DREsS: Dataset for Rubric-based Essay Scoring on EFL Writing

Anonymous ACL submission

Abstract

Automated essay scoring (AES) is a useful tool in English as a Foreign Language (EFL) writing education, offering real-time essay scores for students and instructors. However, previous AES models were trained on essays and scores irrelevant to the practical scenarios of EFL writing education and usually provided a single holistic score due to the lack of appropriate datasets. In this paper, we release DREsS, a large-scale, standard dataset for rubric-based automated essay scoring with 48.9K samples 011 in total. DREsS comprises three sub-datasets: 012 DREsS_{New}, DREsS_{Std.}, and DREsS_{CASE}. We 014 collect DREsS_{New}, a real-classroom dataset with 2.3K essays authored by EFL undergraduate students and scored by English education 017 experts. We also standardize existing rubricbased essay scoring datasets as DREsS_{Std.}. We suggest CASE, a corruption-based augmenta-019 tion strategy for essays, which generates 40.1K synthetic samples of DREsS_{CASE} and improves the baseline results by 45.44%. DREsS will enable further research to provide a more accurate and practical AES system for EFL writing 025 education.¹

1 Introduction

027

In writing education, automated essay scoring (AES) can provide real-time scores of students' essays to both students and instructors. For many students who are hesitant to expose their errors to instructors, the immediate assessment of their essays with AES can create a supportive environment for self-improvement in writing skills (Sun and Fan, 2022). For instructors, AES models can ease the time-consuming process of evaluation and serve as a means to validate their assessments, ensuring consistency in their evaluations.

AES systems can provide either a holistic or an analytic view of essays, but rubric-based, analytical



1. DREsS_New (2,279 samples) EFL classroom data: 1) Student-written essays 2) Rubric-based scores assessed by instructors

2. DREsS_std. (6,515 samples) Unified AES datasets with standardized rubrics under professional consultation

041

042

043

044

045

049

054

057

059

060

061

062

063

064

065

067

068

3. DREsS_case (40,185 samples) Synthetic essay samples generated by CASE, our proposed augmentation strategy

Figure 1: Data construction of DREsS

scores are more preferred in the EFL writing education domain (Ghalib and Al-Hattami, 2015). However, there is only a limited amount of rubric-based datasets available for AES, and the rubrics are not consistent in building generalizable AES systems. Furthermore, AES datasets must be annotated by writing education experts because the scoring task requires pedagogical knowledge of English writing. To date, there is a lack of usable datasets for training rubric-based AES models, as existing AES datasets provide only overall scores and/or make use of scores annotated by non-experts.

In this paper, we release DREsS (Dataset for Rubric-based Essay Scoring on EFL Writing), a large-scale dataset for rubric-based essay scoring using three key rubrics: content, organization, and language. DREsS consists of three datasets: 1) DREsS_{New} with 2,279 essays from English as a foreign language (EFL) learners and their scores assessed by experts, 2) DREsS_{Std.} with 6,515 essays and scores from existing datasets, and 3) DREsS_{CASE} with 40,185 synthetic essay samples. We standardize and rescale existing rubric-based datasets to align our rubrics. We also suggest CASE, a corruption-based augmentation strategy for Essays, employing three rubricspecific strategies to augment the dataset with corruption. DREsS_{CASE} improves the baseline result by 45.44%.

¹We will provide a non-anonymous link to the dataset in the camera-ready version of this manuscript.

		Content	Organization	Language
DREsS _{New}		2,279	2,279	2,279
	ASAP P7	1,569	1,569	1,569
	ASAP P8	723	723	723
DREsS _{Std.}	ASAP++ P1	1,785	1,785	1,785
	ASAP++ P2	1,799	1,799	1,799
	ICNALE EE	639	639	639
DREsS _{CASE}	3	8,307	31,086	792
Total		17,101	39,880	9,586

Table 1: Data statistics of DREsS

2 Related Work

In this section, we describe previous studies in automated essay scoring (AES) in terms of the format of predicted scores: holistic AES (§2.1) and rubric-based AES (§2.2). To date, there is only a limited amount of publicly available AES datasets, and their rubrics are inconsistent. Furthermore, their scores are usually annotated by non-experts lacking pedagogical knowledge in English writing. Here, we introduce DREsS, a publicly available, large-scale, rubric-based, real-classroom dataset, which can be used as training data for rubric-based AES systems.

2.1 Holistic AES

077

081

086

100

101

ASAP Prompt 1-6 ASAP dataset² is widely used in AES tasks, involving eight different prompts. Six out of eight prompt sets (Prompt 1-6) have a single overall score. This holistic AES includes 10K essay scoring data on sourcedependent essay (Prompt 3-6) and argumentative essay (Prompt 1-2). However, these essays are graded by non-expert annotators, though the essays were written by Grade 7-10 students in the US.

TOEFL11 TOEFL11 (Blanchard et al., 2013) corpus from ETS introduced 12K TOEFL iBT essays, which are not publicly accessible now. TOEFL11 only provides a general score for essays in 3 levels (low/mid/high), which is insufficient for building a well-performing AES system.

Models The majority of the previous studies used the ASAP dataset for training and evaluation, aiming to predict the overall score of the essay only (Tay et al., 2018; Cozma et al., 2018; Wang et al., 2018; Yang et al., 2020, *inter alia*). Enhanced AI Scoring Engine (EASE)³ is a commonly used, open-sourced AES system based on feature extraction and statistical methods. In addition, Taghipour and Ng (2016) and Xie et al. (2022) released models based on recurrent neural networks and neural pairwise contrastive regression (NPCR) model, respectively. Still, only a limited number of studies publicly opened their models and codes, highlighting the need for additional publicly available data and further validation of existing models. 102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

2.2 Rubric-based AES

ASAP Prompt 7-8 ASAP includes only two prompts (Prompt 7-8) that are rubric-based. These two rubric-based prompts consist of 1,569 and 723 essays for each respective prompt. The two prompt sets even have distinct rubrics and score ranges, which poses a challenge in leveraging both datasets for training rubric-based models. These essays (Prompt 7-8) are also evaluated by non-expert annotators, similar to ASAP Prompt 1-6.

ASAP++ To overcome the holistic scoring of ASAP Prompt 1-6, Mathias and Bhattacharyya (2018) manually annotated rubric-based scores on those essays. However, most samples in ASAP++ were annotated by a single annotator, who is a nonexpert, including non-native speakers of English. Moreover, each prompt set of ASAP++ has different attributes or rubrics to each other, which need to be more generalizable to fully leverage such dataset for AES model.

ICNALE Edited Essays ICNALE Edited Essays (EE) v3.0 (Ishikawa, 2018) presents rubric-based essay evaluation scores and fully edited versions of

²https://www.kaggle.com/c/asap-aes

³https://github.com/edx/ease

essays written by EFL learners from 10 countries 136 in Asia. Even though the essays are written by 137 EFL learners, the essay is rated and edited only 138 by single annotator per sample. They have five 139 native English speakers, non-experts in the domain 140 of English writing education in total. In addition, 141 it is not openly accessible and only consists of 639 142 samples. 143

The scarcity of publicly available rubric-Models 144 based AES datasets poses significant obstacles to the advancement of AES research. 146 There 147 are industry-driven services such as IntelliMetric® (Rudner et al., 2006) and E-rater® (Blanchard 148 et al., 2013; Attali and Burstein, 2006), but none 149 of them are accessible to the public. Kumar et al. 150 (2022) proposed applying a multi-task learning ap-151 proach in holistic AES with ASAP and ASAP++, 152 using traits as auxiliary tasks. Recent studies have 153 followed up their method, introducing multi-traits AES approaches (Chen and Li, 2023; Do et al., 155 2023, 2024; Lee et al., 2024, inter alia). Still, they 156 shed light on predicting a holistic score only due to limited data and built eight different fine-tuned 158 models due to unconsolidated rubrics by each es-159 say prompt. In order to facilitate AES research in 160 the academic community, it is crucial to release a 161 publicly available rubric-based AES dataset and baseline model. 163

3 DREsS Dataset

164

165

167

168

169

170

171

172

173

174

175

176

177

178

180

181

184

We construct DREsS with 2.3K samples of our newly collected dataset (§3.1), 6.5K standardized samples of existing datasets (§3.2), and 40.1K synthetic samples augmented using CASE (§3.3). The detailed number of samples per rubric is stated in Table 1.

3.1 Dataset Collection

Dataset Details DREsS_{New} includes 2,279 argumentative essays on 22 prompts, having 313.36 words and 21.19 sentences on average. Each sample in DREsS includes students' written essay, essay prompt, rubric-based scores, total score (the sum of three rubric-based scores), and a test type (pre-test, post-test). The essays are scored on a range of 1 to 5, with increments of 0.5, based on the three rubrics: *content*, *organization*, and *language*. We chose such three conventional rubrics as standard criteria for scoring EFL essays, according to previous studies from the language education (Cumming, 1990; Ozfidan and Mitchell, 2022).

Rubric	Description
Content	Paragraph is well-developed and relevant to the argument, sup- ported with strong reasons and ex- amples.
Organization	The argument is very effectively structured and developed, making it easy for the reader to follow the ideas and understand how the writer is building the argument. Paragraphs use coherence devices effectively while focusing on a sin- gle main idea.
Language	The writing displays sophisticated control of a wide range of vocab- ulary and collocations. The essay follows grammar and usage rules throughout the paper. Spelling and punctuation are correct throughout the paper.

Table 2: Rubric explanations

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

209

210

Brief explanations of the rubrics are shown in Table 2. The essays are written by undergraduate students whose TOEFL writing score spans from 15 to 21 and enrolled in EFL writing courses at a college in South Korea from 2020 to 2023. Most students are Korean and their ages span from 18 to 22, with an average of 19.7. During the course, students are asked to write an in-class timed essay for 40 minutes both at the start (pre-test) and the end of the semester (post-test) to measure their improvements.

Annotator Details We collect scoring data from 11 instructors, who serve as the teachers of the students who wrote the essays. Six of them are non-native speakers, and five of them are native speakers. All annotators are experts in English education or Linguistics and are qualified to teach EFL writing courses at a college in South Korea. One instructor was allocated per essay, so the interannotator agreement cannot be measured. It follows that an EFL course is usually led by a single instructor, and the essays from the course are assessed by the instructor in a real-classroom setting. To ensure consistent and reliable scoring across all instructors, they all participate in training sessions with a scoring guide and norming sessions where



Figure 2: Score distribution of DREsS

they develop a consensus on scores using two sample essays. Additionally, there was no significant difference among the score distribution of all instructors tested by one-way ANOVA and Tukey HSD at a *p*-value of 0.05.

3.2 Standardizing the Existing Data

217

218

221

222

223

227

233

240

241

242

243

244

245

247

248

We standardize and unify three existing rubricbased datasets (ASAP Prompt 7-8, ASAP++ Prompt 1-2, and ICNALE EE) to align with the three rubrics in DREsS: content, organization, and *language*. We exclude ASAP++ Prompt 3-6, whose essay type, source-dependent essays, is clearly different from argumentative essays. We create synthetic label based on a weighted average and then rescale the score of all rubrics into a range of 1 to 5. Detailed explanations and rationales behind standardizing weights are described in Appendix C. In the process of consolidating the writing assessment criteria, we sought professional consultation from EFL education experts and strategically grouped together those components that evaluate similar aspects under theoretical considerations.

3.3 Synthetic Data Construction

We construct synthetic data for rubric-based AES to overcome the scarcity of data and provide accurate scores for students and instructors. We introduce a corruption-based augmentation strategy for essays (CASE), which starts with a *well-written* essay and incorporates a certain portion of sentence-level errors into the synthetic essay. In subsequent experiments, we define *well-written* essays as an essay that scored 4.5 or 5.0 out of 5.0 on each criterion.

$$\mathbf{n}(S_c) = \lfloor \mathbf{n}(S_E) * (5.0 - x_i)/5.0 \rceil$$
(1)

 $n(S_c)$ is the number of corrupted sentences in the synthetic essay, and $n(S_E)$ is the number of sentences in the *well-written* essay, which serves as the basis for the synthetic essay. x_i denotes the score of the synthetic essay. In this paper, we generate synthetic data with CASE under ablation study for exploring the optimal number of samples.

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

267

269

270

271

272

273

274

275

276

277

278

279

285

Content We substitute randomly-sampled sentences from *well-written* essays with out-of-domain sentences from different prompts. This is based on an assumption that sentences in *well-written* essays support the given prompt's content, meaning that sentences from the essays on different prompts convey different contents. Therefore, more number of substitutions imply higher levels of corruption in the content of the essay.

Organization We swap two randomly-sampled sentences in *well-written* essays and repeat this process based on the synthetic score, supposing that sentences in *well-written* essays are systematically structured in order. The higher number of swaps implies higher levels of corruption in the organization of the essay.

Language We substitute randomly-sampled sentences into ungrammatical sentences and repeat this process based on the synthetic score. We extract 605 ungrammatical sentences from BEA-2019 data for the shared task of grammatical error correction (GEC) (Bryant et al., 2019). We define ungrammatical sentences with the number of edits of the sentence over 10, which is the 98th percentile. The more substitutions, the more corruption is introduced in the grammar of the essay. We set such a high threshold for ungrammatical sentences because of the limitation of the current GEC dataset that inherent noise may be included, such as erroneous or incomplete correction (Rothe et al., 2021).

3.4 Score Distribution

Figure 2 shows the score distribution of $DREsS_{New}$ and $DREsS_{Std.}$ ranging from 0 to 5. The score distribution of the AES dataset shows a left-skewed bell-shaped curve, following the general trends in real-classroom settings. The scarcity of samples on

Model	Strategy	Content	Organization	Language	Total
EASE (SVR)		-	-	-	0.360
NPCR (Xie et al., 2022)		-	-	-	0.507
ArTS (Do et al., 2024)	SFT w/ DREsS	0.601	0.743	0.592	<u>0.690</u>
BERT (Devlin et al., 2019)		0.642	0.750	0.607	0.685
Llama 3.1 8B (AI@Meta, 2024)		<u>0.631</u>	0.771	0.589	0.691
	(A) zero-shot ICL	0.310	0.322	0.231	0.304
ant-40	(B) five-shot ICL	0.361	0.475	0.367	0.428
866.10	(C) rubric explanation	0.285	0.250	0.200	0.259
	(D) feedback generation	0.313	0.268	0.230	0.290

Table 3: Baseline results of rubric-based automated essay scoring on DREsS (QWK score)

low scores is because instructors are reluctant to give low scores to increase students' self-efficacy and motivate them to learn (Arsyad Arrafii, 2020). To overcome the imbalance of the dataset, we propose CASE, which can generate synthetic data for all score ranges. DREsS_{CASE} has the same number of samples per score.

4 Experimental Result

4.1 Baseline Result on DREsS

Table 3 shows the baseline results of rubric-based AES on DREsS. We use all three subsplits of DREsS as training data, but DREsS_{New}, a subsplit comprising essays and scores from real classroom settings, is used exclusively for the validation and the test sets. In other words, synthetically unified (DREsS_{Std}) or augmented (DREsS_{CASE}) data are reserved for training to avoid incomplete or inaccurate evaluation. Detailed experimental settings are described in Appendix §A. We adopt the quadratic weighted kappa (QWK) scores, a conventional metric to evaluate the consistency between the predicted scores and the gold standard scores.

We provide the baseline results on DREsS using holistic AES models from previous studies (*i.e.*, EASE (SVR), NPCR (Xie et al., 2022), and ArTS (Do et al., 2024)), large language model (*i.e.*, gpt-4o from OpenAI⁴ and Llama 3.1 8B (AI@Meta, 2024) from Meta), and BERT (Devlin et al., 2019). Note that fine-tuned BERT is the model that most state-of-the-art AES systems have leveraged. We train EASE (SVR), NPCR, ArTS, BERT, and Llama 3.1 with DREsS as supervised

fine-tuning (SFT) data. We also test gpt-40 with four different system prompts as follows:	319 320
(A) in-context learning (ICL) with zero-shot	321
(B) in-context learning (ICL) with five-shots of writing prompts and essays	322 323
(C) asking the model to predict essay scores given detailed rubric explanations	324 325

326

327

328

329

331

332

333

334

335

336

337

339

340

341

342

343

344

345

346

347

348

350

351

(D) asking the model to predict essay scores and provide essay feedbacks that support their predicted scores.

The detailed prompts are described in Appendix B.1. Considering the substantial length of writing prompt and essay, we were able to provide a maximum of 5 shots for the prompt to gpt-40. We divided the samples into five distinct score ranges and computed the average total score for each group. Subsequently, we randomly sampled a single essay in each group, ensuring that its total score corresponded to the calculated mean value. Asking gpt-40 to score an essay shows high variances among the essays with the same score, implying their limitations to be applied as AES systems.

4.2 Validation of DREsS_{Std.} and DREsS_{CASE}

Table 4 shows experimental results of rubric-based AES with different language models. We train Longformer (Beltagy et al., 2020) and BigBird (Zaheer et al., 2020), a language model that accepts long input sequences (*i.e.*, 4,096 tokens), considering the substantial length of writing prompts and essays. In addition, we train GPT-NeoX-20B (Black et al., 2022) and Llama 3.1 8B, state-of-the-art LLMs. Nonetheless, exploiting different models does not significantly affect the performance of

311

314

315

316

317

318

287

⁴All following experiments using gpt-40 in this paper was conducted from May 21, 2024 to June 5, 2024 under OpenAI API services.

Model	Strategy	Content	Organization	Language	Total
BERT (Devlin et al., 2019)		0.414	0.311	0.487	0.471
Longformer (Beltagy et al., 2020)		0.409	0.312	<u>0.475</u>	0.463
BigBird (Zaheer et al., 2020)	SFT w/ DREsS _{New}	0.412	0.317	0.473	0.469
GPT-NeoX-20B (Black et al., 2022)		0.410	0.313	0.446	0.475
Llama 3.1 8B (AI@Meta, 2024)		<u>0.413</u>	0.375	0.426	0.466

Table 4: Experimental results of rubric-based AES with different LMs using DREsS_{New}

Model	Strategy	Content	Organization	Language	Total
	SFT w/ DREsS _{New}	0.414	0.311	0.487	0.471
BERT (Devlin et al., 2019)	+ DREsS _{Std.}	0.599	0.593	0.587	0.551
	+ DREsS _{CASE}	0.642	<u>0.750</u>	0.607	<u>0.685</u>
	SFT w/ DREsS _{New}	0.413	0.375	0.426	0.466
Llama 3.1 8B (AI@Meta, 2024)	+ DREsS _{Std.}	0.581	0.608	0.574	0.563
	+ DREsS _{CASE}	<u>0.631</u>	0.771	<u>0.589</u>	0.691

Table 5: Empirical validation of data expansion in DREsS

AES systems. Xie et al. (2022) also observed that leveraging different foundation models has no significant effect on AES performance, and most state-of-the-art AES methods have still leveraged BERT (Devlin et al., 2019). Therefore, based on these observations, we choose BERT and Llama 3.1 (8B) as a representative model to further evaluate and validate the effectiveness of our dataset, particularly focusing on the benefits of data standardization and synthesis.

We validate the practical benefits of data standardization (DREsS_{Std.}) and synthesis (DREsS_{CASE}) with empirical results. Both finetuned BERT and Llama 3.1 exhibit scalable results with the expansion of training data (Table 5). In particular, the model trained with a combination of our approaches outperforms other baseline models by 45.44%, demonstrating the effectiveness of data unification and augmentation using CASE. Interestingly, a state-of-the-art LLM (*i.e.*, gpt-40) does not outperform fine-tuned small-scale language models (*i.e.*, BERT), achieving 0.257 points lower QWK total score. Existing holistic AES models show their inability to compute rubric-based scores.

5 Discussion & Analysis

5.1 Ablation Study

353

354

359

361

371

374

380

We perform an ablation study to find the optimal number of CASE operations per each rubric. In Fig-



Figure 3: Ablation experimental results for CASE. n_{aug} is the number of synthetic data by each class per original data among all classes. The x-axis is a log-arithmetic scale.

381

383

384

385

386

387

389

390

391

392

393

394

395

397

ure 3, we investigate how the number of CASE operations affects the performance over all rubrics for $n_{aug} = \{0.125, 0.25, 0.5, 1, 2, 4, 8\}, \text{ where } n_{aug}$ denotes the number of synthetic data by each class per original data among all classes (i.e., the ratio of augmented data size compared to the source data size). CASE on *content*, *organization*, and language rubrics show their best performances on 0.5, 2, 0.125 of n_{aug} , generating a pair of synthetic essays and corresponding scores in 4.5, 18, 1.125 times, respectively. We suppose that the detailed augmentation strategies for each rubric and the small size of the original data affect the optimal number of CASE operations. Organization, where corruption was made within the essay and irrelevant to the size of the original data, showed the highest n_{auq} . Content, where the corrupted sentences were sampled from 874 *well-written* essays with 21.2 sentences on average, reported higher n_{aug} than language, where the corrupted sentences were sampled from 605 ungrammatical sentences. Leveraging more error patterns in new grammatical error correction (GEC) data will lead to a scalable increase in the size of DREsS_{CASE} for *language*.

5.2 CASE vs. Generative Methods

398

400

401

402

403

404

405

406

407

408

409

410

411

412

413

	Content	Organization	Language
gpt-4o	0.298	0.219	0.158
CASE (Ours)	0.625	0.722	0.635

 Table 6: QWK scores of synthetic essays generated by

 two augmentation methods

We verify the quality of synthetic data using CASE compared to generative methods using LLMs. Here, we use the best-performing baseline rubric-based scoring models trained with DREsS. We measure a quadratic weighted kappa (QWK) score to measure the similarity between the gold label of the synthetic sample and the predicted score by an AES model.

For LLM to generate synthetic essays, we first 414 give the persona of an EFL student taking an En-415 glish writing course in a college for students who 416 get TOEFL scores ranging from 15 to 21 and 417 provide five example essays written by EFL stu-418 dents randomly sampled from five distinct score 419 ranges. We then ask the model to write an essay 420 that matches the rubric-based scores. The detailed 421 prompts to generate synthetic EFL essays are de-422 423 scribed in Appendix B.2. We randomly sample 900 essays (100 samples per score ranging from 1.0 to 424 5.0 with an increment of 0.5) from CASE augmen-425 tation and synthetic samples generated by gpt-40. 426 Table 6 shows QWK scores of synthetic essays, 427 which validate whether the essays match with their 428 scores. We use the best-performing baseline rubric-429 based scoring models in Table 4, which only uses 430 DREsS_{New} as its training and test set. QWK score 431 of CASE augmentation achieves 0.661 (substan-432 *tial agreement*), while the score of the generative 433 method achieves 0.225 on average (slight to fair 434 agreement). Though the detailed persona and ex-435 436 ample essays are given, gpt-40 fails to write an appropriate level of essays. Specifically, the pre-437 dicted rubric-based scores of 900 synthetic essays 438 from gpt-40 across all score ranges are $4.21_{\pm 0.65}$, 439 $4.13_{\pm 0.63}$, and $4.30_{\pm 0.30}$ for *content*, *organization*, 440

and *language*, respectively.

We discuss the benefit of leveraging CASE to generate synthetic essays in EFL writing for three reasons: 1) its difficulty in generating EFL students' essays, 2) low performance in scoring essays, and 3) controllability and interoperability. First of all, LLMs are hardly capable of replicating EFL learners' errors since they are mostly trained with texts from native speakers. The essays of DREsS_{New} written by EFL students reveal various unique characteristics and error patterns of EFL learners. Detailed analysis is described in \S 5.3. Second, we found that the state-of-theart LLM, namely gpt-40, underperforms in essay scoring tasks compared to BERT-based models, as described in Table 3. Lastly, the black-box nature of LLMs poses challenges in terms of controllability and interpretability. In contrast, our proposed CASE method offers enhanced control and interpretability. This mitigates the risks associated with over-reliance on generative methods, fostering a more robust and transparent research approach.

441

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

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

5.3 In-depth Analysis

Table 7 shows quantitative analysis of essays from DREsS_{New} and DREsS_{CASE} compared to gpt-40 augmentation concerning linguistic features. Student-written essays in DREsS_{New} include unique patterns of ELF learners. For instance, essays in DREsS_{New} tend to be longer than synthetic essays from gpt-40, with more number of sentences but easier and shorter sentences, according to Flesch reading ease (Flesch, 1948) and the number of tokens, respectively. Interestingly, EFL students use fewer unique words but frequently use unnecessary stopwords. Essays from EFL students include typos and spelling errors which cannot be made from the generation outputs of LLMs. Note that one of the major strengths of the DREsS dataset is the inclusion of errorful essays written by EFL learners in the real-world classroom.

Table 8 shows two sample essays with a score of 1 under the same writing prompt. The synthetic essay from gpt-40 fails to reflect the EFL learners' errors, generating essays that include *content*, *organization*, and *language* features needed for a wellwritten essay. For *organization*, the essay from gpt-40 is well-structured with the use of appropriate transition signals and an appropriate thesis sentence in the first paragraph (blue text). For *content*, each body paragraph includes detailed examples to support the argument (orange text). For language,

	DREsS _{New}	DREsS _{CASE}	gpt-4o
# of sentences *	$20.96_{\pm 6.66}$	$22.67_{\pm 10.10}$	$16.02_{\pm 2.35}$
# of tokens *	$313.97_{\pm 96.76}$	$327.91 {\scriptstyle \pm 56.01}$	$285.84_{\pm 69.07}$
# of tokens w/o stopwords	$162.64_{\pm 49.97}$	$167.14_{\pm 35.50}$	$165.49_{\pm 47.91}$
Type-token ratio (TTR) *	$0.43_{\pm 0.07}$	$0.43_{\pm 0.06}$	$0.51_{\pm 0.04}$
# of transition signal *	$28.03_{\pm 10.3}$	$28.18_{\pm 14.0}$	$29.61_{\pm 10.02}$
# of typos *	$4.39_{\pm 4.11}$	$6.64_{\pm 8.22}$	$0.59_{\pm 1.22}$
Flesch reading ease (Flesch, 1948) *	$53.66_{\pm 11.84}$	$59.47_{\pm 11.11}$	$29.15_{\pm 19.88}$
US grade level (Kincaid et al., 1975) *	$9.58_{\pm 2.11}$	$9.12_{\pm 2.58}$	$13.64_{\pm 3.5}$

Table 7: Quantitative analysis. The asterisk denotes a statistically significant difference between gpt-40 and others tested under one-way ANOVA and Tukey HSD test.

DREsS _{New}	gpt-4o
In Korea, most of high school students stay with very regular schedule and eat healthy lunch. How- ever, in college they start drinking and often eat fast food. Their health going bad! College students have much free time and they can do their hobby more. Then, spending time to exercise. Exercise can help college student. Many college people want to make girl or boyfriend. If they start exercise and make healthy body, it can be positive effect to people.	If I have the chance to visit a foreign country for two weeks, I would like to go to Japan. Japan is a country that I have always wanted to visit because of its unique culture and beautiful landscapes. One of the main reasons I want to visit Japan is because of its culture. Japan has a very different culture from my country, and I am very interested in learning about it. I want to see the traditional Japanese ceremonies and try wearing a kimono. I also want to learn about Japanese tea ceremonies because they are very famous. Another reason I want to visit Japan is because of its beautiful landscapes. Japan has many beautiful mountains and forests that I would like to explore. I have seen pictures of Mount Fuji, and I think it is very beautiful. I would like to go hiking there and see the cherry blossoms in spring. In conclusion, I want to visit Japan because it has a unique culture and beautiful landscapes. I think it would be a very interesting and enjoyable trip

Table 8: Sample essays whose rubric scores on content, organization, and language are all 1 out of 5. Sample-level analysis on *content* (orange), *organization* (blue), and *language* (red).

the essay does not include any grammatical errors.
 In contrast, the essay from DREsS_{New} lacks transitional signals, a thesis sentence, and supporting examples. The essay also includes a few grammatical errors and awkward phrases (red text), as it is written by EFL learners in a real-world classroom.

6 Conclusion

492

493

494

495

496

497

498

499

502

We release the DREsS, a large-scale, standard rubric-based essay scoring dataset with three subsets: DREsS_{New}, DREsS_{Std.}, and DREsS_{CASE}.
 DREsS_{New} is the first reliable AES dataset with

2.3K samples whose essays are authored by EFL undergraduate students and whose scores are annotated by instructors with expertise. According to previous studies from language education, we also standardize and unify existing rubric-based AES datasets as DREsS_{Std}. We finally suggest CASE, corruption-based augmentation strategies for essays, which generates 40.1K synthetic samples and improves the baseline result by 45.44%. This work aims to encourage further AES research and practical application in EFL education. 503

514 Limitations

524

525

529

530

531

535

537

541

542

543

545

547

548 549

551

553

554

555

557

560

563

515 Our research focuses on learning *English* as a for-516 eign language because there already exist datasets, 517 and the current language models perform the best 518 for English. There are many L2 learners of other 519 languages whose writing classes can also benefit 520 from AES. Our findings can illuminate the direc-521 tions of data collection, annotation, and augmenta-522 tion for L2 writing education in other languages as 523 well. We leave that as future work.

> DREsS_{New} is collected through the EFL writing courses from a college in South Korea, and most of the essays are written by Korean EFL students. EFL students in different cultural and linguistic backgrounds might exhibit different essaywriting patterns, which might affect the distribution of scores and feedback. We suggest a further extension of collecting the DREsS dataset from diverse countries.

Our augmentation strategy primarily starts from *well-written* essays and generates erroneous essays along with corresponding scores; therefore, this approach faces challenges in synthesizing *well-written* essays. However, we believe that *well-written* essays can be reliably produced by LLMs, which have demonstrated strong capabilities in generating high-quality English text. Also, an optimized rationale (*e.g.*, a threshold in corruption, corruption scale) will advance CASE, which we leave for future work.

We acknowledge that the experimental results in Table 3-4 might not fully cover state-of-theart models in AES. Nonetheless, it is noteworthy that those results are a *baseline* for our dataset. We emphasize that the core contribution of this paper is the construction and the public release of a large-scale AES dataset (DREsS), not a proposal for AES model architecture. We believe nine different models-namely, state-of-the-art AESspecialized models (EASE, NPCR, ArTS), LLMs (GPT-40, Llama 3.1, GPT-NeoX), and transformerbased models with different input sizes (BERT, Longformer, BigBird)-sufficiently cover empirical testing of existing models. We leave examining state-of-the-art AES models for future work, with a proposal of and comparison to a novel architecture.

Ethics Statement

All studies in this research project were conducted with the approval of our institutional review board

(IRB). Annotators were fairly compensated (ap-564 proximately USD 18), which exceeds the minimum 565 wage in the Republic of Korea in 2024 (approxi-566 mately USD 7.3). To prevent any potential impact 567 on student scores or grades, we requested students 568 to share their essays only after the end of the EFL 569 courses. We also acknowledged and addressed the 570 potential risk associated with releasing a dataset 571 containing human-written essays, especially con-572 sidering privacy and personal information. To miti-573 gate these risks, we plan to 1) employ rule-based 574 coding and 2) conduct thorough human inspections 575 to filter out all sensitive information. Addition-576 ally, access to our data will be granted only to 577 researchers or practitioners who submit a consent 578 form, ensuring responsible and ethical usage. 579

580

613

References

AI@Meta. 2024. Llama 3 model card.	581
Mohammad Arsyad Arrafii. 2020. Grades and grade	582
inflation: exploring teachers' grading practices in in-	583
donesian efl secondary school classrooms. <i>Pedagogy</i> ,	584
<i>Culture & Society</i> , 28(3):477–499.	585
Yigal Attali and Jill Burstein. 2006. Automated essay scoring with e-rater® v.2. <i>The Journal of Technology, Learning and Assessment</i> , 4(3).	586 587 588
Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. <i>Preprint</i> , arXiv:2004.05150.	589 590 591
Sidney Black, Stella Biderman, Eric Hallahan, Quentin	592
Anthony, Leo Gao, Laurence Golding, Horace	593
He, Connor Leahy, Kyle McDonell, Jason Phang,	594
Michael Pieler, Usvsn Sai Prashanth, Shivanshu Puro-	595
hit, Laria Reynolds, Jonathan Tow, Ben Wang, and	596
Samuel Weinbach. 2022. GPT-NeoX-20B: An open-	597
source autoregressive language model. In <i>Proceed-</i>	598
ings of BigScience Episode #5 – Workshop on Chal-	599
lenges & Perspectives in Creating Large Language	600
Models, pages 95–136, virtual+Dublin. Association	601
for Computational Linguistics.	602
Daniel Blanchard, Joel Tetreault, Derrick Higgins,	603
Aoife Cahill, and Martin Chodorow. 2013. Toefl11:	604
A corpus of non-native english. <i>ETS Research Report</i>	605
<i>Series</i> , 2013(2):i–15.	606
Christopher Bryant, Mariano Felice, Øistein E. Ander-	607
sen, and Ted Briscoe. 2019. The BEA-2019 shared	608
task on grammatical error correction. In <i>Proceedings</i>	609
of the Fourteenth Workshop on Innovative Use of NLP	610
for Building Educational Applications, pages 52–75,	611
Florence, Italy. Association for Computational Lin-	612

guistics.

- 614 615 616
- 618
- 621

tional Linguistics.

51.

guistics.

236.

ies, 25:117-130.

for Computational Linguistics.

Computational Linguistics.

Mădălina Cozma, Andrei Butnaru, and Radu Tudor

Ionescu. 2018. Automated essay scoring with string

kernels and word embeddings. In Proceedings of the

56th Annual Meeting of the Association for Com-

putational Linguistics (Volume 2: Short Papers), pages 503-509, Melbourne, Australia. Association

Alister Cumming. 1990. Expertise in evaluating second

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of

deep bidirectional transformers for language under-

standing. In Proceedings of the 2019 Conference of

the North American Chapter of the Association for

Computational Linguistics: Human Language Tech-

nologies, Volume 1 (Long and Short Papers), pages

4171-4186, Minneapolis, Minnesota. Association for

Heejin Do, Yunsu Kim, and Gary Lee. 2024. Autore-

gressive score generation for multi-trait essay scoring.

In Findings of the Association for Computational Linguistics: EACL 2024, pages 1659-1666, St. Julian's,

Malta. Association for Computational Linguistics.

Heejin Do, Yunsu Kim, and Gary Geunbae Lee. 2023.

Prompt- and trait relation-aware cross-prompt essay

trait scoring. In Findings of the Association for Com-

putational Linguistics: ACL 2023, pages 1538–1551,

Toronto, Canada. Association for Computational Lin-

Rudolph Flesch. 1948. A new readability yardstick. Journal of applied psychology, 32(3):221.

Timnit Gebru, Jamie Morgenstern, Briana Vec-

Thikra K Ghalib and Abdulghani A Al-Hattami. 2015.

Shinichiro Ishikawa. 2018. The icnale edited essays; a

J Peter Kincaid, Robert P Fishburne Jr, Richard L Rogers, and Brad S Chissom. 1975. Derivation of

new readability formulas (automated readability in-

dex, fog count and flesch reading ease formula) for

dataset for analysis of 12 english learner essays based

on a new integrative viewpoint. English Corpus Stud-

Holistic versus analytic evaluation of efl writing: A

case study. English Language Teaching, 8(7):225-

for datasets. Commun. ACM, 64(12):86-92.

chione, Jennifer Wortman Vaughan, Hanna Wallach,

Hal Daumé III, and Kate Crawford. 2021. Datasheets

language compositions. Language Testing, 7(1):31-

- 622 623

- 630
- 631
- 635 636
- 638
- 639

647

652

- 664

- 667

navy enlisted personnel. Institute for Simulation and Training, University of Central Florida. 670

- Yuan Chen and Xia Li. 2023. PMAES: Prompt-Rahul Kumar, Sandeep Mathias, Sriparna Saha, and Pushpak Bhattacharyya. 2022. Many hands make mapping contrastive learning for cross-prompt automated essay scoring. In Proceedings of the 61st light work: Using essay traits to automatically score Annual Meeting of the Association for Computational essays. In Proceedings of the 2022 Conference of Linguistics (Volume 1: Long Papers), pages 1489the North American Chapter of the Association for 1503, Toronto, Canada. Association for Computa-Computational Linguistics: Human Language Technologies, pages 1485-1495, Seattle, United States. Association for Computational Linguistics.
 - J. Richard Landis and Gary G. Koch. 1977. The measurement of observer agreement for categorical data. *Biometrics*, 33(1):159–174.

671

672

673

674

675

676

677

678

679

680

681

682

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

- Sanwoo Lee, Yida Cai, Desong Meng, Ziyang Wang, and Yunfang Wu. 2024. Unleashing large language models' proficiency in zero-shot essay scoring. In Findings of the Association for Computational Linguistics: EMNLP 2024, pages 181-198, Miami, Florida, USA. Association for Computational Linguistics.
- Sandeep Mathias and Pushpak Bhattacharyya. 2018. ASAP++: Enriching the ASAP automated essay grading dataset with essay attribute scores. In Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018), Miyazaki, Japan. European Language Resources Association (ELRA).
- Burhan Ozfidan and Connie Mitchell. 2022. Assessment of students' argumentative writing: A rubric development. Journal of Ethnic and Cultural Studies, 9(2):pp. 121–133.
- Sascha Rothe, Jonathan Mallinson, Eric Malmi, Sebastian Krause, and Aliaksei Severyn. 2021. A simple recipe for multilingual grammatical error correction. In 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), pages 702–707, Online. Association for Computational Linguistics.
- Lawrence M. Rudner, Veronica Garcia, and Catherine Welch. 2006. An evaluation of intellimetricTM essay scoring system. The Journal of Technology, Learning and Assessment, 4(4).
- Bo Sun and Tingting Fan. 2022. The effects of an aweaided assessment approach on business english writing performance and writing anxiety: A contextual consideration. Studies in Educational Evaluation. 72:101123.
- Kaveh Taghipour and Hwee Tou Ng. 2016. A neural approach to automated essay scoring. In Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, pages 1882-1891, Austin, Texas. Association for Computational Linguistics.
- Yi Tay, Minh Phan, Luu Anh Tuan, and Siu Cheung Hui. 2018. Skipflow: Incorporating neural coherence features for end-to-end automatic text scoring. Proceedings of the AAAI Conference on Artificial Intelligence, 32(1).

 Yucheng Wang, Zhongyu Wei, Yaqian Zhou, and Xuanjing Huang. 2018. Automatic essay scoring incorporating rating schema via reinforcement learning. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 791–797, Brussels, Belgium. Association for Computational Linguistics.

727

728

730

731

734

737

738 739

740

741 742

743 744

745

746

747

748

749 750

751

753 754

755

756

- Sara Cushing Weigle. 2002. *Assessing Writing*, volume 1. Cambridge University Press, Englewood Cliffs, NJ.
- Jiayi Xie, Kaiwei Cai, Li Kong, Junsheng Zhou, and Weiguang Qu. 2022. Automated essay scoring via pairwise contrastive regression. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 2724–2733, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
 - Ruosong Yang, Jiannong Cao, Zhiyuan Wen, Youzheng Wu, and Xiaodong He. 2020. Enhancing automated essay scoring performance via fine-tuning pre-trained language models with combination of regression and ranking. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1560–1569, Online. Association for Computational Linguistics.
- Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, and Amr Ahmed. 2020. Big bird: Transformers for longer sequences. In *Advances in Neural Information Processing Systems*, volume 33, pages 17283–17297. Curran Associates, Inc.

Appendix

759

A Experimental Settings

Hyperparameter	Value
Batch Size	32
Number of epochs	10
Early Stopping Patience	5
Learning Rate	2e-5
Learning Rate Scheduler	Linear
Optimizer	AdamW

Table 9: SFT configuration

We split DREsS_{New} into training, validation, and test sets in a 6:2:2 ratio with a random seed of 22. We use DREsS_{Std.} and DREsS_{CASE}, a unified or 762 763 augmented data as training data only. Additionally, we separate the training, validation, and test set first and then apply CASE in Table 3. In other words, training data does not include augmented essays from high-quality essays in the test set, which prevents data leakage. The AES experiments except for ArTS, GPT-NeoX-20B, and Llama 3.1 (8B) in Table 4 were conducted under GeForce RTX 2080 Ti (4 GPUs), 256GiB system memory, and Intel(R) Xeon(R) Silver 4114 CPU @ 2.20GHz (40 CPU cores) with hyperparameters denoted in 773 774 Table 9. Fine-tuning ArTS, GPT-NeoX-20B, and Llama 3.1 (8B) was conducted under Quadro RTX 775 8000 (4 GPUs), 377GiB system memory, and Intel(R) Xeon(R) Silver 4214R CPU @ 2.40GHz (48 CPU cores) with the same hyperparameters. LLM inference uses greedy decoding (i.e., temperature 779 0.0).

B LLM Prompting

This section provides detailed system prompts used for the experiments in this paper.

B.1 Automated Essay Scoring

Table 10 illustrates four different system prompts used in experiments for Table 4.

B.2 Synthetic Essay Generation

You are an English as a foreign language (EFL) learner taking an English writing course in a college for students who get TOEFL scores ranging from 15 to 21.

Examples 1-5: <five pairs of
writing prompts and EFL student's
essays>
Scoring criteria: <three rubrics
explanation>

Write an essay with short paragraphs about the given prompt, of which scores are <score> out of 5.0 for all criteria. Note that the essay should include erroneous patterns or typos from EFL students, according to the score.

Essay prompt: <essay_prompt>

789

790

791

792

793

794

796

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

C Rationale Behind Standardizing

The weights are not arbitrarily chosen but were determined through expert consultation and theoretical considerations. Specifically, ASAP Prompt 7 contains four rubrics—ideas, organization, style, and convention-, while Prompt 8 contains six rubrics-ideas and content, organization, voice, word choice, sentence fluency, and convention. Both sets provide scores ranging from 0 to 3. For language, we first create synthetic labels based on a weighted average. This involves assigning a weight of 0.66 to the style and 0.33 to the convention in ASAP Prompt 7, and assigning equal weights to voice, word choice, sentence fluency, and convention in ASAP Prompt 8. Stylistic features, such as tone, coherence, and voice, are emphasized as higher-order concerns in writing assessment frameworks, while conventions, such as grammar and punctuation, are considered lower-order concerns. This theoretical understanding, combined with consultation with EFL education experts, informs our decision to assign a higher weight to style, particularly for argumentative essays where persuasive and expressive abilities are crucial (Weigle, 2002). For content and organization, we utilize the existing data rubric (idea for content, organization as same) in the dataset. We repeat the same process with ASAP++ Prompt 1 and 2, which have the same attributes as ASAP Prompt 8. Similarly, for ICNALE EE dataset, we unify vocabulary, language use, and mechanics as language rubric with a weight of 0.4, 0.5, and 0.1, respectively.

788

(A)	Please score the essay with three rubrics: content, organization, and language. ### Answer format: {content: Float, organization: Float, language: Float} Note that the float values of scores are within [1.0, 1.5, 2.0, 2.5, 3.0, 3.5, 4.0, 4.5, 5.0]. Please answer only in the above JSON format. ### prompt: <essay prompt=""></essay>
(B)	<pre>### essay: <student's essay=""> Please score the essay with three rubrics: content, organization, and language. ### Answer format: {content: Float, organization: Float, language: Float} Note that the float values of scores are within [1.0, 1.5, 2.0, 2.5, 3.0, 3.5, 4.0, 4.5, 5.0]. Please answer only in the above JSON format. ### Examples 1–5:</student's></pre>
	### prompt: <essay prompt=""> ### essay: <student's essay=""></student's></essay>
(C)	Please score the essay with three rubrics: content, organization, and language. <three explanation="" rubrics=""> ### Answer format: {content: Float, organization: Float, language: Float} Note that the float values of scores are within [1.0, 1.5, 2.0, 2.5, 3.0, 3.5, 4.0, 4.5, 5.0]. Please answer only in the above JSON format.</three>
	<pre>### prompt: <essay prompt=""> ### essay: <student's essay<="" pre=""></student's></essay></pre>
(D)	Please score the essay with three rubrics: content, organization, and language. ### Answer format: {content: Float, organization: Float, language: Float, content_feedback: String, organization_feedback: String, language_feedback: String} Note that the float values of scores are within [1.0, 1.5, 2.0, 2.5, 3.0, 3.5, 4.0, 4.5, 5.0]. Please answer only in the above JSON format, with feedback.
	<pre>### prompt: <essay prompt=""> ### essay: <student's essay=""></student's></essay></pre>

Table 10: Four different prompts for gpt-40 to get rubric-based scores in the last four rows of Table 4

D Additional Annotations

822

While norming session and significance tests are 823 conducted for the collection of DREsS_{New} in tan-824 dem, we agree that involving additional annota-825 tors will enhance the credibility of the scores. 826 Therefore, we recruit additional expert annota-827 tions and conduct further annotations with 10% of the total dataset (227 samples). We recruit five English education experts holding a Secondary School Teacher's Certificate (Grade II) for the En-831 glish Language, licensed by the Ministry of Ed-832 ucation, South Korea. Re-annotation achieves Cohen's Kappa score (κ) of 0.193, representing 834 slight agreement between original scores and re-835

annotations (Landis and Koch, 1977). It is noteworthy that the subjectivity of the essay-scoring task and the broad score range of essays (9 classes on a range of 1 to 5 with increments of 0.5) make this task challenging to reach agreements. Furthermore, the discrepancy in expertise domain between two annotator groups—college-level education and K-12 education—might lead to different scoring criteria and relatively low inter-annotator agreement. 836

837

838

839

840

841

842

843

844

845

846

E Datasheet for Dataset

In this section, we document DREsS following the format of Datasheets for Datasets (Gebru et al., 848 850

2021). The details on the composition and the col-

lection process of the CSRT dataset are described

1. For what purpose was the dataset created?

We aim to construct a large-scale, standard,

rubric-based dataset for automated essay scor-

ing (AES) to build AES systems that meet the

2. Who created the dataset (e.g., which team,

research group) and on behalf of which en-

tity (e.g., company, institution, organiza-

tion)? The authors constructed DREsS by

1) collecting new essays and scores from the

writing courses in their institution, 2) standard-

izing existing works, and 3) synthesizing new

3. Who funded the creation of the dataset?

1. Was any preprocessing/cleaning/labeling

of the data done (e.g., discretization or

bucketing, tokenization, part-of-speech tag-

ging, SIFT feature extraction, removal of

instances, processing of missing values)?

No. Instead, we conduct rule-based postprocessing and human inspection to filter out

2. Was the "raw" data saved in addition to the

support unanticipated future uses)? N/A

3. Is the software that was used to preprocess/-

1. Has the dataset been used for any tasks

2. Is there a repository that links to any or

3. What (other) tasks could the dataset be

used for? DREsS can be used as a training

and evaluation dataset for automated essay

all papers or systems that use the dataset?

clean/label the data available? N/A

preprocessed/cleaned/labeled data (e.g., to

See the Acknowledgments and Disclosure of

needs of both instructors and students.

in the main text.

E.1 Motivation

samples.

Funding section.

sensitive information.

E.2 Preprocessing/cleaning/labeling

E.4

Distribution

form.

website.

under the MIT license.

with the instances? No.

vidual instances? No.

will maintain DREsS.

of the main text.

found in the future.

14

E.5 Maintenance

1. Will the dataset be distributed to third par-

ties outside of the entity (e.g., company, in-

stitution, organization) on behalf of which

the dataset was created? Yes, the dataset is

open to the public who submitted a consent

2. How will the dataset will be distributed

3. Will the dataset be distributed under a

copyright or other intellectual property

(IP) license, and/or under applicable terms

of use (ToU)? The dataset will be distributed

4. Have any third parties imposed IP-based

5. Do any export controls or other regulatory

1. Who will be supporting/hosting/maintain-

2. How can the owner/curator/manager of the

dataset be contacted (e.g., email address)?

The owner/curator/manager(s) of the dataset

are the authors of this paper. They can be

contacted through the emails on the first page

3. Is there an erratum? We will release an

4. Will the dataset be updated (e.g., to correct

labeling errors, add new instances, delete

instances)? Yes, the dataset will be updated

whenever it can be extended to other red-

teaming benchmarks. These updates will be

posted on the main web page for the dataset.

plicable limits on the retention of the data

associated with the instances (e.g., were the

individuals in question told that their data

5. If the dataset relates to people, are there ap-

erratum at the GitHub repository if errors are

ing the dataset? The authors of this paper

restrictions apply to the dataset or to indi-

or other restrictions on the data associated

(e.g., tarball on website, API, GitHub)?

The dataset will be distributed through our

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

- 851
- 852
- 855
- 857

- 862

- 870
- 872
- 873 874
- 876
- 878
- 881

- 887

892

E.3 Uses

N/A

already? No.

scoring tasks.

938	would be retained for a fixed period of time
939	and then deleted)? N/A
940	6. Will older versions of the dataset continue
941	to be supported/hosted/maintained? Yes.
942	7. If others want to extend/augment/build on/-
943	contribute to the dataset, is there a mecha-
944	nism for them to do so? No.