translations. A
 255 deeper investigation on insertion errors found that articles are most often being inserted (see Figure 3
 256 for the most commonly added words by AWS with 20% confusion set errors).

257 While GPT3.5 generally perform better with synthetic dyslexic text, its performance still declines
 258 and could sometimes make serious mistakes due to dyslexic errors. For example, when the baseline
 259 sentence “*The technology is there to do it*” is perturbed to “*The technology is there to do ti.*”, the
 260 translation by GPT3.5 diverges from “*La technologie est là pour le faire*” to “*La technologie le*
 261 *frappe de plein fouet*” (“technology hitting it head on”).

262 4.2 Semantic Divergence

263 While lexical divergence, such as the insertion and deletion of particles, might not significantly impact
 264 the quality of translations, semantic change in the translation of dyslexic text from non-dyslexic
 265 text could have direct user experience consequences. While all audited services demonstrate high
 266 performance with unperturbed text at the semantic dimension (BERTScores and LaBSE scores all
 267 above 0.9), the semantic of the translation diverges as more dyslexic writing errors occur. As shown
 268 in Figure 4, both the BERTScore and LaBSE drops when the percentage of synthetic errors in text
 269 increases. Among all the audited services, the performance of Google and Azure declines most
 270 rapidly, while GPT3.5 maintains a relatively robust level of performance.



Figure 2: Breakdown of WER scores by edit type (20% word perturbed)

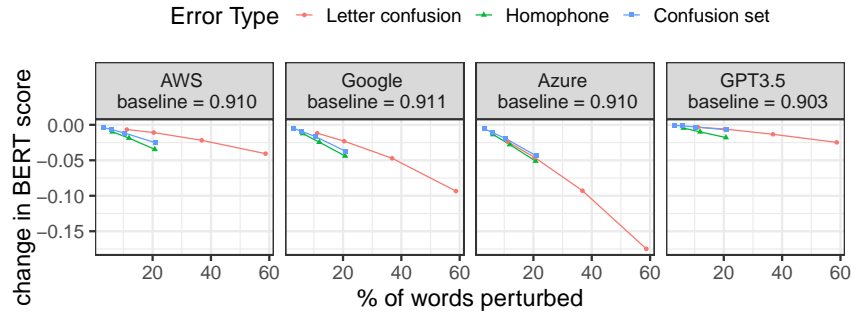


Figure 3: Word cloud of word confusion AWS (20% word modified)

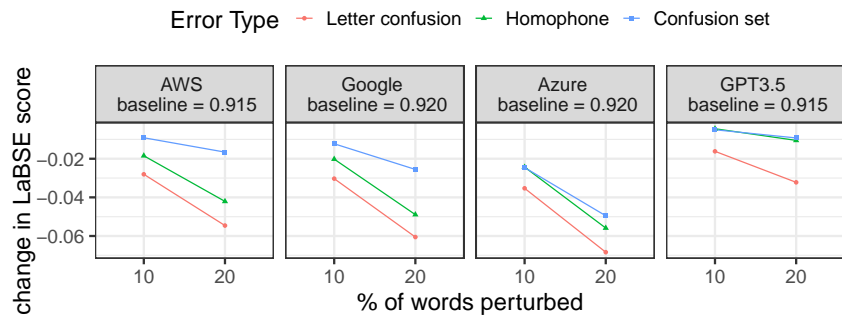
271 Even if the semantic divergence is smaller comparing to the lexical divergence, the disparity between
 272 the baseline and text with 20% dyslexic errors is statistically significant, suggesting a clear gap in
 273 MT service quality for dyslexic users.

274 5 Discussion

275 Our results uncover potential disparities in the quality of MT services for people with and without
 276 dyslexia. As part of the cloud infrastructure, these services have been ubiquitously adopted as
 277 foundation for many other digital products and services. Our work shows how typical dyslexic
 278 writing errors could lead to the degradation of SOTA MT services. Even advanced LLMs, which
 279 have been believed as a solution for dyslexia, struggle with real word errors from homophones and



(a) BERTScore drop as more dyslexic errors occur



(b) LaBSE scores drop as more dyslexic errors occur

Figure 4: Change in semantic metrics for all audited services. Baseline values indicate the metric score for unperturbed text, y-axis shows the change in corresponding metric in comparison to the baseline.

280 confusion set. While LLMs are better than other services in terms of lexical and syntactic mistakes,
 281 they do still produce semantic divergence when translating dyslexic text, and such divergence could
 282 be even harder to be noticed by users with dyslexia, resulting in higher user risk and potentially worse
 283 experience in the long term.

284 6 Limitations and Future Work

285 Although we were able to experiment with a wide variety of configurations with the quantities and
 286 types of dyslexic writing errors, our synthetic datasets are nevertheless limited in their ability to
 287 capture the full heterogeneity of dyslexic writing. Like any other neurodivergence, dyslexia affects
 288 people differently: the way it manifests in writing differs across individuals and situations. More
 289 authentic, real world data from people with dyslexia is required to better represent this community
 290 in AI data in order to develop fair and inclusive NLP models for dyslexia. We also look forward to
 291 extend our methodology to other communities and application domains, making it easier to audit a
 292 wide range of AI models and services using synthetic data about marginalized, sensitive populations.

293 7 Conclusion

294 We proposed a novel method to generate synthetic dyslexia datasets and leveraged them to identify
 295 performance disparities in SOTA machine translation services for people with dyslexia. Our lexical
 296 and semantic metrics allow us to benchmark and better understand existing disparities. Our work
 297 highlights the importance of making NLP and AI more inclusive and equitable to communities most
 298 impacted by such technologies. We call for attention from language technology researchers and
 299 developers to close the equity gap for users with dyslexia.

References

- [1] Ricardo Baeza-Yates and Luz Rello. Estimating dyslexia in the web. page 8, 03 2011. doi: 10.1145/1969289.1969300.
- [2] Ondrej Bojar, Christian Buck, Christian Federmann, Barry Haddow, Philipp Koehn, Johannes Leveling, Christof Monz, Pavel Pecina, Matt Post, Herve Saint-Amand, Radu Soricut, Lucia Specia, and Aleks Tamchyna. Findings of the 2014 workshop on statistical machine translation. In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, pages 12–58, Baltimore, Maryland, USA, June 2014. Association for Computational Linguistics. URL <http://www.aclweb.org/anthology/W/W14/W14-3302>.
- [3] Nicola Brunswick. Unimpaired reading development and dyslexia across different languages. *Reading and dyslexia in different orthographies*, pages 131–154, 2010.
- [4] Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*, pages 77–91. PMLR, 2018.
- [5] Xingyu Chen, Zhengxiong Li, Srirangaraj Setlur, and Wen Yao Xu. Exploring racial and gender disparities in voice biometrics. *Scientific Reports*, 12:3723, 03 2022. doi: 10.1038/s41598-022-06673-y.
- [6] Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. Language-agnostic bert sentence embedding, 2022.
- [7] William Gavia Rojas, Sudnya Diamos, Keertan Kini, David Kanter, Vijay Janapa Reddi, and Cody Coleman. The dollar street dataset: Images representing the geographic and socioeconomic diversity of the world. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 12979–12990. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/5474d9d43c0519aa176276ff2c1ca528-Paper-Datasets_and_Benchmarks.pdf.
- [8] Steven M. Goodman, Erin Buehler, Patrick Clary, Andy Coenen, Aaron Donsbach, Tiffanie N. Horne, Michal Lahav, Robert MacDonald, Rain Breaw Michaels, Ajit Narayanan, Mahima Pushkarna, Joel Riley, Alex Santana, Lei Shi, Rachel Sweeney, Phil Weaver, Ann Yuan, and Meredith Ringel Morris. Lampost: Design and evaluation of an ai-assisted email writing prototype for adults with dyslexia. In *Proceedings of the 24th International ACM SIGACCESS Conference on Computers and Accessibility, ASSETS '22*, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450392587. doi: 10.1145/3517428.3544819. URL <https://doi.org/10.1145/3517428.3544819>.
- [9] Priya Goyal, Adriana Romero Soriano, Caner Hazirbas, Levent Sagun, and Nicolas Usunier. Fairness indicators for systematic assessments of visual feature extractors. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency, FAccT '22*, page 70–88, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450393522. doi: 10.1145/3531146.3533074. URL <https://doi.org/10.1145/3531146.3533074>.
- [10] Laura Gustafson, Megan Richards, Melissa Hall, Caner Hazirbas, Diane Bouchacourt, and Mark Ibrahim. Exploring why object recognition performance degrades across income levels and geographies with factor annotations. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 24729–24753. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/4e3378a8e80af4ffc456c4fa13d46550-Paper-Datasets_and_Benchmarks.pdf.
- [11] S. Gunnell Ingesson. Growing up with dyslexia. *School Psychology International*, 28(5): 574–591, 2007. doi: 10.1177/0143034307085659.

- 348 [12] Tedd Kourkounakis, Amirhossein Hajavi, and Ali Etemad. Fluentnet: End-to-end detection of
349 speech disfluency with deep learning, 2020.
- 350 [13] Seungjun Lee, Jungseob Lee, Hyeonseok Moon, Chanjun Park, Jaehyung Seo, Sugyeong Eo,
351 Seonmin Koo, and Heuseok Lim. A survey on evaluation metrics for machine translation.
352 *Mathematics*, 11(4), 2023. ISSN 2227-7390. doi: 10.3390/math11041006. URL <https://www.mdpi.com/2227-7390/11/4/1006>.
353
- 354 [14] Xiaodong Liu, Kevin Duh, Liyuan Liu, and Jianfeng Gao. Very deep transformers for neural
355 machine translation. *ArXiv*, abs/2008.07772, 2020. URL <https://api.semanticscholar.org/CorpusID:221150462>.
356
- 357 [15] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture
358 models, 2016.
- 359 [16] Mirja Mittermaier, Marium Raza, and Joseph Kvedar. Bias in ai-based models for medical
360 applications: challenges and mitigation strategies. *npj Digital Medicine*, 6, 06 2023. doi:
361 10.1038/s41746-023-00858-z.
- 362 [17] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic
363 evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association
364 for Computational Linguistics*, pages 311–318, 2002.
- 365 [18] Jennifer Pedler. *Computer correction of real-word spelling errors in dyslexic text*. PhD thesis,
366 University of London, 2007.
- 367 [19] Jennifer Pedler and Roger Mitton. A large list of confusion sets for spellchecking assessed
368 against a corpus of real-word errors. In *Proceedings of the Seventh International Conference on
369 Language Resources and Evaluation (LREC'10)*, Valletta, Malta, May 2010. European Lan-
370 guage Resources Association (ELRA). URL [http://www.lrec-conf.org/proceedings/
371 lrec2010/pdf/122_Paper.pdf](http://www.lrec-conf.org/proceedings/lrec2010/pdf/122_Paper.pdf).
- 372 [20] Joseph J Pollock and Antonio Zamora. Automatic spelling correction in scientific and scholarly
373 text. *Communications of the ACM*, 27(4):358–368, 1984.
- 374 [21] Alberto Quattrini Li, Licia Sbattella, and Roberto Tedesco. Polispell: An adaptive spellchecker
375 and predictor for people with dyslexia. In Sandra Carberry, Stephan Weibelzahl, Alessandro
376 Micarelli, and Giovanni Semeraro, editors, *User Modeling, Adaptation, and Personalization*,
377 pages 302–309, Berlin, Heidelberg, 2013. Springer Berlin Heidelberg. ISBN 978-3-642-38844-
378 6.
- 379 [22] Luz Rello, Clara Bayarri, and Azuki Gorriz. What is wrong with this word? dyseggxia: a
380 game for children with dyslexia. In *Proceedings of the 14th International ACM SIGACCESS
381 Conference on Computers and Accessibility*, ASSETS '12, page 219–220, New York, NY, USA,
382 2012. Association for Computing Machinery. ISBN 9781450313216. doi: 10.1145/2384916.
383 2384962. URL <https://doi.org/10.1145/2384916.2384962>.
- 384 [23] Luz Rello, Ricardo A. Baeza-Yates, and Joaquim Llisterri. Dyslist: An annotated resource of
385 dyslexic errors. In *International Conference on Language Resources and Evaluation*, 2014.
- 386 [24] Luz Rello, Miguel Ballesteros, and Jeffrey P. Bigham. A spellchecker for dyslexia. In *Proc.
387 of ASSETS*, 2015. ISBN 978-1-4503-3400-6. doi: 10.1145/2700648.2809850. URL <http://doi.acm.org/10.1145/2700648.2809850>.
388
- 389 [25] Luz Rello, Miguel Ballesteros, and Jeffrey P. Bigham. A spellchecker for dyslexia. In
390 *Proceedings of the 17th International ACM SIGACCESS Conference on Computers & Ac-
391 cessibility*, ASSETS '15, page 39–47, New York, NY, USA, 2015. Association for Com-
392 puting Machinery. ISBN 9781450334006. doi: 10.1145/2700648.2809850. URL <https://doi.org/10.1145/2700648.2809850>.
393

- 394 [26] Lindsay Reynolds and Shamoeli Wu. “i’m never happy with what i write”: Challenges and
395 strategies of people with dyslexia on social media. 2018.
- 396 [27] B. Riddick. *Living With Dyslexia: The social and emotional consequences of specific learning*
397 *difficulties/disabilities*. nasen spotlight. Taylor & Francis, 2009. ISBN 9781135191740. URL
398 <https://books.google.com/books?id=NveMAgAAQBAJ>.
- 399 [28] S. E. Shaywitz and B. A. Shaywitz. Dyslexia (specific reading disability). *Biological Psychiatry*,
400 57:1301–1309, 2005.
- 401 [29] Sally E. Shaywitz, Michael Escobar, Bennett Shaywitz, Jack Fletcher, and Robert Makuch.
402 Evidence that dyslexia may represent the lower tail of a normal distribution of reading ability.
403 326:145–50, 02 1992.
- 404 [30] Anders Søgaard. Should we ban English NLP for a year? In Yoav Goldberg, Zornitsa Kozareva,
405 and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural*
406 *Language Processing*, pages 5254–5260, Abu Dhabi, United Arab Emirates, December 2022.
407 Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.351. URL
408 <https://aclanthology.org/2022.emnlp-main.351>.
- 409 [31] Emma Stanley, Matthias Wilms, and Nils Daniel Forkert. *Disproportionate Subgroup Impacts*
410 *and Other Challenges of Fairness in Artificial Intelligence for Medical Image Analysis*, pages
411 14–25. 12 2022. ISBN 978-3-031-23222-0. doi: 10.1007/978-3-031-23223-7_2.
- 412 [32] Keh-Yih Su, Ming-Wen Wu, and Jing-Shin Chang. A new quantitative quality measure for
413 machine translation systems. In *COLING 1992 Volume 2: The 14th International Conference*
414 *on Computational Linguistics*, 1992.
- 415 [33] US Interagency Committee on Learning Disabilities. *Learning Disabilities: A Report to the*
416 *U.S. Congress*. 1987. URL <https://books.google.com/books?id=0G-7PQAACAAJ>.
- 417 [34] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N
418 Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon,
419 U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, ed-
420 itors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates,
421 Inc., 2017. URL [https://proceedings.neurips.cc/paper_files/paper/2017/file/](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf)
422 [3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).
- 423 [35] Shaomei Wu, Lindsay Reynolds, Xian Li, and Francisco Guzmán. Design and evaluation of a
424 social media writing support tool for people with dyslexia. In *Proceedings of the 2019 CHI*
425 *Conference on Human Factors in Computing Systems*, CHI ’19, page 1–14, New York, NY,
426 USA, 2019. Association for Computing Machinery. ISBN 9781450359702. doi: 10.1145/
427 3290605.3300746. URL <https://doi.org/10.1145/3290605.3300746>.
- 428 [36] Emmanuel J Yannakoudakis and David Fawthrop. The rules of spelling errors. *Information*
429 *Processing & Management*, 19(2):87–99, 1983.
- 430 [37] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. Bertscore:
431 Evaluating text generation with bert, 2020.