# (Dis)improved?! How Simplified Language Affects Large Language Model Performance across Languages

Anonymous ACL submission

## Abstract

002 Simplified language enhances the accessibility and human understanding of texts. However, whether it also benefits large language models (LLMs) remains underexplored. This paper extensively studies whether LLM performance improves on simplified data compared to its 007 original counterpart. Our experiments span six datasets and eight automatic simplification systems across three languages. We show that 011 English models, including GPT-4o-mini, show a weak generalization and exhibit a significant performance drop on simplified data. This in-013 troduces an intriguing paradox: simplified data is helpful for humans but not for LLMs. At the same time, the performance in non-English languages sometimes improves, de-017 pending on the task and quality of the simplifier. Our findings offer a comprehensive view 019 of the impact of simplified language on LLM performance and uncover severe implications for people depending on simple language.

# 1 Introduction

037

041

Automatic Text Simplification (ATS) is the task of rewriting a text using simpler vocabulary while preserving its original meaning. The goal is to increase readability and make information accessible to a broader audience. The primary target group is people with low literacy and mental disabilities, or language learners (Martin et al., 2022). However, previous work has shown that not only people from the target group but even the broad majority of people profit from simplified language (Javourey-Drevet et al., 2022; Murphy Odo, 2022). With this paper, we try to answer if the same holds true for Large Language Models (LLMs). Given that LLMs are approaching human-like capabilities (Grattafiori et al., 2024), it is reasonable to hypothesize that they might also perform better with simplified input or at least show good performance and generalization on this language style.



Figure 1: Text sample from the Sentiment Analysis for Financial News dataset (Malo et al., 2014). We test the generalization of LLMs like Llama3.1 70B from original to automatically simplified data. The sentiment prediction on the original data sample is correct. However, if we use an automatic lexical simplifier that replaced the word "operates" with "works", Llama misclassifies the sample as positive.

043

044

045

047

051

057

060

061

062

063

064

065

To investigate this, we select six labeled datasets across three languages (English, German, and Russian) and simplify their texts using eight pretrained simplification models and LLMs. Then, we benchmark five large language models, including Llama3.1 (Grattafiori et al., 2024), Aya Expanse (Dang et al., 2024), and GPT-4o-mini, on both the original and simplified corpora. Our results show a significant change in performance with a strong performance drop for English (see example in Figure 1). This lack of generalization introduces a severe risk for people who rely on simplified language: If they input prompts or samples in simple language, LLMs may show a worse performance and make more mistakes than with standard English. Especially for tasks with high societal impact, like fake news classification or news summarization, this increases discrimination for already vulnerable target groups. Overall, our contributions can be summarized as follows:

• We present a large-scale multilingual benchmark of LLM generalization on simplified data, including s.o.t.a. models like Llama3.1, Aya Expanse, and GPT-40-mini. The sim-

- 071
- 074

081

084

087

096 097

100 101

103

104 105

108

109

110

111

112

113

114

fying the data at inference time results in a weaker performance than the original data. These results are in contrast to other studies that benchmarked simplification as inference preprocessing. Miyata and Tatsumi (2019) tested Google Translator for Japanese to English translations with sentence splitting and further rule-based simplifications. A human evaluation showed that the simplifications yielded strong improvements 102

plifications are evaluated on a broad range of metrics, covering readability and meaning

• Our results indicate a significant performance

• All code, simplified data, and model predic-

The impact of ATS on NLP tasks has been studied

for many years and for different NLP tasks (Vickrey and Koller, 2008; Schmidek and Barbosa, 2014;

Štajner and Popovic, 2016). However, many of the

older studies could not use transformers or even

large language models and were based on statistical

simplification. Among the more recent studies, we

identify two research directions: text simplification

as data augmentation for pre-training or fine-tuning

and text simplification as a pre-processing step to

improve inference performance. To investigate the

first direction, Van et al. (2021) simplify the train-

ing data for LSTM- and BERT-based classification

models and evaluate the simplification quality with

BLEU only. Results show that different setups of

data augmentation with simplification can improve

the classifiers. However, they also show that simpli-

tions are publicly available for further investi-

decline on English simplified data, but with promising improvements in non-English lan-

preservation, and a human review.

gation and experimentation<sup>1</sup>.

guages.

**Related work** 

2

in the translation outputs. Similarly, Mehta et al. (2020) created an artificial simplification system through back translation and used this system to simplify the machine translation inputs of a lowresource-language translation system. They show improved translation quality across multiple languages. However, the performance changes of the target systems depend on the quality of the ATS systems. As such, Agrawal and Carpuat (2024) investigated how well ATS systems preserve the meaning of the original texts. While human simplifications could improve the performance of a



Figure 2: Structure of our investigations. We compare the performance of the same model between the original inputs and their simplified versions. Red boxes indicate that these parts are investigated under different settings.

pre-trained question-answering model, automatic simplifications worsened the performance. Our work tries to shed light on the contradicting findings of previous work. For this, we extend the existing research by covering more tasks, languages, and simplifiers. We paint a broader picture of the helpfulness of simplification as pre-processing, especially in times of flexible and powerful LLMs.

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

140

141

142

143

144

145

146

147

148

149

A different research direction was chosen by Anschütz et al. (2024), who used human-supervised simplification corpora to investigate how well models generalize between original and simplified data. They are the first ones to include LLMs in their investigations and show that models exhibit an incoherent behavior between original and simplified data. However, they only benchmarked GPT3.5turbo as LLM, and their datasets do not contain ground-truth labels. While they assumed that the human-supervised datasets contain correct simplifications, they cannot measure the actual performance of the classification system without groundtruth labels. We try to overcome this weakness by using labeled datasets and benchmarking the performance of multiple LLMs on these datasets. In addition, we extend the investigation to the task of summarization and not only cover classification tasks.

#### 3 Methodology

Our objective is to compare if the performance of different LLMs changes when the input samples are simplified. For this, we take labeled datasets and simplify the inputs with existing simplifiers. Then, we use pre-trained classification models or LLMs to predict the labels on the original and on the simplified inputs. Finally, we calculate the ac-

<sup>&</sup>lt;sup>1</sup>URL removed for review

curacy and examine whether text simplification at 150 inference can improve the models' performance. 151 An overview of our approach is shown in Figure 2. 152 Our investigations cover three distinct languages 153 with six different datasets, eight simplifiers, and six 154 prediction models, including LLMs like GPT40. 155 All combinations were evaluated independently, 156 and the models did not know if the input text was 157 simplified or not to avoid bias. The different set-158 tings will be discussed in the following subsections. 159

### **3.1** Datasets and tasks

160

161

162

163

165

166

167

168

169

172

173

174

175

176

177

178

We cover the tasks of classification and summarization. The evaluation of text generation is non-trivial since nuances of text and language characteristics need to be covered. In contrast, comparing classification labels is independent of the chosen metric. In addition, ATS systems may struggle to preserve the exact meaning (Säuberli et al., 2024; Agrawal and Carpuat, 2024). Classification tasks like reading comprehension and natural language inference focus on specific text details that can get lost during simplification (Trienes et al., 2024), even though the simplification is of high overall quality. To avoid depending on these details, we focus on more content-related tasks like topic and sentiment prediction. We assume that even if the simplifiers remove minor aspects, the overall content should not change significantly, and thus, the ground-truth labels are still correct for the simplified samples.

The selected datasets are shown in Table 1. We 179 experiment with data in English, German, and Russian. All datasets are from the news domain, a 181 general-purpose domain often targeted by ATS lit-182 183 erature (Ryan et al., 2023). For each of the datasets, we only worked with the test splits. To reduce the fi-184 nancial efforts of the OpenAI API, we created fixed 185 subsets of the AG News and the sentiment dataset and only used these subsets when prompting this 187 API. In the following, results that are based on these subsets are indicated with <sup>†</sup>. Each language 189 contains a multi-task dataset that provides data for 190 topic classification and summarization at the same time to enable a multi-task evaluation. The num-192 ber of classes and granularity of the classes differ 193 among the languages and tasks. The AG News 194 dataset has four very general classes, while the 195 196 TL;DR dataset focuses more on technical news and its subcategories. For the sentiment task, we pur-197 posefully selected a dataset with only three classes 198 (positive, negative, and neutral) to avoid ambiguity. The summarization task is headline generation, 200

where the models create a headline for the respective news snippet. This task has a strongly abstractive nature and is well-suited to evaluate how well the models can retrieve the most important information from the texts (Scialom et al., 2020).

201

202

203

204

206

207

208

209

210

211

212

213

214

215

216

217

218

219

221

222

223

224

225

226

227

228

229

230

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

# 3.2 Simplifiers

We used eight different pre-trained simplification models for our experiments: two multilingual models for all languages and six language-specific models (four for English, one for German, and one for Russian). Our model selection was limited by the availability and reproducibility of existing approaches. Especially unmaintained or weaklydocumented Github repositories make reusing pretrained models challenging (Stodden, 2024; Kew et al., 2023). Nevertheless, the models that we could run give a good variety of approaches, ranging from lexical to paragraph-level simplification, and are trained for general-purpose or specialized domains. For all models, we used the default configurations provided in their repositories or model cards, and we did not add any further preprocessing. We used these simplification models:

**MILES (multiling.)** is a lexical simplification pipeline. It uses frequency-based complex word identification and replaces the complex words with a lexical simplifier similar to LSBert (Qiang et al., 2020). It is available in 22 languages, including our investigated languages.

**GPT40 mini (multiling.)** is one of the state-ofthe-art LLMs by OpenAI and offers support for all three languages. We prompted it in a zero-shot manner to simplify the text samples. The simplification prompts are presented in Appendix C.

**MUSS (EN)** stands for "Multilingual Unsupervised Sentence Simplification" and is one of the most popular pre-trained sentence simplification models (Martin et al., 2022). We used the pretrained muss\_en\_mined checkpoint that utilizes the BART architecture (Lewis et al., 2020). Even though MUSS is multilingual, it does not support all the languages we investigate. Due to the long runtime of MUSS, we create simplifications only on the fixed subsets of the data.

**Cochrane and Medeasi (EN)** are based on the HuggingFace space simplification-model-app. Both utilize a BART model fine-tuned for simplification in the medical domain. The Medeasi checkpoint uses the sentence-level MED-EASi dataset

Language	Dataset	Dataset name	Prediction Task	#samples (sub- set size)	#classes
EN	AG News Sentiment	AG News (Zhang et al., 2015) Sentiment Analysis for Financial News (Malo et al., 2014)	topic sentiment	7600 (760) 4846 (970)	4 3
	TL;DR	tldr_news	topic, summarization	794	5
DE	Gnad10	10k German News Articles Datasets	topic	1028	9
	ML SUM	(Schabus et al., 2017) Multilingual summarization (DE) (Scialom et al., 2020)	topic, summarization	579	12
RU	ML SUM	Multilingual summarization (RU) (Scialom et al., 2020)	topic, summarization	203	9

Table 1: Overview of all datasets and their classification tasks evaluated in this study.

(Basu et al., 2023), while Cochrane is fine-tuned on the paragraph-level data (Devaraj et al., 2021).
SimplifyText (EN) uses the Keep it Simple (KiS) approach by Laban et al. (2021) and is a GPT2-

**DEplain (DE)** is a German simplification model based on mT5 (Stodden, 2024) and fine-tuned on the DEplain-APA corpus (Stodden et al., 2023).

**Russian simplification (RU)** is a Russian sentence simplification model. It is based on ruT5 and was fine-tuned on the RuSimpleSentEval (Sakhovskiy et al., 2021) and the RuAdapt (Dmitrieva and Tiedemann, 2021) datasets.

## 3.3 Classifiers and LLMs

based simplification model.

251

252

255

258

261

264

265

267

269

271

272

276

277

278

281

Our models under test span from DeBERTa-based classification systems to the latest open- and closedsource large language models. Table 2 gives an overview of the models and settings that we investigated.

For each English classification dataset, we finetuned two DeBERTaV3-base classifiers (He et al., 2023). The first classifier was trained on the original data, while the other classifier was fine-tuned on the data simplified with the SimplifyText model. We selected this model for simplification because it received the best scores among the open-source models in our unsupervised simplification evaluation (see subsection 3.4). Every training was conducted for one epoch with a learning rate of  $2 \cdot 10^{-5}$ . We trained the models on the datasets' training splits, so the test splits used for our investigation were still unseen for the models. With this training setup, we can test how much the models adapt to the specific style of simplification and if text simplification as pre-processing or data augmentation during training is beneficial for performance.

Setting	Language(s)
FT Orig	EN
FT Simple	EN
Zero-shot	EN, DE
Zero-shot	EN, DE
Zero-shot	EN, DE, RU
Zero-shot	EN, DE, RU
	Setting FT Orig FT Simple Zero-shot Zero-shot Zero-shot Zero-shot

Table 2: Overview of all models under test. Traditional models are fine-tuned on either the original training data or a simplified version of it. The LLMs are prompted in a zero-shot manner.

287

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

The second part of our study investigated the performance of large language models. For this, we selected four LLMs, three open-source models from Meta's Llama3.1 family (Grattafiori et al., 2024) and Aya Expanse 8B from Cohere for AI (Dang et al., 2024), and the closed-source GPT4o-mini from OpenAI. Llama3.1 is a multilingual LLM with a context of 128k tokens. For our experiments, we use the instruction-tuned versions with 8B and 70B parameters to account for performance differences due to model size. Llama3.1 70B is loaded with bitsandbytes' 8-bit quantization. Unfortunately, Llama is not available in Russian. In contrast, Aya Expanse 8B exhibits powerful multilingual capacities and supports 23 languages, including the three in our study. For GPT, we were limited to fixed subsets to reduce the financial efforts.

For the predictions themselves, we used the same zero-shot prompt for all four models. The prompts per dataset are presented in Appendix D. A native German or Russian speaker created each of the non-English prompts. Even if we told the models to only predict the topic and not provide any reasoning, some of the outputs still contained more

363

311content than the topic. We tried to account for the312most common phrases among them during post-313processing. Therefore, we lower-cased all model314outputs and removed phrases like "The topic of315this snippet is". In addition, some labels were a316combination of multiple terms, e.g., sci/tech in317AG News. If only one part, e.g., only sci, was pre-318dicted, we considered this prediction correct and319replaced it with the proper topic name.

#### 3.4 Unsupervised simplification evaluation

Previous work has investigated the impact of 321 human-supervised simplifications (Anschütz et al., 322 2024), but for our datasets, human supervision is not feasible. In contrast, we investigate the im-324 pact of automatic text simplification, and thus, we need to evaluate the quality of the automatic sim-326 plifications. Our datasets are not targeted to simplification, and hence, no reference simplification 328 exists. Therefore, we based our evaluation on unsupervised metrics that evaluate the simplification against the source instead of comparing it against a reference. While a human evaluation would be the best solution, this is infeasible for our large-scale 333 study setup with multiple languages, datasets, and simplifiers. To still provide an insightful evaluation of the simplifications, we not only evaluate the overall simplification quality but also the readability of the texts and the meaning preservation indepen-338 339 dently. To measure the readability of the texts and the simplicity-gain through simplification, we used the Flesh-Reading-Ease (FRE) (Flesch, 1948). It is 341 a statistical measure based on the number of words per sentence and the average word length. It can 343 be adapted for many languages, including German and Russian. The score ranges from 0 to 100, with a higher score indicating a higher readability. We used the Python textstat package and the German adaption by Amstad (1978).

> The second aspect of our evaluation is the overall simplification quality. For this, we use two different scores, which are LENS\_SALSA (Heineman et al., 2023) and REFeREE (Huang and Kochmar, 2024). Both metrics are learned metrics that were fine-tuned to mimic human annotation scores. LENS\_SALSA is working on the wordand sentence-level and predicts and scores edit annotations that are performed during simplification. In contrast to this, REFeREE employs a multi-step fine-tuning process that aligns the metric scores with traditional metrics like BLEU (Papineni et al., 2002) and performs a multi-aspect evaluation of the

351

353

357

361

fluency and simplicity of the generated text. While LENS\_SALSA ranges from 0 to 100, REFeREE only ranges from -1 to 1. Therefore, we rescale the REFeREE values to make them comparable with the other metrics.

Finally, the third evaluation criterion is testing if the simplification preserves the original text's meaning. This is especially important for content classification tasks, as in our study. Again, we select two metrics to evaluate the factuality of the simplifications. First, we use FactCC (Kryscinski et al., 2020), which has shown the best human correlation on factuality evaluations like the FRANK dataset (Pagnoni et al., 2021). It was originally designed for the evaluation of abstractive summarization, but since some of our simplification systems perform complex operations close to summarization, we consider this metric suitable. FactCC employs a binary classification to predict whether the summary is factually consistent with its source. For our evaluation, we calculate the percentage of samples that are deemed correct to end up with a value between 0 and 100 again. The last metric is MeaningBERT (Beauchemin et al., 2023), which is specifically targeted toward meaning preservation in text simplification.

We provide a detailed evaluation and correlation analysis only for English, as FRE is the only unsupervised metric that we could find for German and Russian simplification.

# 4 Results and Discussion

### 4.1 Simplification evaluation

We evaluate the simplifications in English based on three criteria: the readability of the texts, the overall simplification quality, and the faithfulness of the 396 simplifications. For this, we automatically score 397 the simplifications with five different metrics (see 398 subsection 3.4 for details). Table 3 shows the met-399 rics scores for the English simplifications. In terms 400 of readability, the Cochrane simplifier achieves the 401 highest scores, indicating the biggest simplicity 402 gain. Interestingly, the FRE scores of GPT4o-mini 403 are rather low compared to the other simplifiers, in-404 dicating that it performs rather conservative simpli-405 fication. Nevertheless, it achieves the best overall 406 simplification quality across all datasets. This is 407 probably due to its great fluency and overall capac-408 ities. In terms of faithfulness, MILES has the best 409 scores across almost all datasets. This is expected 410 since it is a lexical simplification system that does 411

Metric	Original	MILES	Cochrane	Medeasi	SimplifyText	MUSS	GPT40 mini		
AG News									
FRE	48.78	54.13	70.22	58.92	65.93	53.64 †	59.11 <sup>†</sup>		
REFeREE	-	36.08	72.48	67.19	71.0	65.35 <sup>†</sup>	87.84 <sup>†</sup>		
LENS_SALSA	-	53.0	66.56	62.41	64.66	60.74 <sup>†</sup>	70.65 <sup>†</sup>		
FactCC	-	91.63	52.37	85.04	60.39	84.87 †	85.53 †		
Meaning_BERT	-	91.56	67.41	85.62	83.29	90.06 †	82.72 †		
Sentiment									
FRE	55.43	61.76	73.34	65.73	65.52	58.97 <sup>†</sup>	61.76 <sup>†</sup>		
REFeREE	-	51.6	56.74	55.49	67.59	65.61 †	75.46 <sup>†</sup>		
LENS_SALSA	-	60.34	65.88	56.42	69.85	64.29 †	69.34 <sup>†</sup>		
FactCC	-	96.22	54.5	91.48	73.85	95.26 †	96.29 <sup>†</sup>		
Meaning_BERT	-	84.84	50.19	85.12	76.74	83.27 †	78.68 $^{\dagger}$		
TL;DR									
FRE	57.27	63.85	76.2	67.74	62.08	60.73	62.32		
REFeREE	-	39.88	75.25	76.0	79.93	79.48	84.64		
LENS_SALSA	-	60.54	72.05	72.9	73.95	72.84	75.74		
FactCC	-	90.93	49.75	87.03	66.37	86.23	88.92		
Meaning_BERT	-	89.11	67.89	70.18	84.22	88.76	87.77		

Table 3: Unsupervised simplification evaluation of the English simplifiers. For all metrics, higher scores indicate better simplification quality. The best scores per metric are bolded. <sup>†</sup>evaluated only on subset

412 not rewrite the sentences but only replaces some
413 complex words within. Overall, all simplification
414 systems show a good performance and can be used
415 for further experiments.

### 4.2 Model performances

416

To investigate if the model performances change 417 when we simplify the input texts, we compare the 418 accuracies of all classification tasks and the rougeL 419 scores (Lin, 2004) for the summarization tasks as 420 implemented in Huggingface evaluate. For each 421 dataset, we report the results of the two fine-tuned 422 DeBERTa classifiers and the four LLMs in a zero-423 424 shot setting. In addition, we tested whether the changes in accuracy were statistically significant. 425 For this, we performed a related t-test with the hy-426 pothesis that the average of the two distributions 427 was the same. If the p-value is smaller than 0.05, 428 we reject this hypothesis and can conclude that the 429 accuracy change is significant. The results for the 430 English classification task are presented in Table 4. 431 Overall, the fine-tuned classifiers (DeBERTa Orig 432 and DeBERTa Simple) show the best accuracies, 433 with GPT-4o-mini coming the closest. However, 434 nearly all models show a decreased classification 435 performance if the inputs are simplified. No per-436 437 formance improvement is statistically significant. However, the majority of the simplifications intro-438 duce a severe performance drop of up to 20 per-439 centage points. The sentiment dataset is the dataset 440 with the most significant performance changes, 441

even though it has the fewest and most distinct classes. The performance decreases are especially remarkable for the DeBERTa classifier, which was fine-tuned on simplified data. This model exhibits a drop in performance even when the same simplifier is used for training and testing. A similar problem can be observed with GPT4o-mini, which exhibits a performance drop even when it is working on its own simplification outputs. However, statistically significant performance changes of the GPT4o-mini simplifications are scarce. 442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

#### 4.3 Human evaluation

Our results show that all classifiers, even powerful LLMs like GPT-40-mini, exhibit a performance decrease when working with simplified inputs. An obvious explanation for this behavior would be that the simplification systems alter the meaning of the input samples. To examine the meaning preservation of the simplifications, we conducted a human evaluation. We randomly selected 12 samples from each of the three datasets and showed the original and simplified versions to a simple language expert (one of the authors). The samples were presented one by one and we randomized the order of the simplifiers so that the annotator did not know which models created the simplification. Overall, we analyzed 216 original-simplified pairs (12 samples across 3 datasets and 6 simplifiers). The annotator graded the samples on three different aspects: content preservation, the existence of a hallucination,

Model	Original	Original (subset)	MILES	Cochrane	Medeasi	Simplify Text	MUSS	GPT4o mini
AG News - Classification (accuracy)								
DeBERTa Orig	94.5	94.34 <sup>†</sup>	-6.58*	-1.07*	-2.79*	-3.71*	-1.58 †	-3.16*†
DeBERTa Simple	90.26	90.26 *	-3.0* †	-0.61*	-0.83*	-1.7*	-1.05 †	-1.32 <sup>†</sup>
AyaExpanse8B	83.13	80.53	-0.14	-2.91*	-1.43*	-1.56*	$+0.39^{\dagger}$	-0.66 †
Llama3.1 8B	80.12	$78.68$ $^{\dagger}$	-1.3*	-1.96*	-1.48*	-1.58* †	+0.27 *	-5.26*†
Llama3.1 70B	79.97	80.26 †	-0.55*	-0.21	+0.08	-0.36	-0.79 †	$+1.45^{\dagger}$
GPT4o-mini	-	84.08 †	-0.66 †	+1.18 <sup>†</sup>	-0.79 †	$\pm$ 0.0 $^{\dagger}$	$\pm$ 0.0 $^{\dagger}$	$-0.53^{\dagger}$
Sentiment - Classification (accuracy)								
DeBERTa Orig	88.16	86.08 †	-6.0*	-13.91*	-1.98*	-5.65*	-0.82 *	$+0.41^{\dagger}$
DeBERTa Simple	87.49	87.53 †	-6.46*	-12.57*	-1.73*	-3.8*	-1.13 †	$-1.24^{\dagger}$
AyaExpanse8B	68.4	$68.76$ $^{\dagger}$	-5.52*	-17.33*	+0.02	-6.47*	$-2.37^{\dagger}$	$-4.12^{\dagger}$
Llama3.1 8B	68.17	68.56 <sup>†</sup>	-8.95*	-20.57*	-1.1	-14.39*	-7.01*†	-6.5*†
Llama3.1 70B	78.23	78.76 $^{\dagger}$	-3.96*	-10.1*	-1.98*	-5.97*	-4.74*†	-1.96†
GPT4o-mini	80.84	80.72 †	-4.09*	-14.76*	-1.01*	-9.8*	-3.19 <sup>†</sup>	$-0.72^{\dagger}$
		TL;	DR - Class	ification (accu	ıracy)			
DeBERTa Orig	76.32	-	-4.91*	-1.39	-15.37*	-0.25	-2.27*	-1.01
DeBERTa Simple	74.56	-	-3.53*	-0.13	-9.07*	+0.25	-0.38	+0.13
AyaExpanse8B	58.19	-	+0.63	-0.76	-0.13	+0.75	+1.51	+0.63
Llama3.1 8B	44.84	-	-3.4*	-1.26	-3.15	+0.75	$\pm 0.0$	-3.91*
Llama3.1 70B	56.55	-	-5.79*	-4.91*	-6.68*	-2.27	-1.01	-1.13
GPT4o-mini	65.74	-	$\pm 0.0$	$\pm 0.0$	-2.39	-2.01	-0.75	-0.75
<b>TL;DR</b> - Summarization ( <i>rougeL</i> )								
AyaExpanse8B	23.09	-	-2.04*	-5.95*	-4.59*	-2.17*	-0.88*	-0.79*
Llama3.1 8B	23.89	-	-3.17*	-6.4*	-6.08*	-2.34*	-1.37*	-0.98*
Llama3.1 70B	27.04	-	-2.81*	-7.43*	-7.04*	-2.9*	-1.62*	-0.76
GPT4o-mini	25.86	-	-2.67*	-7.72*	-6.3*	-1.99*	-2.01*	-0.02

Table 4: Changes in performance across all English datasets. For most of the models and simplifiers, the scores decrease (red boxes). Only a few combinations show improved performance (blue boxes). \* statistically significant change (p < 0.05), significant changes have a darker color, <sup>†</sup>evaluated and compared only on the fixed subset

and whether the simplified sample preserved the original label. The content and label preservation were ranked on a 4-point Likert scale, while the hallucinations received a binary label.

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489 490

491

492

493

494

The most relevant finding is that only nine out of 216 samples changed the original label, i.e., 96%of the analyzed samples preserved the labels and, thus, should receive the same prediction by the classifiers. In contrast, the results from the content and hallucination evaluation paint a less clear picture, as can be seen in Figure 3. While Medeasi, MUSS, and GPT4o-mini preserve most of the content with almost no hallucinations, the Cochrane and SimplifyText simplifiers show some content alternations. MILES is a lexical simplification system that performs minimal changes and shows decent content preservation. Nevertheless, it is among the simplifiers with the strongest performance drops for the classifiers. This indicates that the choice of words in simplified language is more relevant to the classifiers than the sheer amount of edit operations. This aligns with previous research by Anschütz et al. (2024), who find that the Levenshtein distance between original and simplified samples only has a weak correlation with label changes in LLMs.

495

496

497

498

499

500

501

502

503

504

506

507

508

509

510

511

512

513

514

515

516

Overall, human evaluation could verify our assumption from subsection 3.1: While the simplifiers might change small aspects, these changes do not affect the selected classification tasks, and the overall labels are preserved (some examples are presented in Appendix A). Therefore, we reject faithfulness alone as a trivial explanation for the LLM's bad generalization performance.

#### 4.4 Non-English data

Table 5 and Table 6 show the results for German and Russian respectively. First of all, we can see that the FRE scores increase for all ATS systems, indicating that the simplifiers successfully improved the readability of the samples. Again, the GPT4omini simplifications achieve a comparatively small readability improvement. For Russian, we observe hardly any statistically significant changes, except for some strong improvements of Aya Expanse on the classification task. In general, both Russian models show an extremely weak summarization



(a) Content preservation of the simplified versions across the three English datasets

(b) Number of hallucinations per simplifier

Figure 3: Results from human evaluation. GPT40-mini, Medeasi, and MUSS show the best content preservation and the least hallucinations.

Model	Orig.	DEplain	MILES	GPT40 mini				
<b>Gnad10</b> - Classification ( <i>accuracy</i> )								
FRE	46.41	61.34	59.96	52.55				
AyaExpanse8B	26.75	+7.1*	+2.34*	+4.28*				
Llama3.1 8B	50.78	-5.64*	-3.7*	+0.19				
Llama3.1 70B	33.85	+7.4*	-1.85	+7.88*				
GPT4o-mini	58.95	-4.77*	+3.21*	+1.17				
ML SUM DE - Classification (accuracy)								
FRE	48.84	61.06	62.32	53.25				
AyaExpanse8B	49.74	+3.46	-1.73	+3.11				
Llama3.1 8B	62.0	-1.9	-0.51	+2.42				
Llama3.1 70B	61.14	$\pm 0.0$	-6.74*	+5.18*				
GPT4o-mini	77.72	-7.77*	-2.07*	-1.55				
ML SUM DE - Summarization (rougeL)								
AyaExpanse8B	17.46	-10.97*	-3.05*	-1.7*				
Llama3.1 8B	14.78	-9.19*	-1.99*	-0.71				
Llama3.1 70B	15.63	-9.08*	-1.43*	+0.65				
GPT4o-mini	16.1	-9.98*	-1.4*	+0.24				

Table 5: Accuracy changes on German data, \* statistically significant change (p < 0.05)

performance in terms of rougeL score, even for the original data. Therefore, the changes on simplified data have only minor importance as the models don't seem to fulfill the task at all. For German, we observe many improvements, especially for the Gnad10 classification task. In addition, simplifications by GPT40 show the most significant improvements and only one significant performance drop. This is even the case in the summarization task. Our results allow for two interpretations: Most models are primarily trained on English, and they seem to overfit more to the standard language style in their pre-training<sup>2</sup>. Therefore, their performance

517

518

519

521

523

524

525

527

528

529

Model	Orig.	Russian simpl.	MILES	GPT4o mini		
ML SU	MRU -	Classification	on (accurac	y)		
FRE	48.33	51.66	70.74	49.01		
AyaExpanse8B	32.02	+4.93	+8.37*	+14.29*		
GPT4o-mini	67.98	+1.97	-1.97	-0.49		
ML SUM RU - Summarization (rougeL)						
AyaExpanse8B	2.79	+0.16	-0.82	-0.82		
GPT4o-mini	0.99	-0.49	$\pm 0.0$	$\pm 0.0$		

Table 6: Accuracy changes on Russian data, \* statistically significant change (p < 0.05)

on English simplified language drops significantly. However, for languages with weaker LLM support, we expect less overfitting. Thus, these models can benefit from simplifications, especially if they are of high, human-like quality, as with GPT40-mini. 530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

## 5 Conclusion

Experiments across six datasets, eight ATS systems, and three languages show that English LLMs exhibit a severe performance drop when switching from original to simplified language, uncovering a weak generalization to this language style. However, simplified texts can enhance performance at inference time for non-English languages. We thus encourage content creators to prioritize using simple language online as a way to improve LLMs' downstream performance and comprehension and to open their models to a broader audience.

 $<sup>^244.22\%</sup>$  of Llama's instruction-tuning data belongs to the categories code, exam-like, or reasoning and tools (Grattafiori

et al., 2024, Tab. 7). This data uses highly technical terms or long and technical argumentation chains that would not be used in simplified language.

566

570

573

575

577

581

582

584

591

# Limitations

We provide an extensive evaluation of the em-548 ployed simplification models, evaluating them for their simplicity gain, simplification quality, and meaning preservation with automatic metrics. In 551 addition, we conducted a human evaluation to ver-552 ify our label preservation assumption. However, 553 due to the large scope of our experiments with 554 multiple datasets and simplifiers, we could only evaluate 12 samples per dataset and simplifier combination. The results of this evaluation paint a clear picture, with more than 95% of the samples preserving the original label. Nevertheless, this evalua-559 560 tion could be extended to more samples, evaluation aspects, and the non-English languages.

In addition to this, our investigation only covers a limited set of NLP tasks. We selected the sentiment and classification tasks to avoid biases due to automatic evaluation metrics and insufficient meaning preservation of the simplification models. As shown in our human evaluation, this task selection was valuable as the simplifications sometimes altered the content but preserved the original label. In addition, we tested the performance on summarization as a generation task. Nevertheless, it would be interesting to add further NLP tasks to draw a broader picture of LLM generalization on simplified language. Moreover, since the results indicate that simplifications can improve the performance of non-English languages, this research should be extended to further languages.

578 Finally, we used the same prompts for all models and tested them in a zero-shot setting. This 579 could mean that the models could not unfold their 580 full potential and that the performances could be improved further. However, we don't evaluate the models on an absolute scale; rather, we compare the 583 performance of simplified and original texts. All experiments are conducted under the same setting, 585 and thus, the limitations of the zero-shot setting should not affect our overall results. Another problem could be the data contamination. Since our datasets are quite old, it is likely that they were included in the LLM pre-training data. However, our paper measures the generalization of the LLMs on simplified language. Thus, this change in behavior 592 on unseen data is actually part of our investigation, and the potential data contamination does not affect the validity of our findings. 595

# **Ethical considerations**

Our work uncovers novel insights into how LLMs generalize to simplified language. We don't create any new datasets or models, and thus, there is no harm coming from our investigations. However, we find some alarming behavior in most of the LLMs as our results show that they decrease their performance when using simplified language in English. This can have severe implications for people with low literacy or mental disabilities when using platforms like ChatGPT: When a user asks the chatbot for a summarization of a news snippet in plain language, the models are more likely to make mistakes in these interactions. These people are already a vulnerable target group that struggles to verify information on the internet due to information barriers of overly complicated texts. When easy-to-use and trust-evoking platforms like chatbots show a worse performance when interacting with those people, this implies severe discrimination of users of simplified language that we uncovered with this work.

### References

- Sweta Agrawal and Marine Carpuat. 2024. Do text simplification systems preserve meaning? a human evaluation via reading comprehension. Transactions of the Association for Computational Linguistics, 12:432– 448.
- Toni Amstad. 1978. Wie verständlich sind unsere Zeitungen? Ph.D. thesis, Universität Zürich.
- Miriam Anschütz, Edoardo Mosca, and Georg Groh. 2024. Simpler becomes harder: Do LLMs exhibit a coherent behavior on simplified corpora? In Proceedings of the Workshop on DeTermIt! Evaluating Text Difficulty in a Multilingual Context @ LREC-COLING 2024, pages 185–195, Torino, Italia. ELRA and ICCL.
- Chandrayee Basu, Rosni Vasu, Michihiro Yasunaga, and Qian Yang. 2023. Med-easi: Finely annotated dataset and models for controllable simplification of medical texts. Preprint, arXiv:2302.09155.
- David Beauchemin, Horacio Saggion, and Richard Khoury. 2023. Meaningbert: assessing meaning preservation between sentences. Frontiers in Artificial Intelligence, 6.
- John Dang, Shivalika Singh, Daniel D'souza, Arash Ahmadian, Alejandro Salamanca, Madeline Smith, Aidan Peppin, Sungjin Hong, Manoj Govindassamy, Terrence Zhao, Sandra Kublik, Meor Amer, Viraat Aryabumi, Jon Ander Campos, Yi-Chern Tan, Tom Kocmi, Florian Strub, Nathan Grinsztajn, Yannis

596 597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621 622 623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

757

758

759

760

Flet-Berliac, Acyr Locatelli, Hangyu Lin, Dwarak Talupuru, Bharat Venkitesh, David Cairuz, Bowen Yang, Tim Chung, Wei-Yin Ko, Sylvie Shang Shi, Amir Shukayev, Sammie Bae, Aleksandra Piktus, Roman Castagné, Felipe Cruz-Salinas, Eddie Kim, Lucas Crawhall-Stein, Adrien Morisot, Sudip Roy, Phil Blunsom, Ivan Zhang, Aidan Gomez, Nick Frosst, Marzieh Fadaee, Beyza Ermis, Ahmet Üstün, and Sara Hooker. 2024. Aya expanse: Combining research breakthroughs for a new multilingual frontier. *Preprint*, arXiv:2412.04261.

647

666

671

672

675

678

679

703

- Ashwin Devaraj, Iain Marshall, Byron Wallace, and Junyi Jessy Li. 2021. Paragraph-level simplification of medical texts. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4972–4984, Online. Association for Computational Linguistics.
  - Anna Dmitrieva and Jörg Tiedemann. 2021. Creating an aligned Russian text simplification dataset from language learner data. In *Proceedings of the 8th Workshop on Balto-Slavic Natural Language Processing*, pages 73–79, Kiyv, Ukraine. Association for Computational Linguistics.
  - Rudolph Flesch. 1948. A new readability yardstick. *Journal of applied psychology*, 32(3):221.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, and Ava Spataru et al. 2024. The Ilama 3 herd of models. *Preprint*, arXiv:2407.21783.
- Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2023. Debertav3: Improving deberta using electra-style pretraining with gradient-disentangled embedding sharing. *Preprint*, arXiv:2111.09543.
- David Heineman, Yao Dou, Mounica Maddela, and Wei Xu. 2023. Dancing between success and failure: Edit-level simplification evaluation using SALSA. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 3466–3495, Singapore. Association for Computational Linguistics.
- Yichen Huang and Ekaterina Kochmar. 2024. REFeREE: A REference-FREE model-based metric for text simplification. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 13740–13753, Torino, Italia. ELRA and ICCL.
- Ludivine Javourey-Drevet, Stéphane Dufau, Thomas François, Núria Gala, Jacques Ginestié, and Johannes C. Ziegler. 2022. Simplification of literary and scientific texts to improve reading fluency and

comprehension in beginning readers of french. *Applied Psycholinguistics*, 43(2):485–512.

- Tannon Kew, Alison Chi, Laura Vásquez-Rodríguez, Sweta Agrawal, Dennis Aumiller, Fernando Alva-Manchego, and Matthew Shardlow. 2023. BLESS: Benchmarking large language models on sentence simplification. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 13291–13309, Singapore. Association for Computational Linguistics.
- Wojciech Kryscinski, Bryan McCann, Caiming Xiong, and Richard Socher. 2020. Evaluating the factual consistency of abstractive text summarization. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9332–9346, Online. Association for Computational Linguistics.
- Philippe Laban, Tobias Schnabel, Paul Bennett, and Marti A. Hearst. 2021. Keep it simple: Unsupervised simplification of multi-paragraph text. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 6365–6378, Online. Association for Computational Linguistics.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Pekka Malo, Ankur Sinha, Pekka Korhonen, Jyrki Wallenius, and Pyry Takala. 2014. Research article. *Journal of the Association for Information Science and Technology*, 65(4):782 – 796.
- Louis Martin, Angela Fan, Éric de la Clergerie, Antoine Bordes, and Benoît Sagot. 2022. MUSS: Multilingual unsupervised sentence simplification by mining paraphrases. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 1651–1664, Marseille, France. European Language Resources Association.
- Sneha Mehta, Bahareh Azarnoush, Boris Chen, Avneesh Saluja, Vinith Misra, Ballav Bihani, and Ritwik Kumar. 2020. Simplify-then-translate: Automatic preprocessing for black-box translation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):8488–8495.
- Rei Miyata and Midori Tatsumi. 2019. Evaluating the suitability of human-oriented text simplification for

761

- 795
- 796 797
- 799

801

- 805
- 807
- 808

810 811 812

813

814 815

816

(LREC'14), pages 3720–3723, Reykjavik, Iceland. European Language Resources Association (ELRA). 817

1241–1244, Tokyo, Japan.

- machine translation. In Proceedings of the 33rd Pacific Asia Conference on Language, Information and Computation, pages 147–155. Waseda University.
- Dennis Murphy Odo. 2022. The Effect of Automatic Text Simplification on L2 Readers' Text Comprehension. Applied Linguistics, 44(6):1030–1046.
- Artidoro Pagnoni, Vidhisha Balachandran, and Yulia Tsvetkov. 2021. Understanding factuality in abstractive summarization with FRANK: A benchmark for factuality metrics. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 4812–4829, Online. Association for Computational Linguistics.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Jipeng Qiang, Yun Li, Yi Zhu, Yunhao Yuan, and Xindong Wu. 2020. Lsbert: A simple framework for lexical simplification. Preprint, arXiv:2006.14939.
- Michael Ryan, Tarek Naous, and Wei Xu. 2023. Revisiting non-English text simplification: A unified multilingual benchmark. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 4898– 4927, Toronto, Canada. Association for Computational Linguistics.
- Andrey Sakhovskiy, Alexandra Izhevskaya, Alena Pestova, Elena Tutubalina, Valentin Malykh, Ivan Smurov, and Ekaterina Artemova. 2021. Rusimplesenteval-2021 shared task: evaluating sentence simplification for russian. In Proceedings of the International Conference "Dialogue, pages 607–617.
- Andreas Säuberli, Franz Holzknecht, Patrick Haller, Silvana Deilen, Laura Schiffl, Silvia Hansen-Schirra, and Sarah Ebling. 2024. Digital comprehensibility assessment of simplified texts among persons with intellectual disabilities. In Proceedings of the CHI Conference on Human Factors in Computing Systems, CHI '24, New York, NY, USA. Association for Computing Machinery.

Dietmar Schabus, Marcin Skowron, and Martin Trapp.

2017. One million posts: A data set of german on-

line discussions. In Proceedings of the 40th Inter-

national ACM SIGIR Conference on Research and

Development in Information Retrieval (SIGIR), pages

Jordan Schmidek and Denilson Barbosa. 2014. Im-

proving open relation extraction via sentence re-

structuring. In Proceedings of the Ninth International Conference on Language Resources and Evaluation Thomas Scialom, Paul-Alexis Dray, Sylvain Lamprier, Benjamin Piwowarski, and Jacopo Staiano. 2020. MLSUM: The multilingual summarization corpus. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 8051-8067, Online. Association for Computational Linguistics.

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

- Zhenmei Shi, Junyi Wei, Zhuoyan Xu, and Yingyu Liang. 2023. Why larger language models do in-context learning differently? In RO-FoMo: Robustness of Few-shot and Zero-shot Learning in Large Foundation Models.
- Sanja Štajner and Maja Popovic. 2016. Can text simplification help machine translation? In Proceedings of the 19th Annual Conference of the European Association for Machine Translation, pages 230–242.
- Regina Stodden. 2024. Reproduction of German text simplification systems. In Proceedings of the Workshop on DeTermIt! Evaluating Text Difficulty in a Multilingual Context @ LREC-COLING 2024, pages 1–15, Torino, Italia. ELRA and ICCL.
- Regina Stodden, Omar Momen, and Laura Kallmeyer. 2023. DEplain: A German parallel corpus with intralingual translations into plain language for sentence and document simplification. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 16441–16463, Toronto, Canada. Association for Computational Linguistics.
- Jan Trienes, Sebastian Joseph, Jörg Schlötterer, Christin Seifert, Kyle Lo, Wei Xu, Byron C. Wallace, and Junyi Jessy Li. 2024. InfoLossQA: Characterizing and recovering information loss in text simplification. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics, pages 4263-4294.
- Hoang Van, Zheng Tang, and Mihai Surdeanu. 2021. How may I help you? using neural text simplification to improve downstream NLP tasks. In Findings of the Association for Computational Linguistics: EMNLP 2021, pages 4074–4080, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- David Vickrey and Daphne Koller. 2008. Sentence simplification for semantic role labeling. In Proceedings of ACL-08: HLT, pages 344-352, Columbus, Ohio. Association for Computational Linguistics.
- Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. In Advances in Neural Information Processing Systems, volume 28. Curran Associates, Inc.

- 868
- 869 870

875

877

879

884

896

900

901

902

904

905

906

907

909

910

911

913

914

915

917

# A Examples form human evaluation

See Table 7 for examples where the content is altered by the simplifier but the overall label is still preserved.

# **B** Simplifications to augment context

Previous work by Van et al. (2021) experimented with simplification as data augmentation. One of their experiments was to concatenate the original and simplified texts together. Our benchmarking of LLMs is done in a zero-shot setting, so data augmentation at training time is out-of-scope for our work. However, we tested how the models perform when they see the concatenated versions during inference time. For this, we created two additional versions of the data samples. First we concatenated the simplified text to the original one with just a whitespace in between. This version is called *orig+simp*. To ablate whether the accuracy changes are based on the input length or the additional context, we also created a version where the original samples were concatenated to themselves (identified as *orig+orig*). Due to the shorter context window of the DeBERTa models and their fine-tuning to a specific input style and length, we only ran these ablations on the LLMs. In addition, we only tested these settings for the classification tasks. Figure 4 shows the accuracy curves of the original, simple, and the two combined versions. For the larger models like GPT4o-mini and LLama 70B (bottom two), the augmentations seem to make no difference. Moreover, for Aya Expanse, the concatenations on the TL;DR dataset seem to worsen the performance even further. In contrast to this, for the smaller Llama 8B model, the *orig+simple* versions can improve the performance by over 10 percentage points, matching the performance of the larger Llama model. In contrast, we only see minor improvements or even decreased performances of the *orig+orig* concatenations. This implies that the simplifications give additional context or explanations to the original content that can improve the zero-shot performance of some of the smaller language models. This aligns with previous findings that different LLM sizes perform in-contextlearning differently and that smaller models orient themselves more on the task description, while larger models rely on the knowledge they obtained during pre-training (Shi et al., 2023). For our experiments, this means that adding the simplifications to the original samples has a higher impact on the

model performance of the smaller models.

# C LLM simplification prompts

We used GPT4o-mini to create high-quality simplifications. We used the following prompt where sample is replaced by the text to be predicted. For German and Russian, the prompt is translated, respectively. 918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

**Simplify (EN):** {"role": {"system", "content": "You are a helpful assistant. You will be provided with sentences from news articles. Your task is to simplify the texts to enhance readability. You must not alter the meaning and don't provide reasoning." },

{"role": "user", "content": "{sample} - Simplification: "}

**Simplify DE:** {"role": {"system", "content": "Du bist ein hilfreicher Assistent. Du bekommst Sätze aus Nachrichtenartikeln. Deine Aufgabe ist es, die Texte zu vereinfachen, um die Verständlichkeit zu erhöhen. Du darfst den Inhalt nicht verändern und brauchst keine Begründungen angeben." },

{"role": "user", "content": "{sample} - Vereinfachung: "}

Simplify RU: {"role": {"system", "content": "Ты - полезный помощник. Тебе будут предоставлены предложения из новостных статей. Твоя задача - упростить текст, чтобы повысить его читабельность. Ты не должен изменять смысл и приводить аргументы." }, {"role": "user", "content": "{sample} - Упрощение: "}

# **D** LLM Prediction prompts

We used the same system prompts for all four large language models and prompted them in a zero-shot manner. The prompts differ per dataset and language. Below are the prompts we used for the classification and summarization tasks where sample is replaced by the text to be predicted.

AG News (EN): {"role": {"system", "content": "You are a helpful assistant. You will be provided with sentences from news articles. Classify each query into a news topic. There are four possible topics: world, sports, business or sci/tech. You must not choose another topic. Answer only with one single word and do not provide reasoning." },

Original	Simplified	Label
Sudan Peace Talks Resume for South as Tensions Brew KHARTOUM/NAIROBI (Reuters) - Sudan's government resumed talks with rebels in the oil-producing south on Thursday while the United Nations set up a panel to investigate charges of geno- cide in the west of Africa's largest coun- try	Sudan peace talks resume in south as tensions rise KHARTOUM/NAIROBI (Reuters) - Sudan's government held peace talks on Thursday with south-west rebels, while the United Nations set up a panel to investigate allegations of geno- cide in the <b>world's largest country</b> .	world
Operating income rose to EUR 696.4 mn from EUR 600.3 mn in 2009.	This year's <b>net profit more than doubled</b> to EUR 696.4 mn from EUR 600.3 mn in 2009.	positive
All art establishments are concerned with the degradation of paintings. Harmful fac- tors such as sunlight, moisture, and cer- tain volatile organic compounds can accel- erate degradation. Graphene may be the solution to protecting art from exposure to harmful agents. A one-atom-thick sheet of graphene can adhere easily to various substrates and serve as an excellent barrier against oxygen, gases, moisture, and UV light. The graphene sheets can be added to framing glass for artworks with extremely rough surfaces or embossed patterns. The sheets can be removed using a soft rubber eraser.	All art establishments are concerned with the degradation of paintings. Harmful fac- tors such as sunlight, moisture, and certain volatile organic compounds can accelerate the process of deterioration. Graphene, which is made of a variety of materials, can be applied to framing glass to protect against oxygen, gases, and UV light. It can also be used as a barrier against bac- teria and fungi, which can cause skin irritation.	Science & Futuristic Technology

Table 7: Examples from the human evaluation. All simplifications are factually incorrect or introduce hallucinations (bolded parts). Even with these content errors, the original labels are preserved.



Figure 4: Combining original and simplified texts as inputs. Accuracy difference by Aya Expanse 8B, Llama8B, Llama70B, and GPT40 mini. We additionally plot the accuracy of the original data as a dashed grey line to enhance comparability.

{"role": "user", "content": "{sample} - The topic
is"}

964

965

966

967

968 969

970

Sentiment (EN): {"role": {"system", "content": "You are a helpful assistant. You will be provided with sentences from articles. Classify the sentiment of each query. There are three possible sentiments: positive, neutral or negative. You must not choose another sentiment. Answer only with one single word and do not provide reasoning."},

973 {"role": "user", "content": "{sample} - The senti-974 ment is"}

TL;DR (EN): {"role": {"system", "content":
"You are a helpful assistant. You will be provided
with sentences from news articles. Classify each
query into a news topic. There are five possible
topics: 'Sponsor', 'Big Tech & Startups', 'Science
& Futuristic Technology', 'Programming & Design
& Data Science' and 'Miscellaneous'. You must
not choose another topic. Answer only with one
single word and do not provide reasoning." },

984 {"role": "user", "content": "{sample} - The topic 985 is"} **Gnad10 (DE):** {"role": {"system", "content": "Du bist ein hilfreicher Assistent. Du bekommst Sätze aus Nachrichtenartikeln. Ordne jede Anfrage einem Nachrichtenthema zu. Es gibt neun mögliche Themen: Web, Panorama, International, Wirtschaft, Sport, Inland, Etat, Wissenschaft und Kultur. Du darfst kein anderes Thema wählen. Antworte nur mit einem einzigen Wort und gib keine Begründung an." },

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1002

1003

1004

1005

1007

"role": "user", "content": "{sample} - Das Thema ist"}

ML SUM (DE): {"role": {"system", "content": "Du bist ein hilfreicher Assistent. Du bekommst Sätze aus Nachrichtenartikeln. Ordne jede Anfrage einem Nachrichtenthema zu. Es gibt zwölf mögliche Themen: politik, wirtschaft, geld, panorama, sport, muenchen, digital, karriere, bildung, reise, auto und stil. Du darfst kein anderes Thema wählen. Antworte nur mit einem einzigen Wort und gib keine Begründung an." },

{"role": "user", "content": "{sample} - Das Thema
ist"}

14

ML SUM (RU): {"role": {"system", 1008 "content": "Ты - полезный ассистент. Те-1009 бе будут предоставлены предложения из но-1010 востных статей. Классифицируй каждый за-1011 прос в соответствии с темой новости. Темы 1012 даны на английском языке, и есть девять воз-1013 можных тем: science, politics, mosobl, culture, 1014 social, incident, economics, sport, moscow. Ты 1015 не должен выбирать какую-либо другую те-1016 му. Отвечай только одним словом и не объ-1017 ясняй." }, 1018

1019 {"role": "user", "content": "{sample} - Тема"}

1020Summarize (EN): {"role": {"system", "con-1021tent": "You are a helpful assistant. You will be1022provided with sentences from news articles. Your1023task is to create a headline that summarizes the1024content. Answer only with one sentence and don't1025provide reasoning." },

1026 {"role": "user", "content": "{sample} - The head-1027 line is"}

1028Summarize DE:{"role": {"system", "content":1029"Du bist ein hilfreicher Assistent. Du bekommst1030Sätze aus Nachrichtenartikeln. Deine Aufgabe ist1031es, einen Titel zu verfassen, der den Inhalt zusam-1032menfasst. Antworte nur mit einem Satz und gib1033keine Begründung an." },

1034{"role": "user", "content": "{sample} - Der Titel1035ist"}

1036Summarize RU: {"role": {"system", "content":1037"Ты - полезный помощник. Тебе будут предо-1038ставлены предложения из новостных статей.1039Твоя задача - придумать заголовок, кото-1040рый обобщает содержание статьи. Отвечай1041только одним предложением и не приводи1042аргументы." },

1043 {"role": "user", "content": "{sample} - Заголо-1044 вок:"}