

MEDIT: Multilingual Text Editing via Instruction Tuning

Anonymous NAACL submission

Abstract

We introduce MEDIT, a multi-lingual extension to COEDIT – the recent state-of-the-art text editing models for writing assistance. MEDIT models are trained by fine-tuning multilingual large, pre-trained language models (LLMs) via instruction tuning. They are designed to take instructions from the user specifying the attributes of the desired text in the form of natural language instructions, such as “Grammatik korrigieren” (German) or ⓠ 텍스트를 단순화 (Korean). We build MEDIT by curating data from multiple publicly available human-annotated text editing datasets for three text editing tasks (Grammatical Error Correction (GEC), Text Simplification, and Paraphrasing) across diverse languages belonging to six different language families. We detail the design and training of MEDIT models and demonstrate their strong performance on many multilingual text editing benchmarks against other multilingual LLMs. We also find that MEDIT generalizes effectively to new languages over multilingual baselines. We publicly release our code and trained models to the community.

1 Introduction

Large language models (LLMs) have made remarkable progress toward generating fluent and coherent text in a wide variety of tasks and domains to support writing assistance (Brown et al. 2020b; OpenAI 2023; Touvron et al. 2023; *inter alia*). In particular, LLMs have been adapted to perform many complex text editing tasks like GEC (Wu et al., 2023; Coyne and Sakaguchi, 2023; Fang et al., 2023b), text simplification (Baez and Saggon, 2023; Saggion et al., 2022), paraphrasing (Witteveen and Andrews, 2019; Niu et al., 2021), and formality and tone rewriting (Reif et al., 2022; Luo et al., 2023), among others. However, most of these works are restricted to single tasks, with few works adapting LLMs to perform high-quality text editing across multiple tasks (Schick et al.,

Multilingual Editing

Parafrasee la oración: Hoy iré a la escuela a estudiar español.
Hoy asistiré a la escuela para aprender español.

Cross-lingual Editing

文を言い換えてください: Hoy iré a la escuela a estudiar español.
Hoy asistiré a la escuela para aprender español.

Figure 1: Examples illustrating multilingual and cross-lingual text editing. The editing instructions are described in bold. Note that the input and output texts are always in the same language. The monolingual vs. cross-lingual setting is determined by comparing the language of the edit instruction in relation to the language of the input text.

2023; Raheja et al., 2023; Laban et al., 2023). A lot of these improvements have been driven by fine-tuning large language models (LLMs) with task-specific instruction tuning, resulting in remarkable zero-shot generalization abilities (Sanh et al. 2022; Ouyang et al. 2022; Chung et al. 2022; *inter alia*).

At the same time, significant research effort has been dedicated to leveraging and enhancing the multilingual capabilities of LLMs (Lin et al., 2022). These abilities can be improved using methods such as continued pre-training with abundant monolingual data (Yang et al., 2023b; Cui et al., 2023) or language-specific instruction-tuning (Zhu et al., 2023; Li et al., 2023). However, in the case of continued pre-training, the lack of high-quality web-scale data often restricts the ability to improve LLM capabilities in less-represented languages in the same way that English data can be expanded. Moreover, while numerous multilingual instruction-tuned models have been developed (Muennighoff et al., 2023; Workshop, 2023; Xue et al., 2021; Li et al., 2023; Wei et al., 2023), our analyses show that without further task-specific fine-tuning, these models are not suitable for car-

Language	ISO-639-1	Family
Arabic	ar	Semitic
Chinese	zh	Sino-Tibetan
English	en	Germanic
German	de	
Japanese	ja	Japonic
Korean	ko	Koreanic
Spanish	es	Romance

Table 1: **Set of Languages.** The seven languages, along with the ISO-639-1 code and their language family, on which we train and evaluate our models.

066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103

trying out high-quality text-editing tasks (§ 5.1). In the context of text editing tasks, multiple previous works have developed high-quality, general-purpose LLMs on non-English languages, restricting themselves, however, on either specific tasks (Rothe et al., 2021; Sun et al., 2022; Kementchedjhieva and Søgaard, 2023; Ryan et al., 2023; Krishna et al., 2022; Lai et al., 2022) or specific languages (Alhfni et al., 2023; Anschütz et al., 2023). Overall, the aforementioned factors have limited the availability of high-quality *multilingual text editing* (MTE) models, which has limited their usability for writing assistance across multiple tasks in languages beyond English.

We address these gaps with MEDIT, a multitask, multilingual extension of COEDIT (Raheja et al., 2023). MEDIT models can perform text editing operations for three popular tasks: Grammatical Error Correction, Paraphrasing, and Text Simplification, in multilingual and cross-lingual settings (Figure 1) across a diverse set of seven languages, spanning six different language families (Table 1).

To build MEDIT, we fine-tune several multilingual LLMs of varying sizes on carefully curated, largely human-annotated, parallel corpora of over 200k instructional input-output pairs, using publicly available datasets (Table 4) for different text editing tasks. We evaluate the performance of our models extensively on text editing benchmarks in both multilingual and cross-lingual settings to demonstrate their effectiveness.

Our contributions are as follows:

- This work, to the best of our knowledge, is the first to investigate multi-task, multilingual text editing via instruction tuning.
- Our models achieve strong performance on multiple text editing tasks across numerous languages, and are publicly available for fostering further MTE research.¹

¹[anonymous.link](#)

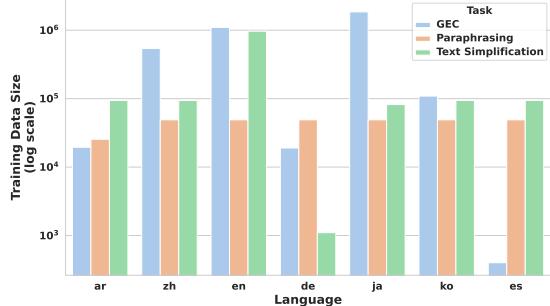


Figure 2: **Data distribution for each of the three tasks and seven languages on which we train.** The plot of the amount of data in log scale to aid visualization.

- 104
105
106
107
108
109
110
- Through a comprehensive set of controlled experiments, we provide insights on how model performance on multilingual text editing tasks is affected by various choices like model architecture, model scale, and training data mixtures.

2 Related Work

Multi-lingual LLMs for Text Editing There is an extensive body of prior literature that has leveraged LLMs for various multi-lingual text editing tasks. These works have proposed models for text editing tasks like GEC (Rothe et al., 2021; Sun et al., 2022), paraphrasing (Chowdhury et al., 2022), formality style transfer (Briakou et al., 2021), and text simplification (Mallinson et al., 2020; Martin et al., 2022; Ryan et al., 2023). However, all of these prior approaches have proposed task-specific multi-lingual models. In contrast, we propose a single unified text-editing model for all the considered tasks by leveraging the power of instruction-tuning and task-specific fine-tuning, which enables our multi-lingual models to generalize to multiple text-editing tasks.

Multi-lingual Instruction-Tuning While numerous multi-lingual instruction fine-tuned models like Muennighoff et al. (2023); Wei et al. (2023) and Li et al. (2023) have been developed, they are not focused or tailored for text editing tasks, which we address by task-specific fine-tuning. Specific to text editing, many prior works have explored instruction tuning capable of performing multiple text editing tasks with a single model, such as GEC, simplification, sentence fusion, style transfer, and paraphrasing, to name a few (Mallinson et al., 2022; Du et al., 2022; Kim et al., 2022; Schick et al., 2023; Raheja et al., 2023). However, while they are able to support multi-task text editing, they are generally mono-lingual and typically restricted to

143 specific languages (predominantly English). Our
144 work, thus, addresses this gap by proposing multi-
145 lingual, instruction-tuned models for multiple text
146 editing and revision tasks.

147 3 MEDIT

148 3.1 Tasks and Languages

149 We chose a broad set of languages to ensure cover-
150 age and chose text editing tasks that had multilingual,
151 publicly available human-annotated datasets
152 to ensure high data quality. Another criteria was to
153 choose languages at the intersection of the publicly
154 available corpora we could find across a large set
155 of languages for all the tasks we considered. We
156 refer to this as the MEDIT dataset.

157 **Table 1** describes the languages covered in our
158 work, whereas **Figure 2** depicts the amounts of
159 training datasets that were available for all tasks
160 and languages we considered. **Appendix A** details
161 all the training and testing datasets.

162 3.2 Models

163 We fine-tune different versions of pre-trained multi-
164 lingual LLMs (both encoder-decoder/sequence-to-
165 sequence (Seq2Seq) and decoder-only/causal lan-
166 guage models (CLM)) on the MEDIT dataset using
167 cross-entropy loss. The details of the MEDIT mod-
168 els are described in § 4.2, whereas the training
169 details are summarized in § 4.3.

170 4 Experiments

171 4.1 No-Edits Baseline

172 We first evaluate a no-edits baseline, where the out-
173 put is simply a copy of the source input without the
174 instruction. This strategy performs reasonably well
175 on tasks where the target output largely overlaps
176 with the input (e.g., GEC).

177 4.2 Multilingual LLMs

178 **mT5** ([Xue et al., 2021](#)) is a multilingual variant
179 of T5 ([Raffel et al., 2020](#)), trained on the mC4
180 dataset, a multilingual variant of the C4 dataset
181 extended to 101 languages. We experiment with
182 three variants of mT5 – LARGE (770M), XL (3B),
183 and XXL (13B) parameters.

184 **mT0** ([Muennighoff et al., 2023](#)) is a family of
185 multilingual Seq2Seq models capable of zero-shot
186 following human instructions in dozens of lan-
187 guages. We use the mt0-LARGE (1.2B), mt0-XL
188 (3.7B), and mt0-XXL (13B) models. These models

189 are constructed by fine-tuning mT5 models on the
190 xP3 cross-lingual task mixture dataset, which con-
191 sists of multilingual datasets with English prompts.
192 As a result, mT0 models are better suited for follow-
193 ing English prompts. We also use the mt0-XXL-MT
194 variant, which is fine-tuned on the xP3mt dataset
195 and is better suited for prompting in non-English.

196 **BLOOMZ** ([Muennighoff et al., 2023](#)) is a family
197 of multilingual Causal Language Models (CLMs)
198 constructed by fine-tuning BLOOM ([Workshop,](#)
199 [2023](#)) on the xP3 dataset. We use BLOOMZ-3b
200 and BLOOMZ-7b1 models for our experiments.

201 **PolyLM** ([Wei et al., 2023](#)) is a set of multilin-
202 gual LLMs trained on 640B tokens. Furthermore,
203 we also use the **PolyLM-MultiAlpaca-13B** model,
204 which is PolyLM fine-tuned on the MULTIALPACA
205 dataset, consisting of 132k samples of multilingual
206 instructions.

207 **Bactrian-X** ([Li et al., 2023](#)) is a collection of
208 lightweight adapters for LLaMA (7B and 13B)
209 ([Touvron et al., 2023](#)) and BLOOM (7B) ([Work-
210 shop, 2023](#)) on the Bactrian-X dataset, which is
211 a multilingual parallel dataset comprising 3.4 mil-
212 lion instruction-response pairs across 52 languages.
213 For simplicity, we only compare against its higher-
214 performant LLaMA-adapted versions.

215 4.2.1 Large-Pretrained Decoder-only Models

216 We also conduct zero-shot evaluations against state-
217 of-the-art decoder-only LLMs that have shown im-
218 pressive multilingual capabilities on a variety of
219 NLP tasks leveraging the power of in-context learn-
220 ing ([Lai et al., 2023; OpenAI, 2023](#)).

221 **GPT3.5** (also referred to as ChatGPT),² is an im-
222 proved version of GPT3 ([Brown et al., 2020a](#)) op-
223 timized for chat. We use the gpt-3.5-turbo0613
224 model from the OpenAI API.³

225 **GPT4** ([OpenAI, 2023](#)) is the latest iteration of
226 the GPT models and is also optimized for chat. We
227 use the gpt-4-0613 model from the OpenAI API.

228 While we recognize that these models may not
229 be explicitly optimized or trained for multi-lingual
230 settings, considering that they have been trained on
231 massive amounts of web-scale data, these models
232 have been shown to have multi-lingual capabilities
233 ([Lai et al., 2023](#)), hence, we consider them as one
234 of our baseline groups.

²<https://openai.com/blog/chatgpt>

³<https://api.openai.com>

235 4.3 Training Setup

236 We perform instruction tuning for all our models
237 by crafting custom prompts for each of the
238 21 task-language combinations (seven languages,
239 three tasks). Similar to COEDIT, for each task-
240 language combination, depending on the number
241 of ways the instructions can be translated without
242 altering the meaning, we write between 14 and
243 27 instructions by automatically translating each
244 one from English and verifying the accuracy of the
245 translations by asking native language speakers to
246 evaluate and correct them. The total number of task-
247 language instructions is 231, which can be found
248 in Appendix D. We explore three different multilin-
249 gual and cross-lingual instructional settings,
250 depending on the language of the prompt, where
251 the editing instruction could be in (a) English, (b)
252 the same language as the text being edited (Native),
253 and (c) a random language which may or may not
254 be the same as the language of the text being edited
255 (Random). With this definition, English and Random
256 are cross-lingual text editing tasks, and Native
257 is a multilingual text editing task. In all settings,
258 the input-output pairs are in the same language, but
259 only the language of the instruction changes.

260 We train all models on 8xA100 80G GPU
261 instances for five epochs. For the PolyLM and
262 Bactrian-X models (>7B parameters), we also use
263 LoRA (Hu et al., 2022) to lower the number of
264 trainable parameters and increase the batch size.

265 4.4 Evaluation

266 For GEC evaluation, we follow prior work on each
267 language we report on and use the appropriate GEC
268 metric accordingly. Mainly, we use the MaxMatch
269 (M^2) Scorer (Dahlmeier and Ng, 2012), ERRANT
270 (Bryant et al., 2017), and GLEU (Napoles et al.,
271 2015, 2016) as our evaluation metrics. The M^2
272 Scorer and ERRANT compare the edits made by a
273 GEC system against annotated reference edits and
274 calculate the precision (P), recall (R), and $F_{0.5}$ (i.e.,
275 weighing precision twice as much as recall). GLEU
276 computes the precision of the n-grams that overlap
277 with the references but not the original texts and
278 penalizes n-grams that overlap with the original
279 texts but not the references.

280 For simplification, we follow Ryan et al. (2023)
281 and use SARI (Xu et al., 2016a) and BLEU (Pap-
282 ineni et al., 2002) for evaluation. SARI is the aver-
283 age of the F1 score for adding, keeping, and delet-
284 ing n-grams ($n \in 1, 2, 3, 4$) and has been shown

285 to correlate with human judgments of simplicity
286 (Xu et al., 2016a). BLEU, on the other hand, is a
287 common metric in machine translation and is used
288 as a check for grammatical and meaning preser-
289 vation. We compute all metrics using the EASSE
290 evaluation suite (Alva-Manchego et al., 2019).

291 We evaluate paraphrasing on two criteria and
292 metrics: diversity and semantic similarity. For di-
293 versity, we use Self-BLEU (Zhu et al., 2018) to
294 measure the diversity of the paraphrases relative
295 to the given source and reference texts. We use
296 Semantic Similarity to measure meaning preser-
297 vation. Specifically, we use mUSE (Yang et al.,
298 2020) for this, as it is the best-performing mul-
299 tilingual sentence similarity model that supports
300 all the languages in our work. We also consid-
301 ered other notable works that have made significant
302 progress on multilingual sentence similarity, such
303 as Multilingual-SBERT (Reimers and Gurevych,
304 2020) and LaBSE (Feng et al., 2022). However,
305 we found them unsuitable for our purposes as they
306 were either limited by the languages they support or
307 suffered from lower performance for multilingual
308 sentence similarity.

309 5 Quantitative Results: MTE Quality

310 We split the models into three main groups. The
311 first group (a) consists of the “no edits” baseline,
312 the second group (b) is the untrained baseline,
313 where models are evaluated in a zero-shot setting
314 without any task-specific fine-tuning, while the
315 third group (c) is our set of multi-lingual models
316 trained on task-specific datasets. For all our exper-
317 iments, we aggregate the models’ performance by
318 text editing tasks. Specifically, we aggregate the
319 metrics for each task using the harmonic means
320 of its constituents. Specifically, we use (1-Self-
321 BLEU)⁴ and Semantic Accuracy for Paraphrasing,
322 SARI and BLEU for Simplification, and $F_{0.5}$ and
323 GLEU for GEC. We scale all metrics to lie between
324 0 and 100. We show full results on all models in
325 Appendix C for the best-performing setup.

326 5.1 Baselines

327 In Figure 3, we report the results of our trained
328 models against various baselines by aggregating
329 the performance on all tasks (as detailed in § 5).

330 **No Edits (Copy) Baseline** We observe that not
331 making any edits leads to a performance that is on

332 ⁴We subtract Self-BLEU from 1 because lower is better in
333 terms of making changes to the source text.

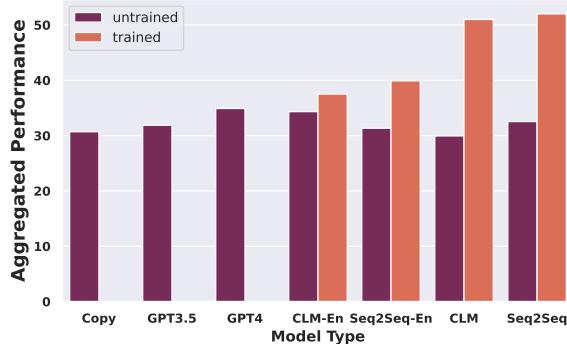


Figure 3: **Overall performance comparison of all baselines against trained models.** We calculate the aggregated performance across all tasks using the harmonic mean of task-specific scores. Baselines are “No Edits (Copy)”, “English-only” (“-En,”), and our trained models are marked as “CLM” and “Seq2Seq,” respectively. The aggregated performance is calculated as described in § 5.

332 par with the untrained versions of all models, which
 333 highlights the limitations of the n-gram overlap-
 334 based metrics.

335 **Untrained Baseline** Similar to Raheja et al.
 336 (2023), a core contribution of this work is to push
 337 the performance of small- ($\sim 1B$ parameters) to
 338 medium-sized (1-15B parameters) LLMs for com-
 339 mon text editing tasks across multiple languages.
 340 This drives the need for fine-tuning task-specific
 341 and language-specific datasets. For this work, we
 342 compare our fine-tuned models against their non-
 343 fine-tuned counterparts. We find a substantial gap
 344 between the untrained models and their trained
 345 counterparts, highlighting the impact of task- and
 346 language-specific fine-tuning for the tasks.

347 **English-Only Baseline** In this experiment, we
 348 analyze the ability of multilingual LLMs to adapt
 349 to different text editing tasks across different lan-
 350 guages by fine-tuning them in the most prominent
 351 high-resource language (English). We fine-tune
 352 all the multi-lingual models on just the English
 353 subsets of the training data, as it is the largest in
 354 terms of quality and quantity. This experiment tests
 355 the need for language-specific data. Similar to the
 356 previous result, we observe that the gap between
 357 the untrained versions of English-only models is
 358 relatively small vs. the ones trained on the full
 359 dataset; it increases significantly with language-
 360 specific fine-tuning.

361 **State-of-the-art LLMs** Additionally, we also
 362 evaluate against the most powerful commercially
 363 available LLMs on their ability to perform MTE.

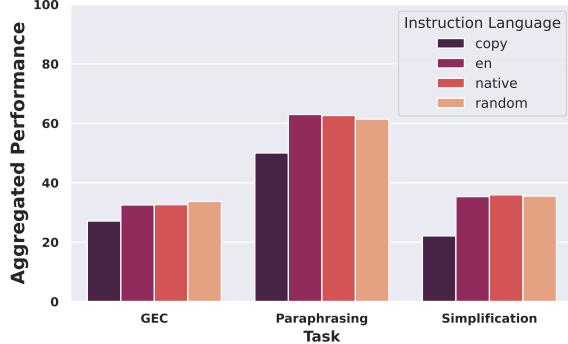


Figure 4: **Aggregated performance on different tasks broken down by instruction language.** Apart from some minor fluctuation there is no significant impact of instruction language on our results.

364 Specifically, we evaluate GPT3.5 and GPT-4 in a
 365 zero-shot setting. Although these models have been
 366 shown to exhibit strong zero-shot performance on
 367 a variety of NLP tasks, we find that the overall
 368 performance of both models is close to most un-
 369 trained baselines. This can be attributed to the
 370 rather limited multilingual capabilities of GPT3.5,
 371 which often lead to outputs being generated in other
 372 languages (English in particular). To some extent,
 373 the verbosity of responses is highly detrimental to
 374 the performance, especially for GPT4 as it gets
 375 penalized by the automatic metrics (especially for
 376 GEC).

377 The rest of this section analyzes different aspects
 378 of the quantitative performance of our models.

5.2 Model Performance by Language

380 In this section, we analyze the performance of dif-
 381 ferent MTE models by language (Figure 5). It is in-
 382 teresting to note that Paraphrasing exhibits a rather
 383 steady performance across languages. This can
 384 partially be attributed to the weakness of the evalua-
 385 tion metrics as they rely mostly on n-gram over-
 386 lap, on the multilingual pre-training of the LLMs
 387 where they are exposed to medium-large corpora
 388 of nearly all the languages, but also that with the
 389 increase in the model size and fine-tuning, they
 390 tend to make fewer changes, thus, leading to higher
 391 scores. For simplification, the variance in the per-
 392 formance across languages can be attributed to the
 393 amount and quality of training data available for
 394 each task. For instance, for German, only 1.1k
 395 training data points were available, which leads
 396 to models not only showing a great improvement
 397 in performance with fine tuning, but also a great
 398 variance. Similarly, the training data is very noisy
 399 for Japanese, which leads to a similar effect. For

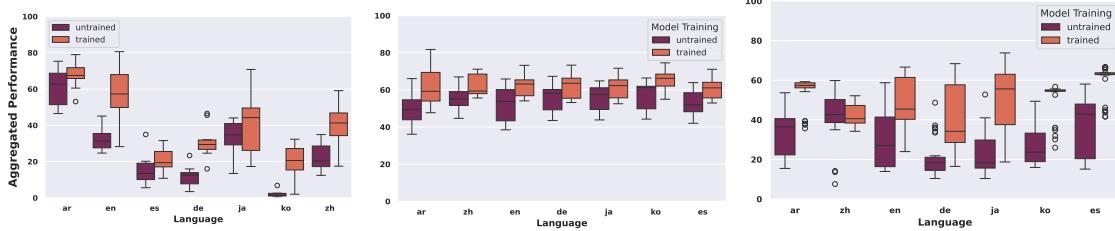


Figure 5: **Aggregated model performance by language** for (a) GEC, (b) Paraphrasing, and (c) Simplification. For each task, we aggregate the relevant metrics as described in § 5 and split them by model training.

GEC, we observe that the performance varies a lot by language, indicating the challenging nature of the task. This can be partially attributed to the frequency and the type of errors in each dataset, a phenomenon we see in Arabic datasets. Moreover, the quality of the training GEC data available also leads to varying performance across languages.

5.3 Language of Instruction

Here, we analyze the effect of the language of the instruction used to instruct the model. As mentioned in § 4.3, we have three configurations for this set of experiments:

English-language Instructions We train the first set of multilingual MEDIT models with just English instructions. These are MTE models capable of performing cross-lingual text editing, trained on data where the instruction is always in English.

Native-language Instructions We train the next set of multi-lingual MEDIT models with instructions in their native language. These models are capable of performing MTE, where the language of edit instructions is the same as the language of texts being edited.

Randomized-language Instructions Finally, we also explore the cross-lingual text editing (Figure 1) abilities of multi-lingual LLMs. To do so, we modify our dataset by appending an edit instruction from a randomly chosen language (different from the language of the edited source-target text pair) and train our models on this cross-lingual-prompted dataset.

Figure 4 shows the effect of instruction language on performance on all three tasks. We note that there is no significant difference in performance between the different settings. This is likely because in each setting, owing to its multilingual instructional pre-training, the model is able to adapt well to the language of the instruction in the fine-tuning phase, hence focusing mostly on the specific tasks.

5.4 Effect of Model Architecture

We present the task-specific results across each model type in Figure 6, observing that CLMs generally are either on par or outperform the rest of the models with GPT3.5, yielding the lowest results.

BLEU relies on n-gram overlap artificially boosting the scores, which highlights the disadvantage of the Copy baseline. Also, GPT3.5 and GPT4 consistently perform poorly in comparison to the rest of the models, which is especially surprising given GPT4’s multilingual capability. In addition, this could be an artifact of the metrics since RLHF can produce excellent results that have little overlap with the references, and future human studies may shed more light on the issue.

CLMs and Seq2Seq models perform similarly on GEC and Paraphrasing, while Seq2Seq performs better on simplification. We posit that this discrepancy happens due to shorter generations from the Seq2Seq models since we observe that the seq2seq models generate significantly shorter sequences than the expected distribution ($D = 0.06, p < 0.001$)⁵, increasing the BLEU scores, which are sensitive to the prediction length, whereas CLMs tend to generate longer sequences ($D = 0.27, p < 0.001$). Looking at the SARI scores, we observe that the two model types do not differ significantly ($p > .05$), indicating similar performance overall.

5.5 Effect of Model Scale

In Figure 8, we describe the overall performance of MTE models by size aggregated over the three tasks. It is evident that scaling model size generally increases overall performance significantly, thus reinforcing the effectiveness of model scaling. We also note that all three tasks display similar trends, with relative improvements being the greatest in GEC (28.8%) and Paraphrasing (14.8%).

⁵Using a two-tailed Kolmogorov-Smirnov test, compared against the distribution of lengths of the reference texts.

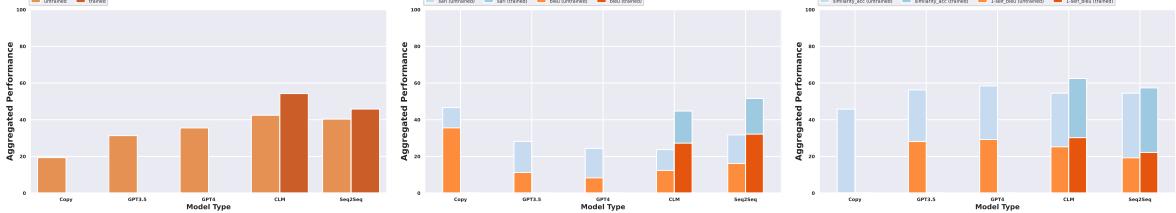


Figure 6: **Aggregated performance by model type.** For each task, we aggregate the relevant metrics as described in § 4.4 and split them by model type (CLM vs Seq2Seq), including the copy baseline.

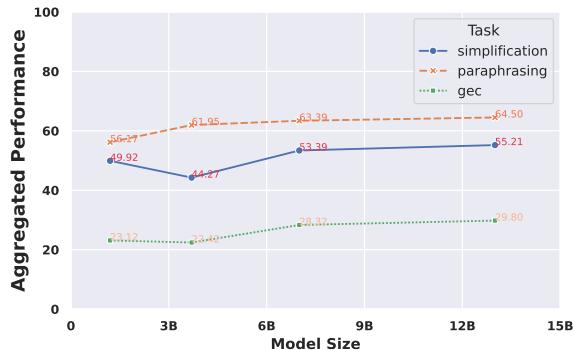


Figure 7: **Aggregated model performance on different tasks broken down by parameter size.** For visualization reasons, we group the 1.2B and 1.7B models and the 7B and 7.1B models together.

5.6 Effect of Task-specific Data

To understand the effect of task-specific data on model performance, we systematically ablate the proportion of training data for each task. Specifically, we conduct three groups by varying the amounts of training data across a given task between 0%, 10%, 50%, and 100% while keeping the amount of training data across other tasks at 100%, the results of which are shown in Figure 8. We observe that: (a) Performance on the ablated task generally improves as the amount of training data for that task increases, as expected. As the proportion of training data increases, so does the performance on the specific task. (b) We also note a synergistic relationship between some tasks where training data from one task helps improve performance on a different task. For example, as we add training examples from GEC, we also notice an improvement in model performance on simplification (58.90 vs 45.20) and paraphrasing (65.30 vs 58.34). Similar trends also hold when we add data for the other two tasks. We believe our model (which is inherently multi-task) enables us to leverage such synergy between text editing tasks better, as compared to task-specific models.

Task	Language	Dataset	Test Size	Metric	MEDIT	SOTA
GEC	Romanian	RoGEC	1518	GLEU	45.58	–
	Hindi	HiWikEd	13187	ERRANT GLEU	32.61 68.91	49.4 80
	Italian	Simpitiki PaCSSS-IT	176 63007	SARI BLEU	47.84 41.11	24.27 36.92
Simplification	Hindi	IndicSS	42771	SARI Rouge-L	40.08 19.92	45.57
	French	PAWS-X	903	Self-BLEU SA	69.06 98.38	–
Paraphrasing	Hindi	IndicPara.	10000	Self-BLEU SA iBLEU	23.91 92.06 14.2	– 18.55

Table 2: **Zero-shot evaluation results on the language generalization experiments.** We present the scores achieved by our best-performing model (**Our score**) along with the current **SOTA** results. Wherever possible, we report the metrics reported in the SOTA papers, and, if not available, we report commonly used ones by the literature. Note that we focus only on languages that the LLMs have seen during pre-training § 5.7.⁶

5.7 Generalization to new languages

We also explore the capabilities of our MEDIT models on new languages. For every task, we chose two new languages not present in the training set: one related to the language families covered in our training dataset (Table 1) and one belonging to a language family not present in the training dataset. We only considered the languages that the underlying LLM had as part of its pre-training corpora so as to ensure the models had some understanding of the languages in question. Table 2 provides a summary of the languages considered and the datasets they were sourced from, as well as the results of the language generalization experiments. We follow the same metrics for all the tasks as in § 4.4.

MEDIT models are competitive on many unseen languages as compared to the monolingual state-of-the-art, especially on it-Simplification, and hi-GEC and Paraphrasing.⁷

GEC: (Sonawane et al., 2020), IT-SIMP: (Tonelli et al., 2016; Brunato et al., 2016), HI-SIMP: (Kumar et al., 2022), FR-PARA: (Yang et al., 2019), HI-PARA: (Kumar et al., 2022). For datasets that contain more than 1k rows, we extract a subset of 1k rows for our generalization experiments.

⁷Since no multilingual models have attempted to perform

⁶Datasets include: RO-GEC: (Cotet et al., 2020), HI-

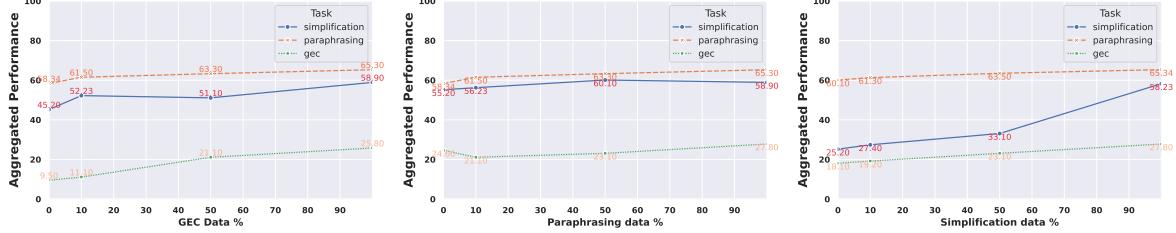


Figure 8: **Aggregated model performance by varying amounts of data samples:** (0% to 100%) by task (in the order: GEC, Paraphrasing, Simplification). We aggregate the scores as described in § 5.

	Language	V. Good	Good	Neutral	Bad	V. Bad
FLUENCY	Arabic	25.00	14.29	17.86	14.29	28.57
	Chinese	56.67	13.33	23.33	3.33	3.33
	English	56.67	23.33	13.33	3.33	3.33
	German	30.0	56.67	3.33	6.67	3.33
	Japanese	50.0	4.54	22.72	13.63	9.09
	Korean	39.13	21.74	17.39	13.04	8.70
ADEQUACY	Spanish	63.33	10.00	13.33	10.00	3.33
	Arabic	21.43	14.29	10.71	25.00	28.57
	Chinese	56.67	16.67	6.67	10	10.0
	English	62.33	18.32	9.09	9.09	1.16
	German	33.33	63.33	0.0	3.33	0.0
	Japanese	63.63	4.55	4.55	18.18	9.09
ACCURACY	Korean	41.67	16.67	12.50	20.83	8.33
	Spanish	60.0	6.67	6.67	13.33	13.33
	Arabic	21.43	3.57	10.71	35.71	28.57
	Chinese	3.45	13.79	17.24	51.72	13.79
	English	37.93	32.23	8.33	18.18	3.33
	German	30.0	40.0	23.33	6.67	0.0
	Japanese	34.48	10.34	3.45	20.69	31.03
	Korean	18.52	18.52	11.11	14.81	37.04
	Spanish	37.93	24.14	6.90	3.45	27.59

Table 3: **Results of the human evaluation of the model output across three criteria.** For each of the criteria, expert human annotators rate the system output, and we note the frequency of their rating (%).

6 Human Evaluations

We conduct human evaluations of our model outputs by proficient linguists (and native language speakers of the respective languages) on 50 test inputs (per task- language) to ensure they meet the instructional task-specific constraints across the various languages since text editing is often subjective, and automatic metrics are often limited in their effectiveness and accuracy. We evaluate our best-performing MEDIT model⁸ based on its strong performance (Table 5). Specifically, we conduct a qualitative evaluation where each annotator is shown an instructional input and output from the model and asked to rate the quality of the output on three criteria: **fluency** (Is the output grammatically correct and sound like it was written by a native speaker of the language?), **adequacy** (does the output preserve the meaning of the input?), and

the considered tasks in the respective languages, we are unable to report the metrics on others. We also do not provide any comparisons if the language-specific metrics either do not exist or are reported using different metrics.

⁸bactrian-x-llama-13b-merged

accuracy (did the model make the desired edits according to the given edit instructions?) of the edited texts, on a Likert scale ranging from *Very Bad* to *Very Good*. We collect two annotations for each data point and adjudicate the conflicting judgments with the annotators. The annotation guidelines are provided in Appendix E. Table 3 shows the results of the evaluation. The expert annotators generally rate the model outputs as *Good* or *Very Good* across nearly all languages. For Arabic, the preferences are more balanced across the scale and sometimes even leaning towards the *Bad* or *Very Bad* across all criteria, which confirms our findings on the automatic metrics as well, indicating that while our model performs very well on languages such as English, German, Chinese, and Spanish, it still has some way to go in terms of performance for languages such as Arabic in terms of quality, and on Japanese, Korean, and Spanish on accuracy.

7 Conclusions

We present MEDIT – an open-sourced dataset and set of multilingual instruction-tuned large language models capable of following natural language instructions in seven languages to perform various textual editing tasks. It is the first publicly available set of models that heavily outperforms numerous multilingual LLMs on multiple tasks in multiple languages. Positive feedback from human evaluations shows that MEDIT can assist writers with various aspects of the text revision process at scale by following natural language instructions in multiple languages. Experiments on various multilingual NLP tasks demonstrate that MEDIT models outperform both their corresponding non-fine-tuned variants and models fine-tuned on other multilingual instruction datasets for text editing, in addition to achieving strong performance on languages unseen during the task-specific fine-tuning. Making our data/models available, we hope to help advance multilingual intelligent writing assistants.

578 Limitations

579 Despite our efforts to develop and evaluate
580 instruction-tuned LLMs capable of text editing in
581 multiple languages, our work suffers from several
582 limitations that can be improved in future work.

583 First, although we have attempted to cover a di-
584 verse set of seven languages, there are still a large
585 number of languages that are not considered in our
586 work, primarily due to the lack of high-quality text
587 editing data. Future work can extend our system to
588 include more languages, especially low-resource
589 languages, to gain a more comprehensive under-
590 standing of the LLM generalization to new lan-
591 guages and building even more accessible and ubiq-
592 uitous writing assistants.

593 Secondly, we rely primarily on data available
594 in specific languages for training and evaluation,
595 and sometimes we generate instructions in English
596 and translate them into multiple languages using
597 Google Translate (for Simplification tasks, as an
598 example). In spite of the fact that our approach
599 enables us to extend to multiple languages with
600 affordable development costs, data generated and
601 translated might contain unexpected noise. More-
602 over, they might not give the best representation of
603 expert-annotated edits in different languages. In or-
604 der to further improve multilingual LLMs for text
605 editing, future research can use human-generated
606 data for training and evaluation.

607 Thirdly, our system leverages numerous LLMs
608 with billions of parameters. Considering the com-
609 puting resources required for running and develop-
610 ing these models, replicating the results may prove
611 difficult (which we try to address by sharing our
612 models publicly).

613 Finally, our evaluations only investigate the per-
614 formance of the models on benchmark datasets for
615 text editing, which focus on measuring superficial
616 characteristics based on n-gram overlaps. These
617 evaluations are limited as they do not test for the
618 more nuanced aspects of text editing, such as flu-
619 ency, coherence, and meaning preservation. The
620 lack of human evaluations also makes it difficult to
621 assess these nuanced characteristics. Future work
622 in this direction could look at robust and scalable
623 evaluation metrics for multilingual text editing.

624 Ethics Statement

625 While our models offer several advantages to make
626 intelligent writing assistance more accessible, we
627 do recognize their potential limitations. Since our

628 work mainly focuses on text editing, we are able
629 to avoid many issues involving generating harm-
630 ful text. Although there is still a possibility of
631 small meaning changes for stylistic tasks due to the
632 lack of user-specific context (Kulkarni and Raheja,
633 2023), we try to reduce the chance of hallucinations
634 by constraining the generation to strictly editing
635 tasks in order to reduce the chance of adding any
636 new information or perpetuating biases.

637 Moreover, due to the multilingual settings, there
638 is a risk of our models generating responses that are
639 discriminatory, biased, or contain false information.
640 Hence, our models, when fine-tuned on the text
641 editing datasets, may inadvertently learn or propa-
642 gate these problematic patterns. To address these
643 concerns and minimize potential harm, we are ded-
644 icated to mitigating the risks associated with the
645 use of our models in future research. We strongly
646 advocate for the responsible use of our models to
647 prevent any unintended negative consequences.

648 References

- Bashar Alhafni, Nizar Habash, and Houda Bouamor. 2022. *The Arabic parallel gender corpus 2.0: Extensions and analyses*. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 1870–1884, Marseille, France. European Language Resources Association.
- Bashar Alhafni, Go Inoue, Christian Khairallah, and Nizar Habash. 2023. *Advancements in Arabic grammatical error detection and correction: An empirical investigation*. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 6430–6448, Singapore. Association for Computational Linguistics.
- Marwah Alian, Arafat Awajan, Ahmad Al-Hasan, and Raeda Akuzhia. 2019. *Towards building arabic paraphrasing benchmark*. In *Proceedings of the Second International Conference on Data Science, E-Learning and Information Systems*, DATA ’19, New York, NY, USA. Association for Computing Machinery.
- Fernando Alva-Manchego, Louis Martin, Antoine Borde, Carolina Scarton, Benoît Sagot, and Lucia Specia. 2020. *ASSET: A dataset for tuning and evaluation of sentence simplification models with multiple rewriting transformations*. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4668–4679, Online. Association for Computational Linguistics.
- Fernando Alva-Manchego, Louis Martin, Carolina Scarton, and Lucia Specia. 2019. *EASSE: Easier automatic sentence simplification evaluation*. In *Proceedings of the 2019 Conference on Empirical Methods*

681	<i>in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP): System Demonstrations</i> , pages 49–54, Hong Kong, China. Association for Computational Linguistics.	740
682		741
683		742
684		743
685		744
686		745
687		746
688	Miriam Anschütz, Joshua Oehms, Thomas Wimmer, Bartłomiej Jezierski, and Georg Groh. 2023. <i>Language models for German text simplification: Overcoming parallel data scarcity through style-specific pre-training</i> . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 1147–1158, Toronto, Canada. Association for Computational Linguistics.	747
689		748
690		749
691		750
692		751
693		752
694		753
695	Anthony Baez and Horacio Saggion. 2023. <i>LSLLama: Fine-tuned LLaMA for lexical simplification</i> . In <i>Proceedings of the Second Workshop on Text Simplification, Accessibility and Readability</i> , pages 102–108, Varna, Bulgaria. INCOMA Ltd., Shoumen, Bulgaria.	754
696		755
697		756
698		757
699		758
700		759
701		760
702	Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. <i>A large annotated corpus for learning natural language inference</i> . In <i>Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing</i> , pages 632–642, Lisbon, Portugal. Association for Computational Linguistics.	761
703		762
704		763
705		764
706		765
707		766
708		767
709	Adriane Boyd. 2018. <i>Using Wikipedia edits in low resource grammatical error correction</i> . In <i>Proceedings of the 2018 EMNLP Workshop W-NUT: The 4th Workshop on Noisy User-generated Text</i> , pages 79–84, Brussels, Belgium. Association for Computational Linguistics.	768
710		769
711		770
712		771
713	Eleftheria Briakou, Di Lu, Ke Zhang, and Joel Tetreault. 2021. <i>Olá, bonjour, salve! XFORMAL: A benchmark for multilingual formality style transfer</i> . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 3199–3216, Online. Association for Computational Linguistics.	772
714		773
715		774
716		775
717		776
718		777
719		
720	Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020a. <i>Language models are few-shot learners</i> . In <i>Advances in Neural Information Processing Systems</i> , volume 33, pages 1877–1901. Curran Associates, Inc.	778
721		779
722		780
723		781
724		782
725		
726	Dominique Brunato, Andrea Cimino, Felice Dell’Orletta, and Giulia Venturi. 2016. <i>PaCCSS-IT: A parallel corpus of complex-simple sentences for automatic text simplification</i> . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 351–361, Austin, Texas. Association for Computational Linguistics.	783
727		784
728		785
729		
730	Christopher Bryant, Mariano Felice, Øistein E. Andersen, and Ted Briscoe. 2019. <i>The BEA-2019 shared task on grammatical error correction</i> . In <i>Proceedings of the Fourteenth Workshop on Innovative Use of NLP for Building Educational Applications</i> , pages 52–75, Florence, Italy. Association for Computational Linguistics.	786
731		787
732		788
733		
734	Christopher Bryant, Mariano Felice, and Ted Briscoe. 2017. <i>Automatic annotation and evaluation of error types for grammatical error correction</i> . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 793–805, Vancouver, Canada. Association for Computational Linguistics.	789
735		790
736		791
737		792
738		793
739		794
740	Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. 2017. <i>SemEval-2017 task 1: Semantic textual similarity multilingual and crosslingual focused evaluation</i> . In <i>Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)</i> , pages 1–14, Vancouver, Canada. Association for Computational Linguistics.	795
741		
742	Jishnu Ray Chowdhury, Yong Zhuang, and Shuyi Wang. 2022. Novelty controlled paraphrase generation with retrieval augmented conditional prompt tuning. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 36, pages 10535–10544.	796
743		797
744		798
745	Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. <i>arXiv preprint arXiv:2210.11416</i> .	799
746		800
747	Teodor-Mihai Cotet, Stefan Ruseti, and Mihai Dascalu. 2020. <i>Neural grammatical error correction for romanian</i> . In <i>2020 IEEE 32nd International Conference on Tools with Artificial Intelligence (ICTAI)</i> , pages 625–631.	801
748		802
749		803
750		804
751		805
752		806
753		807
754	Steven Coyne and Keisuke Sakaguchi. 2023. An analysis of gpt-3’s performance in grammatical error correction. <i>arXiv preprint arXiv:2303.14342</i> .	808
755		809
756		810
757		811
758	Yiming Cui, Ziqing Yang, and Xin Yao. 2023. <i>Efficient and effective text encoding for chinese llama and alpaca</i> . <i>arXiv preprint arXiv:2304.08177</i> .	812
759		813
760		814
761		815
762	Daniel Dahlmeier and Hwee Tou Ng. 2012. <i>Better evaluation for grammatical error correction</i> . In <i>Proceedings of the 2012 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 568–572, Montréal, Canada. Association for Computational Linguistics.	816
763		817
764		818
765		819
766		820
767		821
768		822
769		823
770		824
771		825
772		826
773		827
774		828
775		829
776		830
777		831

796	Sam Davidson, Aaron Yamada, Paloma Fernandez Mira, Agustina Carando, Claudia H. Sanchez Gutierrez, and Kenji Sagae. 2020. Developing NLP tools with a new corpus of learner Spanish. In <i>Proceedings of the Twelfth Language Resources and Evaluation Conference</i> , pages 7238–7243, Marseille, France. European Language Resources Association.	852
797		853
798		854
799		855
800		
801		
802		
803	Yupei Du, Qi Zheng, Yuanbin Wu, Man Lan, Yan Yang, and Meirong Ma. 2022. Understanding gender bias in knowledge base embeddings. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1381–1395, Dublin, Ireland. Association for Computational Linguistics.	856
804		857
805		858
806		859
807		860
808		861
809		
810	Tao Fang, Jinpeng Hu, Derek F. Wong, Xiang Wan, Lidia S. Chao, and Tsung-Hui Chang. 2023a. Improving grammatical error correction with multi-modal feature integration. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 9328–9344, Toronto, Canada. Association for Computational Linguistics.	862
811		863
812		864
813		865
814		866
815		867
816		868
817	Tao Fang, Shu Yang, Kaixin Lan, Derek F. Wong, Jinpeng Hu, Lidia S. Chao, and Yue Zhang. 2023b. Is chatgpt a highly fluent grammatical error correction system? a comprehensive evaluation.	869
818		870
819		871
820		872
821	Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2022. Language-agnostic BERT sentence embedding. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 878–891, Dublin, Ireland. Association for Computational Linguistics.	873
822		874
823		875
824		876
825		
826		
827		
828	Simon Flachs, Felix Stahlberg, and Shankar Kumar. 2021. Data strategies for low-resource grammatical error correction. In <i>Proceedings of the 16th Workshop on Innovative Use of NLP for Building Educational Applications</i> , pages 117–122, Online. Association for Computational Linguistics.	877
829		878
830		879
831		880
832		881
833		882
834	Nizar Habash and David Palfreyman. 2022. ZAEBUC: An annotated Arabic-English bilingual writer corpus. In <i>Proceedings of the Thirteenth Language Resources and Evaluation Conference</i> , pages 79–88, Marseille, France. European Language Resources Association.	883
835		884
836		
837		
838		
839	Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In <i>International Conference on Learning Representations</i> .	885
840		886
841		
842		
843		
844	Chao Jiang, Mounica Maddela, Wuwei Lan, Yang Zhong, and Wei Xu. 2020. Neural CRF model for sentence alignment in text simplification. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7943–7960, Online. Association for Computational Linguistics.	887
845		888
846		889
847		
848		
849		
850	Akihiro Katsuta and Kazuhide Yamamoto. 2018. Crowdsourced corpus of sentence simplification with core vocabulary. In <i>Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)</i> , Miyazaki, Japan. European Language Resources Association (ELRA).	890
851		891
852		892
853		893
854		894
855		895
856		
857		
858		
859		
860		
861		
862	Zae Myung Kim, Wanyu Du, Vipul Raheja, Dhruv Kumar, and Dongyeop Kang. 2022. Improving iterative text revision by learning where to edit from other revision tasks. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 9986–9999, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	896
863		897
864		898
865		899
866		
867		
868		
869	Aomi Koyama, Tomoshige Kiyuna, Kenji Kobayashi, Mio Arai, and Mamoru Komachi. 2020. Construction of an evaluation corpus for grammatical error correction for learners of Japanese as a second language. In <i>Proceedings of the Twelfth Language Resources and Evaluation Conference</i> , pages 204–211, Marseille, France. European Language Resources Association.	900
870		901
871		902
872		903
873		904
874		905
875		906
876		
877	Kalpesh Krishna, Deepak Nathani, Xavier Garcia, Bidisha Samanta, and Partha Talukdar. 2022. Few-shot controllable style transfer for low-resource multilingual settings. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 7439–7468, Dublin, Ireland. Association for Computational Linguistics.	907
878		908
879		909
880		910
881		911
882		912
883		913
884		
885	Vivek Kulkarni and Vipul Raheja. 2023. Writing assistants should model social factors of language.	914
886		915
887	Aman Kumar, Himani Shrotriya, Prachi Sahu, Amogh Mishra, Raj Dabre, Ratish Puduppully, Anoop Kunchukuttan, Mitesh M. Khapra, and Pratyush Kumar. 2022. IndicNLG benchmark: Multilingual datasets for diverse NLG tasks in Indic languages. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 5363–5394, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	916
888		917
889		918
890		919
891		920
892		921
893		922
894		923
895		
896	Philippe Laban, Jesse Vig, Marti A Hearst, Caiming Xiong, and Chien-Sheng Wu. 2023. Beyond the chat: Executable and verifiable text-editing with llms. <i>arXiv preprint arXiv:2309.15337</i> .	924
897		925
898		926
899		927
900	Huiyuan Lai, Antonio Toral, and Malvina Nissim. 2022. Multilingual pre-training with language and task adaptation for multilingual text style transfer. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 262–271, Dublin, Ireland. Association for Computational Linguistics.	928
901		929
902		930
903		931
904		932
905		933
906		
907	Viet Dac Lai, Nghia Trung Ngo, Amir Pouran Ben Veyseh, Hieu Man, Franck Dernoncourt, Trung Bui,	934
908		935

909	and Thien Huu Nguyen. 2023. Chatgpt beyond en-	Behrang Mohit, Alla Rozovskaya, Nizar Habash, Wa-	966
910	glish: Towards a comprehensive evaluation of large	jdi Zaghouani, and Ossama Obeid. 2014. <i>The first</i>	967
911	language models in multilingual learning.	<i>QALB shared task on automatic text correction</i>	968
912	<i>arXiv:2304.05613</i> .	<i>for Arabic</i> . In <i>Proceedings of the EMNLP 2014</i>	969
913		<i>Workshop on Arabic Natural Language Processing</i>	970
914		(ANLP)	971
915	Haonan Li, Fajri Koto, Minghao Wu, Alham Fikri Aji,	, pages 39–47, Doha, Qatar. Association for	971
916	and Timothy Baldwin. 2023. <i>Bactrian-x : A multi-</i>	Computational Linguistics.	972
917	<i>lingual replicable instruction-following model with</i>		
918	<i>low-rank adaptation</i> .		
919			
920	Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu	Niklas Muennighoff, Thomas Wang, Lintang Sutawika,	973
921	Wang, Shuhui Chen, Daniel Simig, Myle Ott, Na-	Adam Roberts, Stella Biderman, Teven Le Scao,	974
922	man Goyal, Shruti Bhosale, Jingfei Du, Ramakanth	M Saiful Bari, Sheng Shen, Zheng Xin Yong, Hai-	975
923	Pasunuru, Sam Shleifer, Punit Singh Koura, Vishrav	ley Schoelkopf, Xiangru Tang, Dragomir Radev,	976
924	Chaudhary, Brian O’Horo, Jeff Wang, Luke Zettle-	Alham Fikri Aji, Khalid Almubarak, Samuel Al-	977
925	moyer, Zornitsa Kozareva, Mona Diab, Veselin Stoy-	banie, Zaid Alyafeai, Albert Webson, Edward Raff,	978
926	anov, and Xian Li. 2022. <i>Few-shot learning with</i>	and Colin Raffel. 2023. <i>Crosslingual generaliza-</i>	979
927	<i>multilingual generative language models</i> . In <i>Proce-</i>	<i>tion through multitask finetuning</i> . In <i>Proceedings</i>	980
928	<i>of the 2022 Conference on Empirical Methods</i>	<i>of the 61st Annual Meeting of the Association for</i>	981
929	<i>in Natural Language Processing</i> , pages 9019–9052,	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	982
930	Abu Dhabi, United Arab Emirates. Association for	, pages 15991–16111, Toronto, Canada. Associa-	983
931	Computational Linguistics.	tion for Computational Linguistics.	984
932			
933	Guoqing Luo, Yu Han, Lili Mou, and Mauajama Firdaus.	Courtney Napoles, Keisuke Sakaguchi, Matt Post, and	985
934	2023. <i>Prompt-based editing for text style transfer</i> .	Joel Tetreault. 2015. <i>Ground truth for grammatical</i>	986
935	In <i>Findings of the Association for Computational</i>	<i>error correction metrics</i> . In <i>Proceedings of the 53rd</i>	987
936	<i>Linguistics: EMNLP 2023</i> , pages 5740–5750.	<i>Annual Meeting of the Association for Computational</i>	988
937		<i>Linguistics and the 7th International Joint Confer-</i>	989
938	Jonathan Mallinson, Jakub Adamek, Eric Malmi, and	<i>ence on Natural Language Processing (Volume 2:</i>	990
939	Aliaksei Severyn. 2022. <i>EdiT5: Semi-autoregressive</i>	<i>Short Papers</i> , pages 588–593, Beijing, China. Asso-	991
940	<i>text editing with t5 warm-start</i> . In <i>Findings of the</i>	ciation for Computational Linguistics.	992
941	<i>Association for Computational Linguistics: EMNLP</i>		
942	2022		
943	, pages 2126–2138, Abu Dhabi, United Arab		
944	Emirates. Association for Computational Linguistics.		
945	Jonathan Mallinson, Rico Sennrich, and Mirella Lapata.	Courtney Napoles, Keisuke Sakaguchi, Matt Post, and	993
946	2020. <i>Zero-shot crosslingual sentence simplification</i> .	Joel Tetreault. 2016. <i>Gleu without tuning</i> .	994
947	In <i>Proceedings of the 2020 Conference on Empirical</i>		
948	<i>Methods in Natural Language Processing (EMNLP)</i> ,		
949	pages 5109–5126, Online. Association for Compu-		
950	tational Linguistics.		
951			
952	Louis Martin, Angela Fan, Éric de la Clergerie, Antoine	Courtney Napoles, Keisuke Sakaguchi, and Joel	995
953	Bordes, and Benoît Sagot. 2022. <i>MUSS: Multilin-</i>	Tetreault. 2017. <i>JFLEG: A fluency corpus and bench-</i>	996
954	<i>gual unsupervised sentence simplification by mining</i>	<i>mark for grammatical error correction</i> . In <i>Proce-</i>	997
955	<i>paraphrases</i> . In <i>Proceedings of the Thirteenth Lan-</i>	<i>dings of the 15th Conference of the European Chapter</i>	998
956	<i>guage Resources and Evaluation Conference</i> , pages	<i>of the Association for Computational Linguistics:</i>	999
957	1651–1664, Marseille, France. European Language	<i>Volume 2, Short Papers</i> , pages 229–234, Valencia,	1000
958	Resources Association.	Spain. Association for Computational Linguistics.	1001
959			
960	Takumi Maruyama and Kazuhide Yamamoto. 2018.	Tong Niu, Semih Yavuz, Yingbo Zhou, Nitish Shirish	1002
961	<i>Simplified corpus with core vocabulary</i> . In <i>Pro-</i>	Keskar, Huan Wang, and Caiming Xiong. 2021. <i>Un-</i>	1003
962	<iceedings conference="" eleventh="" i="" international="" of="" on<="" the=""></iceedings>	<i>supervised paraphrasing with pretrained language</i>	1004
963	<i>Language Resources and Evaluation (LREC 2018)</i> ,	<i>models</i> . In <i>Proceedings of the 2021 Conference on</i>	1005
964	Miyazaki, Japan. European Language Resources As-	<i>Empirical Methods in Natural Language Processing</i> ,	1006
965	sociation (ELRA).	pages 5136–5150, Online and Punta Cana, Domini-	1007
966		cian Republic. Association for Computational Lin-	1008
967		guistics.	1009
968			
969	Tomoya Mizumoto, Mamoru Komachi, Masaaki Nagata,	OpenAI. 2023. <i>GPT-4 Technical Report</i> .	1010
970	and Yuji Matsumoto. 2011. <i>Mining revision log of</i>		
971	<i>language learning SNS for automated Japanese error</i>	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	1011
972	<i>correction of second language learners</i> . In <i>Proce-</i>	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	1012
973	<i>dings of 5th International Joint Conference on Natural</i>	Sandhini Agarwal, Katarina Slama, Alex Ray, et al.	1013
974	<i>Language Processing</i> , pages 147–155, Chiang Mai,	2022. Training language models to follow instruc-	1014
975	Thailand. Asian Federation of Natural Language Pro-	<i>tions with human feedback</i> . <i>Advances in Neural</i>	1015
976	cessing.	<i>Information Processing Systems</i> , 35:27730–27744.	1016
977			
978			
979			
980			
981			
982			
983			
984			
985			
986			
987			
988			
989			
990			
991			
992			
993			
994			
995			
996			
997			
998			
999			
1000			
1001			
1002			
1003			
1004			
1005			
1006			
1007			
1008			
1009			
1010			
1011			
1012			
1013			
1014			
1015			
1016			
1017			
1018			
1019			
1020			
1021			

1022	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. <i>Journal of Machine Learning Research</i> , 21(140):1–67.	International Conference on Learning Representations.	1079
1023			1080
1024			
1025			
1026			
1027			
1028	Vipul Raheja, Dhruv Kumar, Ryan Koo, and Dongyeop Kang. 2023. Coedit: Text editing by task-specific instruction tuning. <i>arXiv preprint arXiv:2305.09857</i> .		
1029			
1030			
1031	Emily Reif, Daphne Ippolito, Ann Yuan, Andy Coenen, Chris Callison-Burch, and Jason Wei. 2022. A recipe for arbitrary text style transfer with large language models. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 837–848, Dublin, Ireland. Association for Computational Linguistics.		
1032			
1033			
1034			
1035			
1036			
1037			
1038	Nils Reimers and Iryna Gurevych. 2020. Making monolingual sentence embeddings multilingual using knowledge distillation. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 4512–4525, Online. Association for Computational Linguistics.		
1039			
1040			
1041			
1042			
1043			
1044	Sascha Rothe, Jonathan Mallinson, Eric Malmi, Sebastian Krause, and Aliaksei Severyn. 2021. A simple recipe for multilingual grammatical error correction. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)</i> , pages 702–707, Online. Association for Computational Linguistics.		
1045			
1046			
1047			
1048			
1049			
1050			
1051			
1052	Alla Rozovskaya, Houda Bouamor, Nizar Habash, Wajdi Zaghouani, Ossama Obeid, and Behrang Mohit. 2015. The second QALB shared task on automatic text correction for Arabic. In <i>Proceedings of the Second Workshop on Arabic Natural Language Processing</i> , pages 26–35, Beijing, China. Association for Computational Linguistics.		
1053			
1054			
1055			
1056			
1057			
1058			
1059	Michael Ryan, Tarek Naous, and Wei Xu. 2023. Revisiting non-English text simplification: A unified multilingual benchmark. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 4898–4927, Toronto, Canada. Association for Computational Linguistics.		
1060			
1061			
1062			
1063			
1064			
1065			
1066	Horacio Saggion, Sanja Štajner, Daniel Ferrés, Kim Cheng Sheang, Matthew Shardlow, Kai North, and Marcos Zampieri. 2022. Findings of the TSAR-2022 shared task on multilingual lexical simplification. In <i>Proceedings of the Workshop on Text Simplification, Accessibility, and Readability (TSAR-2022)</i> , pages 271–283, Abu Dhabi, United Arab Emirates (Virtual). Association for Computational Linguistics.		
1067			
1068			
1069			
1070			
1071			
1072			
1073			
1074	Victor Sanh, Albert Webson, Colin Raffel, Stephen H Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, et al. 2022. Multitask prompted training enables zero-shot task generalization. In <i>ICLR 2022-Tenth International Conference on Learning Representations</i> .		
1075			
1076			
1077			
1078			
1079	Timo Schick, Jane A. Yu, Zhengbao Jiang, Fabio Petroni, Patrick Lewis, Gautier Izacard, Qingfei You, Christoforos Nalmpantis, Edouard Grave, and Sebastian Riedel. 2023. PEER: A collaborative language model. In <i>The Eleventh International Conference on Learning Representations</i> .		
1080			
1081	Haitham Seelawi, Ahmad Mustafa, Hesham Al-Bataineh, Wael Farhan, and Hussein T. Al-Natsheh. 2019. NSURL-2019 task 8: Semantic question similarity in Arabic. In <i>Proceedings of the First International Workshop on NLP Solutions for Under Resourced Languages (NSURL 2019) co-located with ICNLSP 2019 - Short Papers</i> , pages 1–8, Trento, Italy. Association for Computational Linguistics.		
1082			
1083			
1084			
1085	Laura Seiffe, Fares Kallel, Sebastian Möller, Babak Naderi, and Roland Roller. 2022. Subjective text complexity assessment for German. In <i>Proceedings of the Thirteenth Language Resources and Evaluation Conference</i> , pages 707–714, Marseille, France. European Language Resources Association.		
1086			
1087	Ankur Sonawane, Sujeet Kumar Vishwakarma, Bhavana Srivastava, and Anil Kumar Singh. 2020. Generating inflectional errors for grammatical error correction in Hindi. In <i>Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing: Student Research Workshop</i> , pages 165–171, Suzhou, China. Association for Computational Linguistics.		
1088			
1089	Xin Sun, Tao Ge, Shuming Ma, Jingjing Li, Furu Wei, and Houfeng Wang. 2022. A unified strategy for multilingual grammatical error correction with pre-trained cross-lingual language model. In <i>Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22</i> , pages 4367–4374. International Joint Conferences on Artificial Intelligence Organization. Main Track.		
1090			
1091			
1092			
1093			
1094			
1095	Toshikazu Tajiri, Mamoru Komachi, and Yuji Matsumoto. 2012. Tense and aspect error correction for ESL learners using global context. In <i>Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 198–202, Jeju Island, Korea. Association for Computational Linguistics.		
1096			
1097			
1098			
1099			
1100			
1101	Yasuhiro Tanaka. 2001. Compilation of a multilingual parallel corpus. <i>Proceedings of PACLING 2001</i> , pages 265–268.		
1102			
1103			
1104			
1105			
1106			
1107			
1108			
1109			
1110	Sara Tonelli, Alessio Palmero Aprilio, and Francesca Saltori. 2016. Simpitiki: a simplification corpus for italian. <i>Proceedings of the Third Italian Conference on Computational Linguistics (CLiC-it) 2016</i> , pages 291–296.		
1111			
1112			
1113			
1114			
1115			
1116			
1117	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro,		
1118			
1119			
1120			
1121			
1122			
1123			
1124			
1125			
1126			
1127			
1128			
1129			
1130			
1131			
1132			
1133			
1134			
1135			

1136	Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .	1190
1137		1191
1138		1192
1139		1193
1140	Xiangpeng Wei, Haoran Wei, Huan Lin, Tianhao Li, Pei Zhang, Xingzhang Ren, Mei Li, Yu Wan, Zhiwei Cao, Binbin Xie, Tianxiang Hu, Shangjie Li, Binyuan Hui, Bowen Yu, Dayiheng Liu, Baosong Yang, Fei Huang, and Jun Xie. 2023. Polym : An open source polyglot large language model.	1194
1141		1195
1142		1196
1143		1197
1144		1198
1145	Sam Witteveen and Martin Andrews. 2019. Paraphrasing with large language models. In <i>Proceedings of the 3rd Workshop on Neural Generation and Translation</i> , pages 215–220, Hong Kong. Association for Computational Linguistics.	1199
1146		1200
1147		1201
1148		1202
1149		1203
1150	BigScience Workshop. 2023. Bloom: A 176b-parameter open-access multilingual language model.	1204
1151		1205
1152	Haoran Wu, Wenxuan Wang, Yuxuan Wan, Wenxiang Jiao, and Michael Lyu. 2023. Chatgpt or grammarly? evaluating chatgpt on grammatical error correction benchmark.	1206
1153		1207
1154		
1155		
1156	Wei Xu, Chris Callison-Burch, and Courtney Napoles. 2015. Problems in current text simplification research: New data can help. <i>Transactions of the Association for Computational Linguistics</i> , 3:283–297.	1208
1157		1209
1158		1210
1159		1211
1160	Wei Xu, Courtney Napoles, Ellie Pavlick, Quanze Chen, and Chris Callison-Burch. 2016a. Optimizing statistical machine translation for text simplification. <i>Transactions of the Association for Computational Linguistics</i> , 4:401–415.	1212
1161		1213
1162		1214
1163		1215
1164		
1165	Wei Xu, Courtney Napoles, Ellie Pavlick, Quanze Chen, and Chris Callison-Burch. 2016b. Optimizing statistical machine translation for text simplification. <i>Transactions of the Association for Computational Linguistics</i> , 4:401–415.	1216
1166		1217
1167		1218
1168		1219
1169		1220
1170	Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mT5: A massively multilingual pre-trained text-to-text transformer. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 483–498, Online. Association for Computational Linguistics.	1221
1171		
1172		
1173		
1174		
1175		
1176		
1177		
1178	Shiping Yang, Renliang Sun, and Xiaojun Wan. 2023a. A new dataset and empirical study for sentence simplification in Chinese. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8306–8321, Toronto, Canada. Association for Computational Linguistics.	1228
1179		1229
1180		1230
1181		1231
1182		1232
1183		1233
1184		1234
1185	Wen Yang, Chong Li, Jiajun Zhang, and Chengqing Zong. 2023b. Bigtranslate: Augmenting large language models with multilingual translation capability over 100 languages. <i>arXiv preprint arXiv:2305.18098</i> .	1240
1186		1241
1187		1242
1188		1243
1189		
1190	Yinfei Yang, Daniel Cer, Amin Ahmad, Mandy Guo, Jax Law, Noah Constant, Gustavo Hernandez Abrego, Steve Yuan, Chris Tar, Yun-hsuan Sung, Brian Strope, and Ray Kurzweil. 2020. Multilingual universal sentence encoder for semantic retrieval. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations</i> , pages 87–94, Online. Association for Computational Linguistics.	1244
1191		1245
1192		1246
1193		
1194		
1195		
1196		
1197		
1198		
1199	Yinfei Yang, Yuan Zhang, Chris Tar, and Jason Baldridge. 2019. PAWS-X: A cross-lingual adversarial dataset for paraphrase identification. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3687–3692, Hong Kong, China. Association for Computational Linguistics.	1247
1200		1248
1201		1249
1202		1250
1203		1251
1204		1252
1205		1253
1206		1254
1207		1255
1208	Soyoung Yoon, Sungjoon Park, Gyuwan Kim, Junhee Cho, Kihyo Park, Gyu Tae Kim, Minjoon Seo, and Alice Oh. 2023. Towards standardizing Korean grammatical error correction: Datasets and annotation. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 6713–6742, Toronto, Canada. Association for Computational Linguistics.	1256
1209		1257
1210		1258
1211		1259
1212		1260
1213		1261
1214		1262
1215		1263
1216	Xingxing Zhang and Mirella Lapata. 2017. Sentence simplification with deep reinforcement learning. In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 584–594, Copenhagen, Denmark. Association for Computational Linguistics.	1264
1217		1265
1218		1266
1219		1267
1220		1268
1221		1269
1222	Ying Zhang, Hidetaka Kamigaito, and Manabu Okumura. 2023. Bidirectional transformer reranker for grammatical error correction. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 3801–3825, Toronto, Canada. Association for Computational Linguistics.	1270
1223		1271
1224		1272
1225		1273
1226		1274
1227		1275
1228	Yuan Zhang, Jason Baldridge, and Luheng He. 2019. PAWS: Paraphrase adversaries from word scrambling. In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 1298–1308, Minneapolis, Minnesota. Association for Computational Linguistics.	1276
1229		1277
1230		1278
1231		1279
1232		1280
1233		1281
1234		1282
1235		1283
1236	Yuanyuan Zhao, Nan Jiang, Weiwei Sun, and Xiaojun Wan. 2018. Overview of the nlpcc 2018 shared task: Grammatical error correction. In <i>Natural Language Processing and Chinese Computing</i> .	1284
1237		1285
1238		1286
1239		1287
1240	Wenhao Zhu, Yunzhe Lv, Qingxiu Dong, Fei Yuan, Jingjing Xu, Shujian Huang, Lingpeng Kong, Jiajun Chen, and Lei Li. 2023. Extrapolating large language models to non-english by aligning languages.	1288
1241		1289
1242		1290
1243		1291
1244	Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, and Yong Yu. 2018. Texxygen: A benchmarking platform for text generation models.	1292
1245		1293
1246		1294

1247
1248
1249
1250
In *The 41st International ACM SIGIR Conference on*
Research & Development in Information Retrieval,
SIGIR ’18, page 1097–1100, New York, NY, USA.
Association for Computing Machinery.

1251 A Training and Evaluation Datasets

1252 A.1 Grammatical Error Correction

1253 Arabic We report on three publicly available
1254 Arabic GEC datasets. The first two come from the
1255 QALB-2014 (Mohit et al., 2014) and QALB-2015
1256 (Rozovskaya et al., 2015) shared tasks. The third is
1257 the newly created ZAEBUC dataset (Habash and
1258 Palfreyman, 2022; Alhafni et al., 2023). QALB-
1259 2014 consists of native/L1 user comments from
1260 the Aljazeera news website, whereas QALB-2015
1261 consists of essays written by Arabic L2 learners
1262 with various levels of proficiency. It is worth noting
1263 that the QALB-2015 dataset has two test sets con-
1264 sisting of L1 and L2 data. In this work, we report
1265 results on the L1 test set. The ZAEBUC dataset
1266 comprises essays written by native Arabic speak-
1267 ers, which were manually corrected. We use the
1268 MaxMatch (M^2) Scorer (Dahlmeier and Ng, 2012)
1269 for the evaluation.

1270 English When it comes to English, we use the
1271 Write & Improve + LOCNESS (W&I) corpus re-
1272 leased in the Building Educational Applications
1273 (BEA) shared task on GEC (Bryant et al., 2019).
1274 We also use the NAIST Lang-8 corpus (Tajiri et al.,
1275 2012), which is one of the largest and most widely
1276 used datasets for English GEC. To test our systems,
1277 we use the JFLEG (Napoles et al., 2017) dataset.
1278 We use GLEU (Napoles et al., 2015, 2016) for the
1279 evaluation.

1280 German For German, we use the Falko-
1281 MERLIN corpus (Boyd, 2018), which consists of
1282 sentences written by L2 learners that were manu-
1283 ally corrected. We use the MaxMatch (M^2) Scorer
1284 (Dahlmeier and Ng, 2012) for the evaluation.

1285 Spanish For Spanish, we use the publicly avail-
1286 able COWS-L2H (Davidson et al., 2020) dataset.
1287 COWS-L2H consists of essays written by Spanish
1288 L2 learners at the university level in the United
1289 States. We use ERRANT (Bryant et al., 2017) for
1290 the evaluation.

1291 Chinese For Chinese, we use the data that is part
1292 of the NLPCC18 shared task (Zhao et al., 2018).
1293 The training data used in the shared task was col-
1294 lected from the NAIST Lang-8 corpus (Tajiri et al.,
1295 2012), whereas the test data consists of manually

1296 corrected sentences written by Chinese L2 learners.
1297 We use GLEU (Napoles et al., 2015, 2016) for the
1298 evaluation.

1299 Japanese For Japanese, we use the NAIST Lang-
1300 8 corpus (Mizumoto et al., 2011) to train our sys-
1301 tems. For evaluation, we use the Japanese L2 TEC-
1302 JL dataset (Koyama et al., 2020). We use GLEU
1303 (Napoles et al., 2015, 2016) for the evaluation.

1304 Korean We use the recently created Kor-Union
1305 dataset (Yoon et al., 2023). Kor-Union was cre-
1306 ated by collecting and combining GEC data from
1307 various sources. This includes essays written by
1308 Korean native/L1 speakers and L2 learners. We use
1309 the MaxMatch (M^2) Scorer (Dahlmeier and Ng,
1310 2012) for the evaluation.

1311 A.2 Paraphrasing

1312 Arabic For training, we use the Arabic SemEval
1313 Paraphrasing (ASEP) corpus, which sourced three
1314 existing Arabic semantic similarity datasets re-
1315 leased during SemEval 2017 Task 1 (Cer et al.,
1316 2017), consisting of roughly 1100 sentence pairs.
1317 For our purposes, similar to them, we only keep the
1318 sentence pairs with a semantic similarity score \geq
1319 3.25, which leads to 603 pairs. We also inverted the
1320 pairs for training, leading to a total of 1.2k training
1321 pairs. For evaluation, we use the evaluation dataset
1322 that was used for SemEval 2017 Track 1, but with
1323 the same similarity threshold as the training data,
1324 consisting of 67 sentence pairs. This evaluation set
1325 consists of sentences from the SNLI Corpus (Bow-
1326 man et al., 2015) that were human-translated into
1327 Arabic, provided by CMU-Qatar by native Arabic
1328 speakers with strong English skills.

1329 We also source from the Arabic Question Simi-
1330 larity (Shared Task 8) organized at the Workshop
1331 on NLP Solutions for Under-Resourced Languages
1332 (NSURL 2019) (Seelawi et al., 2019). The dataset
1333 was developed by mawdoo, and consists of 12k
1334 pairs for training and 3715 for testing. For both
1335 training and evaluation, we filter the semantically
1336 similar pairs (similarity score of 1), which leaves
1337 us with 10.7k training and 1.7k test pairs.

1338 We also use the Arabic Paraphrasing Benchmark
1339 (APB) dataset (Alian et al., 2019), which consists
1340 of 1010 Arabic sentence pairs that are collected
1341 from different Arabic books. Paraphrasing was per-
1342 formed manually using six transformation proce-
1343 dures (i.e., addition, deletion, expansion, permuta-
1344 tion, reduction, and replacement). Similar to other
1345 evaluation sets, we only keep the sentence pairs

Task	Language	Dataset	Split	Size
Grammatical Error Correction	en	BEA (Bryant et al., 2019) JFLEG (Napoles et al., 2017)	Train Dev, Test	1.1M 754, 747
	ar	QALB-2014 (Mohit et al., 2014) QALB-2015 (Rozovskaya et al., 2015) ZAEBCU (Habash and Palfrayman, 2022)	Train, Dev, Test	19k, 1k, 968 310, 154, 920 150, 33, 31
	de	Falko-MERLIN (Boyd, 2018)	Train, Dev, Test	19k, 2.5k, 2.3k
	es	COWS-L2H (Davidson et al., 2020)	Train, Dev, Test	398, 85, 86
	ja	Lang-8 (Mizumoto et al., 2011) TEC-JL (Koyama et al., 2020)	Train Test	1.85M 1.9k
	ko	Kor-Union (Yoon et al., 2023)	Train, Dev, Test	108.9k, 23.3k, 23.3k
Paraphrasing	zh	NLPCC-2018 (Zhao et al., 2018)	Train, Dev, Test	540k, 53.5k, 2k
	en	PAWS (Zhang et al., 2019)	Train, Dev, Test	49k, 8k, 8k
	ar	SemEval 2017 - Task 1 (Cer et al., 2017) NSURL 2019 - Task 8 (Seelawi et al., 2019) APB (Alian et al., 2019)	Train, Test Test	1.2k, 67 24k, 3.7k 286
	de			
	es			
	fr			
Simplification	ja	PAWS-X (Yang et al., 2019)	Train, Dev, Test	49k, 2k, 2k
	ko			
	zh			
	en	WikiLarge (Zhang and Lapata, 2017) WikiAuto (Jiang et al., 2020) NEWSLEA (Xu et al., 2015) ASSET (Alva-Manchego et al., 2020)	Train, Dev, Test Train Dev, Test	296k, 2k, 359 576k, 5k, 5k 94k 2000, 359
	ar	NEWSLEA-Auto-AR ASSET-Auto-AR	Train Dev, Test	94k 100, 359
	de	GEOlino (Mallinson et al., 2020) TextComplexityDE (Seiffe et al., 2022)	Train, Dev, Test	958, 122, 118 200, 25, 25
1346	es	NEWSLEA-Auto-ES ASSET-Auto-ES	Train Dev, Test	94k 100, 359
	ja	EasyJapanese (Maruyama and Yamamoto, 2018) EasyJapanese Extended (Katsuta and Yamamoto, 2018)	Train, Dev, Test Train, Test	48k, 1k, 1k 34k, 731
	ko	NEWSLEA-Auto-KO ASSET-Auto-KO	Train Dev, Test	94k 100, 359
	zh	NEWSLEA-Auto-ZH CSS (Yang et al., 2023a)	Train Dev, Test	94k , 383

Table 4: **Datasets used to train and evaluate MEDIT.** With the exceptions of Spanish GEC and German Simplification, every other dataset contains >10k examples for all our experiments.

with a semantic similarity score ≥ 3.25 , which leads to 286 pairs.

English Paraphrase Adversaries from Word Scrambling (PAWS) is a dataset that contains pairs of sentences with a high lexical overlap (Zhang et al., 2019). We use the PAWS dataset for training and evaluation.

German, Spanish, Japanese, Korean, Chinese We use the Cross-lingual Paraphrase Adversaries from Word Scrambling (Yang et al., 2019) dataset (PAWS-X), which was created by translating a subset of the PAWS validation and test sets to six other languages by professional translators.

A.3 Simplification

We draw on a variety of existing text simplification datasets in various languages. Table 4 shows the different simplification datasets we draw on in our

work and also outlines the training, development, and test settings.

A major issue with text simplification is the absence of publicly available, human-annotated, sentence-level parallel corpora for some of the languages we considered, such as Arabic, Spanish, and Korean. Therefore, we addressed this by translating the Text Simplification datasets for English to these three languages, in which the parallel data is absent. One potential limitation of this approach could be the poor quality of the translation models, which could negatively impact the overall data quality. Therefore, we use the latest Google Translate API⁹ to construct the translated data, and further verify the quality of the translated text with human annotators (native speakers) for a subset of the data. We chose the Google API since it performed best

⁹<https://cloud.google.com/translate/docs/advanced/translating-text-v3>

1380 amongst the other open-source machine translation
1381 models and APIs we tested.¹⁰

1382 **English** For English, we used Wikilarge (Zhang
1383 and Lapata, 2017), WikiAuto (Jiang et al., 2020),
1384 and Newsela (Xu et al., 2015) datasets for train-
1385 ing. WikiAuto is a neural CRF-aligned corpus of
1386 original and simple Wikipedia documents that are
1387 automatically aligned to generate sentence pairs,
1388 whereas the Newsela (Xu et al., 2015) dataset con-
1389 tains automatically aligned sentence pairs from doc-
1390 uments that are generated by professional writers
1391 at Newsela for various grade levels. For testing, we
1392 use ASSET (Alva-Manchego et al., 2020), which
1393 contains ten high-quality human written simplifi-
1394 cations for each of the 2,390 sentences from the
1395 TurkCorpus (Xu et al., 2016b).

1396 **German** We use the GEOLino (Mallinson et al.,
1397 2020) and TextComplexityDE (Seiffe et al., 2022)
1398 datasets for both training and testing. GEOLinoTest
1399 contains text about nature, physics, and
1400 people from GeoLino, a children’s magazine that
1401 was manually simplified by a German linguist to
1402 a five to seven-year-old reading level. TextCom-
1403 plexityDE contains 250 complex sentences from
1404 German Wikipedia that native speakers manually
1405 simplified.

1406 **Japanese** We use the EasyJapanese (Maruyama
1407 and Yamamoto, 2018) and EasyJapaneseExtended
1408 (Maruyama and Yamamoto, 2018) datasets for
1409 training and testing. EasyJapanese contains 50k
1410 sentence pairs that were manually created by five
1411 students by simplifying text from the Tanaka cor-
1412 pus (Tanaka, 2001). The EasyJapaneseExtended
1413 dataset contains an additional 34.4k sentences
1414 from the Tanaka corpus with simplifications crowd-
1415 sourced.

1416 **Arabic, Spanish, Korean** For Arabic, Spanish
1417 and Korean, as there were no publicly available
1418 sentence-level parallel datasets available, we trans-
1419 lated the English simplification datasets. Specif-
1420 ically, we translated the English Newsela dataset for
1421 training and ASSET for testing using the Google
1422 Translate API, giving us 94k and 359 examples for
1423 training and testing, respectively.

1424 **Chinese** We found no publicly available dataset
1425 for training Chinese Simplification. Therefore,
1426 we again translated the English Newsela training
1427 dataset into Chinese. However, for the testing set,

¹⁰<https://libretranslate.com/>

1428 we use the CSS (Yang et al., 2023a) dataset. CSS
1429 consists of two human-written simplifications for
1430 each of the 383 original sentences from the PFR
1431 corpus.¹¹

B Data Preparation

1432 For Seq2Seq models, we prepend the task-specific
1433 instructions to the input to the encoder for each lan-
1434 guage, performing a full parameter update on the
1435 entire sequence. We construct the CLM datasets
1436 by wrapping each example in model-specific in-
1437 structions¹² computing the loss *only* on the target
1438 text; hence, in the Native and Random settings, the
1439 model does not optimize for the specific “translated”
1440 instructions.

1441 We randomly sample 10k examples from the
1442 original datasets for each language-task combina-
1443 tion and keep the original validation and test tests
1444 (see Table 4). We chose this quantity as a bal-
1445 ance between computational cost and qualitative
1446 performance based on the insights from (Raheja
1447 et al., 2023). Moreover, in our experiments, we
1448 did not find a significant impact of increasing the
1449 data quantity per task-language combination be-
1450 yond 10k. In the Spanish GEC task, we only have
1451 398 data points, so this portion of our data is con-
1452 siderably smaller than the rest of the languages.
1453 Furthermore, in the GEC and Simplification train-
1454 ing data for all languages, we reserve 20% of the
1455 set as input-output pairs without any edits to avoid
1456 over-corrections by the model.

C Results

1457 We present the full set of results for our best-
1458 performing random-language setting. In addition to
1459 the marginally higher performance that all trained
1460 models demonstrate in this setting, we choose this
1461 as our representative setting due to the fact that it is
1462 the most capable setting for our models, allowing
1463 them to perform cross-lingual editing.

1464 We present the results for all trained models on
1465 all datasets, as well as the No Edits (Copy) baseline
1466 and State-of-the-art LLMs. In the interest of space
1467 and interpretability, we skip the detailed results for
1468 the untrained baseline as we already show them
1469 to be massively underperformant relative to the
1470 trained models. We also compare our results to the

¹¹<https://www.heywhale.com/mw/dataset>

¹²PolyLM and Bactrian-X follow different prompt tem-
1471 plates. Details in Appendix D.

1473 previously reported SOTA results across tasks and
1474 languages:

1475 **Grammatical Error Correction** For GEC, we
1476 compare against the following SOTA results: Ara-
1477 bic (Alhafni et al., 2022), English (Zhang et al.,
1478 2023), German (Fang et al., 2023a), Korean (Yoon
1479 et al., 2023), Spanish (Flachs et al., 2021), and
1480 Japanese (Koyama et al., 2020).

1481 **Simplification** For simplification, we take the
1482 results from the fine-tuned mT5 models from (Ryan
1483 et al., 2023). They compare multiple settings in
1484 their work, such as using data from a single training
1485 dataset, from a single language, and from multiple
1486 languages. We pick the score from any setting that
1487 gives the highest score. Thus, the SARI score for
1488 German might be picked from one setting while the
1489 BLEU score for German from some other setting,
1490 and so on. This ensures that we take their strongest
1491 models for each use case. We still see that our
1492 models perform better than them in most languages
1493 and datasets.

1494 **Paraphrasing** To paraphrase, most of the related
1495 works that utilize the PAWS-X dataset (Yang et al.,
1496 2019) have used it for paraphrase detection.

1497 D Task Verbalizers

1498 We present the full list of our manually curated
1499 task-specific verbalizers used for training and eval-
1500 uations in Table 6, Table 7, and Table 8.

1501 E Human Evaluation Annotation 1502 Guidelines

1503 The human experts were asked to rate the **fluency**,
1504 **adequacy**, and **accuracy** separately, following the
1505 guidelines in Figure 9, Figure 10 and Figure 11,
1506 respectively.

Given the instruction to edit the text and the text to be edited, rate the output on the following dimensions:

Fluency: is the output valid, free from errors, and like how a native speaker of the language would write?
Please use the following scale for your evaluations along every dimension:

Very Good: The output is of high quality, valid, and correct, like a native speaker.

Good: The output is acceptable with minor errors.

Average: The output is relevant but has significant errors.

Bad: The output is subpar quality.

Very Bad: The output is unusable.

Figure 9: Annotation guidelines for human evaluations for Fluency.

Given the instruction to edit the text and the text to be edited, rate the output on the following dimensions:

Adequacy: does the output preserve the meaning of the original sentence?

Please use the following scale for your evaluations along every dimension:

Very Good: The output fully preserves the meaning of the text.

Good: The output is semantically similar to the input with minor errors.

Average: The output is semantically similar to the input but has significant errors.

Bad: The output is barely similar to the input.

Very Bad: The output has opposite meaning to the input.

Figure 10: Annotation guidelines for human evaluations for Adequacy.

Given the instruction to edit the text and the text to be edited, rate the output on the following dimensions:

Accuracy: how well do the edits made in the output follow the given instructions?

Please use the following scale for your evaluations along every dimension:

Very Good: The output follows the instructions exactly.

Good: The output generally follows the instructions with minor errors.

Average: The instructions are partially followed but has errors.

Bad: The output did not follow the instructions but did not significantly make the text unusable.

Very Bad: The instructions were completely ignored, and wrong edits were made to make the text unusable.

Figure 11: Annotation guidelines for human evaluations for Accuracy.

Model	Version	Size	en	de	es	ja	ko	zh	ar		
			JFLEG	Falko-MERLIN	CowsL2H	TEC-JL	Ko-Union	NLPCC-18	QALB-14	QALB-15-L1	ZAEBUC
Copy	–	–	40.47	35.23	31.36	49.09	48.2	56.89	1.91	0.07	0.0
GPT3.5	gpt-3.5-turbo@613	–	39.54	39.41	38.89	38.56	27.5	12.5	28.08	35.01	45.17
GPT4	gpt-4-@613	–	35.43	40.23	38.23	40.55	29.34	11.05	41.78	44.55	47.9
Multilingual SOTA	–	–	61.97	76.3	57.32	73.6	31.70	–	79.6	80.3	83.1
mT5	mt5-large mt5-xl mt5-xxl	1.2B 3.7B 13B	36.28 40.72 41.56	12.98 40.21 51.4	40.89 52.98 39.68	12.98 26.33 39.68	35.17 36.14 37.11	50.76 52.45 55.56	64.76 68.36 67.83	64.93 68.53 67.31	68.56 68.95 67.23
mT0 / Bloomz	mt0-large bloomz-3b mt0-xl mt0-xxl	1.2B 3.7B 3.7B 13B	38.25 7.25 40.3 40.65	32.32 4.71 39.6 40.75	43.39 31.20 49.19 52.31	9.14 11.35 10.22 14.14	24.04 21.33 36.53 37.06	51.54 52.45 52.25/ 52.96 52.96	64.74 66.21 67.45 68.15	64.34 65.68 67.18 67.84	69.78 76.42 65.98 66.85
mT0 / Bloomz (mt)	bloomz-7b1-mt mt0-xxl-mt	1.2B 13B	7.1B 37.58	30.67 41.75	9.91 53.23	33.13 46.35	10.91 57.54	24.53 53.67	30.15 70.19	68 70.19	66.65 69.95
PolyLM	polylm-1.7b	1.7B	38.35	35.45	46.87	43.22	22.3	57.78	54.1	51.5	51.26
Bactrian-X	bx-llama-7b bx-llama-13b	7B 13B	58.67 59.55	60.07 64.11	45.32 49.0	47.1 53.41	25.56 26.66	56.09 67.54	64.42 69.92	64.03 70.13	67.63 71.45

Model	Version	Size	en	de	es	ja	ko	zh	ar		
			PAWS	PAWS-X	PAWS-X	PAWS-X	PAWS-X	PAWS-X	NSURL	ASEP	APP
Copy	–	–	0.0 / 100.0	0.0 / 100.0	0.0 / 100.0	0.0 / 100.0	0.0 / 100.0	0.0 / 100.0	0.0 / 100.0	0.0 / 100.0	0.0 / 100.0
GPT3.5	gpt-3.5-turbo@613	–	40.85 / 98.78	41.49 / 97.8	45.17 / 97.98	82.97 / 42.88	82.97 / 42.88	82.38 / 42.19	88.2 / 41.33	88.57 / 45.7	99.55 / 58.94
GPT4	gpt-4-@613	–	36.84 / 97.2	48.5 / 95.96	44.09 / 95.96	83.45 / 48.38	44.81 / 91.47	46.60 / 68.49	38.98 / 77.23	80.0 / 37.23	78.80 / 49.36
mT5	mt5-large mt5-xl mt5-xxl	1.2B 3.7B 13B	25.77 / 100.00	39.72 / 97.29	24.83 / 96.54	35.80 / 91.13	28.01 / 88.76	36.79 / 94.87	59.94 / 93.79	38.61 / 86.83	14.81 / 68.25
mT0 / Bloomz	bloomz-3b mt0-xl mt0-xxl	1.2B 3.7B 3.7B 13B	27.63 / 100.00 32.99 / 96.31 32.56 / 98.28 45.39 / 100.00	34.01 / 97.22 43.15 / 91.59 30.17 / 98.18 42.74 / 98.27	40.44 / 88.58 44.04 / 89.60 46.31 / 92.05 41.66 / 98.22	32.35 / 95.81 32.35 / 95.81 33.68 / 89.62 52.68 / 93.12	45.92 / 95.84 45.92 / 95.84 38.79 / 94.34 50.57 / 92.04	57.68 / 88.85 57.68 / 88.85 53.22 / 94.96 57.68 / 88.85	64.51 / 86.83 64.51 / 86.83 35.43 / 86.89 74.99 / 94.45	65.19 / 86.89 65.19 / 86.89 17.26 / 68.25 71.14 / 86.89	63.95 / 68.25 63.95 / 68.25 36.98 / 68.25 36.98 / 68.25
mT0 / Bloomz (mt)	bloomz-7b1-mt mt0-xxl-mt	1.2B 13B	41.30 / 100.00 42.51 / 100.00	43.02 / 97.88 42.23 / 98.12	42.23 / 98.12 43.67 / 91.12	48.20 / 89.76 48.20 / 89.76	54.58 / 95.82 54.01 / 95.82	70.62 / 94.74 70.62 / 94.74	70.14 / 86.89 70.14 / 86.89	67.37 / 68.25 50.67 / 68.25	
PolyLM	polylm-1.7b	1.7B	47.57 / 100.00	43.15 / 94.74	29.33 / 86.89	35.96 / 68.25	30.22 / 98.21	28.32 / 98.27	75.32 / 92.05	77.41 / 89.66	83.04 / 95.82
Bactrian-X	bx-llama-7b bx-llama-13b	7B 13B	51.91 / 100.00 53.49 / 100.00	50.07 / 98.32 51.50 / 100.00	46.67 / 98.11 48.91 / 98.21	55.86 / 93.11 58.47 / 94.23	54.88 / 89.66 59.88 / 89.66	54.18 / 95.82 55.22 / 95.89	75.37 / 94.74 76.17 / 94.74	94.76 / 86.89 93.04 / 86.89	94.60 / 68.25 93.42 / 68.25

Model	Version	Size	en	ar	es	de	ja	ko	zh			
			ASSET	WikiAuto	ar-ASSET	es-ASSET	GeoLine	TCD-E	EasyJ	EasyIE	ko-ASSET	CS
Copy	–	–	20.73 / 92.81	20.93 / 45.40	17.91 / 86.75	21.17 / 92.77	27.45 / 69.86	15.42 / 26.77	29.66 / 75.91	22.00 / 48.47	16.45 / 82.32	29.27 / 90.42
GPT3.5	gpt-3.5-turbo@613	–	38.69 / 53.53	38.99 / 22.00	36.90 / 20.21	43.17 / 51.85	27.81 / 14.77	38.04 / 10.04	20.35 / 12.01	27.88 / 6.35	38.10 / 18.95	21.82 / 15.23
GPT4	gpt-4-@613	–	39.74 / 46.04	39.64 / 19.55	36.77 / 17.76	40.41 / 35.97	24.28 / 9.05	38.43 / 8.47	15.35 / 5.52	26.32 / 4.96	35.81 / 9.84	18.73 / 9.34
Multilingual SOTA	–	–	42.77 / 88.26	42.48 / 37.95	– / –	– / –	50.75 / 71.9	41.15 / 24.53	70.95 / 68.12	53.49 / 35.67	– / –	– / –
mT5	mt5-large mt5-xl mt5-xxl	1.2B 3.7B 13B	33.10 / 91.04 31.01 / 92.13 32.78 / 91.65	37.30 / 43.17 37.16 / 46.27 38.06 / 44.90	36.60 / 76.04 38.44 / 78.93 38.93 / 75.76	35.60 / 90.62 37.82 / 88.01 39.47 / 81.96	40.96 / 59.02 31.07 / 73.12 50.93 / 70.65	31.49 / 22.88 31.86 / 29.27 52.89 / 71.70	41.85 / 74.93 65.78 / 79.19 32.68 / 26.44	31.23 / 49.24 65.38 / 59.36 60.44 / 61.30	33.08 / 76.20 36.75 / 73.42 38.33 / 66.93	32.63 / 39.32 32.63 / 39.32 31.95 / 33.22
mT0 / Bloomz	mt0-large bloomz-3b mt0-xl mt0-xxl	1.2B 3.7B 3.7B 13B	30.28 / 92.81 29.92 / 60.20 29.63 / 90.96 32.78 / 91.65	34.08 / 46.62 21.21 / 30.50 34.70 / 46.58 38.93 / 88.54	35.37 / 93.12 21.28 / 66.15 36.99 / 77.40 39.83 / 85.54	44.82 / 69.98 27.88 / 44.05 47.07 / 68.70 50.93 / 70.65	27.84 / 26.20 30.69 / 26.90 60.62 / 78.02 33.92 / 27.29	40.06 / 75.04 30.01 / 19.55 50.58 / 54.03 68.22 / 79.69	25.84 / 47.18 17.35 / 34.56 33.71 / 77.04 61.63 / 62.77	29.18 / 80.14 29.41 / 69.21 32.71 / 41.33 37.51 / 67.01	34.04 / 53.82 29.41 / 69.21 32.71 / 41.33 31.95 / 33.22	
mT0 / Bloomz (mt)	bloomz-7b1-mt mt0-xxl-mt	1.2B 13B	20.88 / 66.85 20.92 / 60.20	21.09 / 34.24 21.21 / 30.50	17.98 / 60.62 17.99 / 57.87	21.24 / 71.10 21.28 / 66.15	27.79 / 47.49 27.88 / 44.05	15.60 / 20.43 15.73 / 18.40	29.81 / 37.77 30.01 / 19.55	22.38 / 27.17 23.01 / 17.98	17.20 / 54.66 17.35 / 34.56	29.33 / 74.98 29.41 / 69.21
PolyLM	polylm-1.7b	1.7B	21.12 / 22.33	21.22 / 11.54	18.00 / 18.23	21.43 / 22.77	27.50 / 13.77	15.49 / 7.94	29.69 / 6.15	22.50 / 4.58	16.86 / 19.84	29.53 / 21.18
Bactrian-X	bx-llama-7b bx-llama-13b	7B 13B	41.05 / 90.79 41.63 / 91.63	43.88 / 46.96 44.19 / 47.26	40.76 / 74.33 41.75 / 73.97	43.57 / 89.80 44.03 / 88.59	56.33 / 64.73 63.77 / 70.28	41.00 / 27.42 40.13 / 25.44	67.47 / 74.64 68.31 / 74.88	55.06 / 48.20 56.89 / 49.89	43.30 / 58.35 39.91 / 70.50	41.12 / 48.11

Table 5: Full set of results on the best-performing setting of 10k random-language-prompted data. For GEC, we report GLEU or $F_{0.5}$ depending on the metric as described in Appendix A. For Paraphrasing, the first quantity is 1 – Self-BLEU and the second one is the accuracy of semantic similarity as calculated by m USE (explained in § 4.4). Finally, for Simplification, the first quantity is SARI, and the second one is BLEU. In terms of models, bx-llama-7b denotes the MBZUAI/bactrian-x-llama-7b-merged checkpoint, and bx-llama-7b denotes the MBZUAI/bactrian-x-llama-13b-merged one.

Language Verbalizers

Arabic	اصلح القواعد النحوية في هذه الجملة	اصلح القواعد النحوية في هذه الجملة	اصلح القواعد النحوية في الجملة
	اصلح الأخطاء النحوية	اصلح الأخطاء النحوية	اصلح الأخطاء النحوية في هذه الجملة
	اصلح جميع الأخطاء النحوية في هذه الجملة	اصلح الأخطاء النحوية في هذه الجملة	اصلح القواعد النحوية للجملة
	اصلح الأخطاء النحوية في هذه الجملة	اصلح القواعد النحوية في هذه الجملة	اجعل الجملة نحوية
	اصلح التناقضات النحوية في الجملة	اصلح الجملة نحوية	اجعل جميع الأخطاء النحوية من هذا النص
	اصلح الأخطاء في هذا النص	حدث لإزالة الأخطاء النحوية	حسن القواعد النحوية لهذا النص
	حسن قواعد هذا النص	حسن القواعد	ازل الأخطاء النحوية
	حسن القواعد النحوية لهذه الجملة	تحسينات نحوية	اصلح الأخطاء النحوية
	ازل الأخطاء النحوية	اصلح الأخطاء النحوية	
Chinese	修复语法	修正这句话的语法	修复句子中的语法
	修复语法错误	修正语法错误	修复语法
	修正所有语法错误	修正这句话中的语法错误	修正这句话中的语法错误
	修正这句话中的语法错误	修正这句话的语法问题	修正句子的语法
	修正句子中的不连贯之处	使句子符合语法	使句子流畅
	修正本文中的错误	更新以删除语法错误	删除本文中的所有语法错误
	改进本文的语法	提高语法性	提高文本的语法性
	提高这句话的语法性	语法改进	删除语法错误
English	Fix grammar	Fix grammar in this sentence	Fix grammar in the sentence
	Fix grammar errors	Fix grammatical errors	Fix grammaticality
	Fix all grammatical errors	Fix grammatical errors in this sentence	Fix grammar errors in this sentence
	Fix grammatical mistakes in this sentence	Fix grammaticality in this sentence	Fix grammaticality of the sentence
	Fix disfluencies in the sentence	Make the sentence grammatical	Make the sentence fluent
	Fix errors in this text	Update to remove grammar errors	Remove all grammatical errors from this text
	Improve the grammar of this text	Improve the grammaticality	Improve the grammaticality of this text
	Remove grammatical mistakes	Grammar improvements	Remove grammar mistakes
German	Grammatik korrigieren	Grammatik in diesem Satz korrigieren	Grammatik im Satz korrigieren
	Grammatikfehler beheben	Grammatikfehler beheben	Grammatikalität korrigieren
	Alle Grammatikfehler beheben	Grammatikfehler in diesem Satz korrigieren	Grammatikfehler in diesem Satz korrigieren
	Grammatikfehler in diesem Satz korrigieren	Grammatikalität in diesem Satz korrigieren	Grammatikalität des Satzes korrigieren
	Unstimmigkeiten im Satz beheben	Machen Sie den Satz fließend	Fehler in diesem Text beheben
	Update zum Entfernen von Grammatikfehlern	Entfernen Sie alle Grammatikfehler aus diesem Text	Verbessern Sie die Grammatik dieses Textes
	Verbessern Sie die Grammatikalität	Verbessern Sie die Grammatikalität dieses Textes	Verbessern Sie die Grammatikalität dieses Satzes
	Grammatikverbesserungen	Grammatikfehler entfernen	Grammatikfehler entfernen
Japanese	文法を修正する	この文の文法を修正してください	文内の文法を修正してください
	文法エラーを修正	文法上の誤りを修正してください	文法を修正する
	文法上の誤りをすべて修正してください	この文の文法上の誤りを修正してください	この文の文法性を修正する
	この文の文法上の間違いを修正してください	この文の文法を修正してください	文章を流暢にしましょう
	文章の不一致を修正する	文章を文法的にする	このテキストから文法上の誤りをすべて削除してください
	このテキストのエラーを修正してください	文法エラーを削除するためには更新	このテキストの文法性を改善してください
	このテキストの文法を改善してください	文法を改善する	文法の改善
	この文の文法性を改善してください	文法の改善	文法上の間違いを取り除く
Korean	문법수정	이문장의 문법수정	문장의 문법수정
	문법오류수정	문법오류수정	문법수정
	모든문법오류수정	이문장의 문법오류수정	이문장의 문법오류수정
	이문장의 문법오류수정	이문장의 문법수정하십시오	문장의 문법수정
	문장에서수정	문장을 문법적으로 만드십시오	문장을 유창하게 만드십시오
	이텍스트의오류수정	문법오류를 제거하기 위한 업데이트	이텍스트에서 모든 문법 오류 제거
	이텍스트의문법개선	문법성 향상	이텍스트의 문법성을 개선하십시오
	이문장의문법성을 개선하십시오	문법 개선	문법 오류 제거
Spanish	Corregir gramática	Corrigir la gramática en esta oración	Arreglar la gramática en la oración
	Corregir errores gramaticales	Corregir errores gramaticales	Corregir la gramaticalidad
	Corregir todos los errores gramaticales	Corregir errores gramaticales en esta oración	Corregir errores gramaticales en esta oración
	Corregir errores gramaticales en esta oración	Corrigir la gramaticalidad en esta oración	Corregir la gramaticalidad de la oración
	Corregir disfluencias en la oración	Haz que la oración sea gramatical	Haz que la oración sea fluida
	Corregir errores en este texto	Actualizar para eliminar errores gramaticales	Eliminar todos los errores gramaticales de este texto
	Mejorar la gramática de este texto	Mejorar la gramaticalidad	Mejorar la gramaticalidad de este texto
	Mejorar la gramaticalidad de esta oración	Mejoras gramaticales	Eliminar errores gramaticales
Eliminar errores gramaticales	Corrige los errores gramaticales	Corrige los errores gramaticales	

Table 6: **Grammatical Error Correction instruction verbalizers.** For every language, we craft 27 GEC-specific instructions, increasing their diversity when the model is trained. For this and subsequent tables, we verify the validity of the instructions with native language speakers (§ 4.3).

Language	Verbalizers		
Arabic	<p>بسط الجملة</p> <p>أكتب نسخة أبسط للجملة</p> <p>أعد كتابة هذه الجملة من أجل التبسيط</p> <p>اجعل الجملة أبسط</p> <p>بسط</p> <p>بسط هذه الفقرة</p> <p>اجعل هذا أبسط لفهمهم</p>	<p>بسط هذه الجملة</p> <p>أعد كتابة الجملة لتكون أبسط</p> <p>أعد كتابة هذا بصيغة أبسط</p> <p>اجعل هذا النص أقل تمقيداً</p> <p>تبسيط</p> <p>بسط هذا النص</p>	
Chinese	<p>简化句子</p> <p>为该句子写一个更简单的版本</p> <p>为简单起见重写这句话</p> <p>让句子变得更容易理解</p>	<p>简化这句话</p> <p>将句子改写得更简单</p> <p>用更简单的措辞重写这个</p> <p>让这段文字不那么复杂</p>	<p>简化这段文字</p> <p>用更简单的方式重写这句话</p> <p>让句子变得简单</p> <p>让这件事变得更简单</p> <p>改为更简单的措辞</p> <p>使用更简单的措辞</p>
English	<p>Simplify the sentence</p> <p>Write a simpler version for the sentence</p> <p>Rewrite this sentence for simplicity</p> <p>Make the sentence simpler</p> <p>Simplify</p> <p>Simplify this paragraph</p> <p>Make this easier to understand</p>	<p>Simplify this sentence</p> <p>Rewrite the sentence to be simpler</p> <p>Rewrite this with simpler wording</p> <p>Make this text less complex</p> <p>Simplification</p> <p>Simplify this text</p>	<p>Simplify this text</p> <p>Rewrite this sentence in a simpler manner</p> <p>Make the sentence simple</p> <p>Make this simpler</p> <p>Change to simpler wording</p> <p>Use simpler wording</p>
German	<p>Vereinfachen Sie den Satz</p> <p>Schreiben Sie eine einfachere Version des Satzes</p> <p>Formulieren Sie diesen Satz der Einfachheit halber um</p> <p>Machen Sie den Satz einfacher</p> <p>Vereinfachen</p> <p>Vereinfachen Sie diesen Absatz</p> <p>Machen Sie es verständlicher</p>	<p>Vereinfachen Sie diesen Satz</p> <p>Formulieren Sie den Satz um, damit er einfacher ist</p> <p>Formulieren Sie dies mit einer einfacheren Formulierung um</p> <p>Machen Sie diesen Text weniger komplex</p> <p>Vereinfachung</p> <p>Vereinfachen Sie diesen Text</p>	<p>Vereinfachen Sie diesen Text</p> <p>Formulieren Sie diesen Satz einfacher um</p> <p>Machen Sie den Satz einfach</p> <p>Machen Sie es einfacher</p> <p>Änderung zu einer einfacheren Formulierung</p> <p>Verwenden Sie einfachere Formulierungen</p>
Japanese	<p>文を簡略化してください</p> <p>文のより簡単なバージョンを書いてください</p> <p>わかりやすくするためにこの文を書き直してください</p> <p>文をもっと簡単にしてください</p> <p>簡略化する</p> <p>この段落を簡略化してください</p> <p>これをもっとわかりやすくしてください</p>	<p>この文を簡単にしてください</p> <p>文章をもっと簡単に書き直してください</p> <p>これをもっと簡単な表現で書き直してください</p> <p>このテキストをより複雑にしないでください</p> <p>単純化</p> <p>このテキストを簡略化してください</p>	<p>このテキストを簡略化してください</p> <p>この文をもっと簡単に書き直してください</p> <p>文章を簡単にしてしまう</p> <p>これをもっとシンプルにしてください</p> <p>より簡単な表現に変更します</p> <p>より簡単な表現を使用してください</p>
Korean	<p>문장을간단하게</p> <p>문장의간단한번작성</p> <p>단순화를위해이문장을다시작성하십시오</p> <p>문장을간단하게만드십시오</p> <p>단순화</p> <p>이단락을단순화</p> <p>이해하기쉽게해주세요</p>	<p>이문장을단순화하십시오</p> <p>문장을더간단하게다시쓰세요</p> <p>간단한표현으로다시작성</p> <p>이텍스트를덜복잡하게만들기</p> <p>단순화</p> <p>이텍스트를단순화</p>	<p>이텍스트를단순화</p> <p>이문장을더간단한방식으로다시작성하십시오</p> <p>문장을간단하게만드십시오</p> <p>이것을더간단하게만드십시오</p> <p>간단한문구로변경</p> <p>간단한표현사용</p>
Spanish	<p>Simplifica la oración</p> <p>Escribe una versión más simple para la oración</p> <p>Reescribe esta oración para simplificar</p> <p>Hacer la oración más simple</p> <p>Simplificar</p> <p>Simplificar este párrafo</p> <p>Haz que esto sea más fácil de entender</p>	<p>Simplifica esta oración</p> <p>Reescribe la oración para que sea más simple</p> <p>Reescribe esto con una redacción más simple</p> <p>Hacer este texto menos complejo</p> <p>Simplificación</p> <p>Simplificar este texto</p>	<p>Simplificar este texto</p> <p>Reescribe esta oración de una manera más simple</p> <p>Haz la oración simple</p> <p>Haz esto más simple</p> <p>Cambiar a una redacción más simple</p> <p>Usar una redacción más simple</p>

Table 7: Simplification instruction verbalizers. For the simplification task, we generate 19 instructions per language, taking care to not change the meaning of the instruction. For more information see § 4.3 and Table 6.

Language	Verbalizers	Arabic	Chinese	English	German	Japanese	Korean	Spanish									
	أعد صياغة الجملة اخرج النص أكتب إعاده صياغة مختلفة استخدم صياغة مختلفة أعد صياغة هذه الجملة أعد صياغة هذا النص	أعد صياغة هذه الجملة اكتب إعاده صياغة لجملة استخدم صياغة مختلفة أعد صياغة هذه الجملة أعد صياغة هذا النص	解释一下句子 释义 用不同的措辞重写句子 改写这句话 重写此文本	解释一下这句话 为这句话写一个解释 使用不同的措辞 改写这句话 改写这段文字	Paraphrase the sentence Paraphrase Rewrite the sentence with different wording Reword this sentence Reword this text	Paraphrase this sentence Write a paraphrase for the sentence Use different wording Rephrase this sentence Rephrase this text	Umschreiben Sie den Satz Paraphrase Schreiben Sie den Satz mit einem anderen Wortlaut um Formulieren Sie diesen Satz um Diesen Text umformulieren	Umschreiben Sie diesen Satz Schreiben Sie eine Paraphrase für den Satz Andere Formulierungen verwenden Formulieren Sie diesen Satz um Formulieren Sie diesen Text neu	Paraphrase this text Write a paraphrased version of the sentence Rewrite this sentence Rewrite this text	Umschreiben Sie diesen Text Schreiben Sie eine paraphrasierte Version des Satzes Schreiben Sie diesen Satz um Diesen Text umschreiben	文を言い換える 言い換え 別の表現で文章を書き直してください この文を言い直してください このテキストを書き直してください	この文を言い換えてください 文の言い換えを書いてください 別の表現を使用する この文を言い換えてください このテキストを言い換えてください	문장을의역 의역 다른단어로문장을다시작성 이문장을바꾸십시오 이텍스트를바구십시오	이문장은의역 문장에대한의역쓰기 다른문구사용 이문장을바꿔보세요 이텍스트를다음말로바꿔보세요	Paraphrasear la oración Paráfrasis Reescribir la oración con una redacción diferente Reformular esta oración Reformular este texto	Parafaseando esta oración Escribe una paráfrasis de la oración Usar una redacción diferente Reformula esta oración Reformular este texto	أعد صياغة هذا النص أكتب نسخة معاد صياغتها من الجملة أعد كتابة هذه الجملة أعد كتابة هذا النص

Table 8: Paraphrasing instruction verbalizers. In this case, we create 14 instructions per language so as to not alter the meaning of the task. For more information see § 4.3 and Table 6.