# **Rethinking Evaluation Metrics for Grammatical Error Correction:** Why Use a Different Evaluation Process than Human?

#### **Anonymous ACL submission**

### Abstract

One of the goals of automatic evaluation met-001 rics in grammatical error correction (GEC) is to rank GEC systems such that it matches 004 human preferences. However, current automatic evaluations are based on procedures that diverge from human evaluation. Specif-007 ically, human evaluation derives rankings by aggregating sentence-level relative evaluation results, e.g., pairwise comparisons, using a rating algorithm, whereas automatic evalua-011 tion averages sentence-level absolute scores to obtain corpus-level scores, which are then sorted to determine rankings. In this study, we propose an aggregation method for exist-015 ing automatic evaluation metrics which aligns with human evaluation methods to bridge this gap. We conducted experiments using vari-017 ous metrics, including edit-based metrics, ngram based metrics, and sentence-level metrics, 019 and show that resolving the gap improves results for the most of metrics on the SEEDA benchmark. We also found that even BERTbased metrics sometimes outperform the metrics of GPT-4. We publish our unified implementation of the metrics and meta-evaluations: https://anonymized\_for\_review.

#### 1 Introduction

027

037

041

Grammatical error correction (GEC) task aims to automatically correct grammatical errors and surface errors such as spelling and orthographic errors in text. Various GEC systems have been proposed based on sequence-to-sequence models (Katsumata and Komachi, 2020; Rothe et al., 2021), sequence tagging (Awasthi et al., 2019; Omelianchuk et al., 2020), and language models (Kaneko and Okazaki, 2023; Loem et al., 2023), and it is crucial to rank those systems based on automatic evaluation metrics to select the best system matching user's demands. Automatic evaluation is expected to rank GEC systems aligning with human preference, as evidenced by meta-evaluations of automatic met-



Figure 1: An overview of current human and automatic evaluation when ranking three GEC systems based on a dataset containing two sentences. Each system output represents edits for simplicity.

rics that assess their agreement with human evaluation (Grundkiewicz et al., 2015; Kobayashi et al., 2024b). For example, one can compute Spearman's rank correlation coefficient between the rankings produced by automatic and human evaluation, considering a metric with a higher correlation as a better metric.

However, despite the clear goal of reproducing human evaluation, current automatic evaluation is based on procedures that diverge from human evaluation. Figure 1 illustrates the evaluation procedure for ranking three GEC systems using a dataset comprising two sentences. In human evaluation, corrected sentences generated for the same input sentence are compared relatively across system outputs, i.e., pairwise comparison, and the re-

sults are aggregated as rankings using rating algo-058 rithms such as TrueSkill (Herbrich et al., 2006). In 059 contrast, automatic evaluation estimates sentence-060 wise scores, then averages them at the corpus level and determines rankings by sorting these averaged scores. As such, current automatic evaluation fol-063 lows a procedure that deviates from human evalua-064 tion, contradicting the goal of reproducing human judgment. Intuitively, it would be desirable for automatic evaluation to follow the same procedure as 067 human evaluation.

In this study, we hypothesize that resolving this gap will more closely align automatic evaluation to human evaluation. Based on this hypothesis, we propose computing rankings in automatic evaluation using the same procedure as human evaluation, e.g., using TrueSkill after deriving pairwise estimates based on sentence-wise scores when human evaluation is employing TrueSkill. In our experiments, we conducted a meta-evaluation on various existing automatic evaluation metrics using the SEEDA dataset (Kobayashi et al., 2024b) that is a representative meta-evaluation benchmark. The results show that bridging the identified gap improves ranking capability for many metrics and that BERTbased (Devlin et al., 2019) automatic evaluation metrics can even outperform large language models (LLMs), GPT-4, in evaluation. Furthermore, we discuss the use and development of automatic evaluation metrics in the future, emphasizing that sentence-level relative evaluation is particularly important for developing new evaluation metrics.

## 2 Gap Between Human and Automatic Evaluation

#### 2.1 Background

071

073

076

079

084

090

091

100

101

102

103

104

105

107

Human evaluation has been conducted by Grundkiewicz et al. (2015), who manually evaluated systems submitted to the CoNLL-2014 shared task, and by Kobayashi et al. (2024b), who included state-of-the-art GEC systems such as LLMs in their dataset. In both studies, system rankings were derived by applying a rating algorithm to sentence-level pairwise comparisons. Commonly used rating algorithms include Expected Wins and TrueSkill (Herbrich et al., 2006; Sakaguchi et al., 2014). Grundkiewicz et al. (2015) adopted Expected Wins as their final ranking method, whereas Kobayashi et al. (2024b) used TrueSkill to determine the final ranking. Kobayashi et al. (2024b) also pointed out the importance of aligning the granularity of evaluation between automatic evaluation and human evaluation but did not mention the procedure for converting sentence-level evaluation into system rankings. 108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

Automatic evaluation is conducted using various evaluation metrics, including reference-based and reference-free approaches, as well as sentencelevel and edit-based metrics. Most of these metrics follow a procedure in which each sentence is assigned an absolute score, which is then aggregated into a corpus-level evaluation score. For example, sentence-level metrics such as SOME (Yoshimura et al., 2020) and IMPARA (Maeda et al., 2022) aggregate scores by averaging, while edit-based metrics such as ERRANT (Felice et al., 2016; Bryant et al., 2017) and GoToScorer, as well as n-grambased metrics such as GLEU (Napoles et al., 2015, 2016) and GREEN (Koyama et al., 2024), aggregate scores by accumulating the number of edits or *n*-grams. The corpus-level scores obtained through these methods can be converted into system rankings by sorting.

#### 2.2 How to Resolve the Gap?

Given that the SEEDA dataset uses TrueSkill as the aggregation method, we will close the gap by using TrueSkill for automated evaluation as well. First, since existing automatic evaluation metrics compute sentence-wise scores, we convert these scores into pairwise comparison results. For example, in the case illustrated in Figure 1, the evaluation scores of 0.8, 0.7, and 0.9 for corrected sentences corresponding to the first sentence ("He play a tennis") can be compared to produce pairwise comparison results similar to those in human evaluation. Next, we compute system rankings by applying TrueSkill to the transformed pairwise comparison results. In this study, we consider all combinations of pairwise comparisons for system set. That is, given N systems, a total of N(N-1) comparisons are performed per sentence, and system rankings are computed based on these results.

The similar method was employed by Kobayashi et al. (2024a), but they did not mention the gap. Also, their experiments used the TrueSkill aggregation for their proposed LLM-based metrics, but used conventional aggregation methods, e.g., averaging, for other metrics. We discuss and organize the gap between human and automatic evaluation in detail, and then solve the gap by applying TrueSkill to all metrics for fair comparison.

	SEEDA-S				SEEDA-E			
	Base		+Fluency		Base		+Fluency	
Metrics	r (Pearson)	$\rho$ (Spearman)	r	$\rho$	r	$\rho$	r	ho
w/o TrueSkill								
ERRANT	0.545	0.343	-0.591	-0.156	0.689	0.643	-0.507	0.033
PTERRANT	0.700	0.629	-0.546	0.077	0.788	0.874	-0.470	0.231
GLEU+	0.886	0.902	0.155	0.543	0.912	0.944	0.232	0.569
GREEN	0.925	0.881	0.185	0.569	0.932	0.965	0.252	0.618
SOME	0.892	0.867	0.931	0.916	0.901	0.951	0.943	0.969
IMPARA	0.916	0.902	0.887	0.938	0.902	0.965	0.900	0.978
Scribendi	0.620	0.636	0.604	0.714	0.825	0.839	0.715	0.842
w/ TrueSkill								
ERRANT	0.763	0.706	-0.463	0.095	0.881	0.895	-0.374	0.231
PTERRANT	0.870	0.797	-0.366	0.182	0.924	0.951	-0.288	0.279
GLEU+	0.863	0.846	0.017	0.393	0.909	0.965	0.102	0.486
GREEN	0.855	0.846	-0.214	0.327	0.912	0.965	-0.135	0.420
SOME	0.932	0.881	0.971	0.925	0.893	0.944	0.965	0.965
IMPARA	0.939	0.923	0.975	0.952	0.901	0.944	0.969	0.965
Scribendi	0.674	0.762	0.745	0.859	0.837	0.888	0.826	0.912
GPT-4-E (fluency)	0.844	0.860	0.793	0.908	0.905	0.986	0.848	0.987
GPT-4-S (fluency)	0.913	0.874	0.952	0.916	0.974	0.979	0.981	0.982
GPT-4-S (meaning)	0.958	0.881	0.952	0.925	0.911	0.960	0.976	0.974

Table 1: Correlation with human evaluation using the SEEDA dataset. *w/o TrueSkill* refers to the conventional evaluation procedure, while *w/ TrueSkill* represents the proposed evaluation procedure. Improvements over the conventional procedure are underlined, and the highest value in each column is highlighted in bold. The GPT-4 results refer to those reported in Kobayashi et al. (2024b).

#### **3** Experiments

158

159

160

161

162

163

164

166

167

168

169

170

171

172

173

174

#### 3.1 Automatic Evaluation Metrics

We provide more detailed experimental settings for each metric in Appendix A.

**Edit-based metrics** We use ERRANT (Felice et al., 2016; Bryant et al., 2017) and PT-ERRANT (Gong et al., 2022). Both are reference-based evaluation metrics that assess at the edit level. When multiple references are available, the reference that yields the highest  $F_{0.5}$  score is selected for each sentence.

*n*-gram based metrics We use GLEU+ (Napoles et al., 2015, 2016) and GREEN (Koyama et al., 2024). The *n*-gram overlap is checked among the input sentence, hypothesis sentence, and reference sentence. When multiple references are available, the average score across all references is used.

Sentence-level metrics SOME (Yoshimura et al., 175 2020), IMPARA (Maeda et al., 2022), and 176 Scribendi Score (Islam and Magnani, 2021) are 177 used. All of them are based on small neural mod-179 els such as BERT<sub>base</sub> (Devlin et al., 2019) and designed as a reference-free metric that considers 180 the correction quality estimation score as well as 181 the meaning preservation score between the input and corrected sentences. 183

#### 3.2 Meta-Evaluation Method

We use SEEDA dataset (Kobayashi et al., 2024b) for meta-evaluation. Meta-evaluation results are reported based on human evaluation results using TrueSkill for both the sentence-level humanevaluation, SEEDA-S, and the edit-level humanevaluation, SEEDA-E. Additionally, we also report results for both the Base configuration, which excludes reference sentences and GPT-3.5 outputs that allow for larger rewrites, and the +Fluency configuration, which includes them.

Furthermore, we evaluate the robustness of the calculated rankings using window analysis (Kobayashi et al., 2024a). The window analysis computes correlation coefficients only for consecutive N systems, after sorting systems based on human evaluation results. This allows us to analyze whether automatic evaluation can correctly assess a set of systems that appear to have similar performance from the human evaluation. In this study, we perform it with N = 8 for 14 systems corresponding to the +Fluency configuration, and report both Pearson and Spearman correlation coefficients. That is, correlation coefficients are computed for the rankings 1 to 8, 2 to 9, ..., and 7 to 14 from human evaluation. 185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

201

202

203

204

205

208

209

#### **3.3 Experimental Results**

210

211

213

214

215

216

217

218

219

229

230

234

235

241

242

244

246

247

248

251

253

257

259

Table 1 shows the results of the meta-evaluation. The upper group presents evaluation results based on the conventional method of averaging or summing, and the bottom group presents results evaluated using TrueSkill, which follows the same evaluation method as human evaluation. The bottom group includes the evaluation results based on GPT-4 reported by Kobayashi et al. (2024a), which correspond to the state-of-the-art metrics.

The overall trend indicates that using TrueSkillbased evaluation improves the correlation coefficients for most of metrics. In particular, the results of IMPARA in the SEEDA-S and +Fluency setting outperformed those of GPT-4 results. Additionally, ERRANT showed an improvement of more than 0.2 points in many configurations. These results show that using automatic evaluation metrics with the same evaluation procedure as human evaluation make the ranking closer to human evaluation. In other words, the existing automatic evaluation metrics were underestimated in the prior reports due to the gap of the meta evaluation procedure. On the other hand, no effect was observed in ngram-based metrics such as GLEU+ and GREEN. Since human evaluation is not conducted based on *n*-grams, aligning the evaluation procedure likely led to negative effects due to differences in the granularity of evaluation. This result is consistent with Kobayashi et al.'s (2024b) claiming that aligning the granularity of evaluation between automatic and human evaluation is important.

Figure 2 shows the results of the window analysis for IMPARA and ERRANT measured on SEEDA-S and SEEDA-E, respectively. From Figure 2a, it can be seen that IMPARA particularly aligns with human evaluation in the lower ranks. The Pearson correlation coefficient also showed an improvement in the evaluation results for the top systems as well. Since the top systems include GEC systems that are largely rewritten, such as GPT-3.5, this characteristic is useful considering that LLM-based correction methods will become popular in the future. Figure 2b shows that ER-RANT consistently showed improved correlation coefficients with the proposed method, but still struggled with evaluating the top systems. For editbased evaluation metrics, it is still considered difficult to assess such GEC systems even with the evaluation method aligned with human evaluation  $^{1}$ .



(b) ERRANT

Figure 2: The results of the window analysis for N = 8 are shown. The x-axis represents the starting rank of human evaluation. For example, x = 2 shows the results for the systems ranked 2nd to 10th in human evaluation.

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

278

279

280

281

283

#### 4 Conclusion

In this study, we focused on the fact that human evaluation aggregates sentence-level scores into system rankings based on TrueSkill, while automatic evaluation uses a different evaluation, and we proposed to use TrueSkill in automatic evaluation as well. Results with various existing metrics showed improvements of correlations with human evaluation for many of the metrics, indicating that agreement on the aggregation method is important. We also release extensible implementations and expect aggressive development of metrics in the future<sup>2</sup>.

Given the results so far, we recommend transitioning the aggregation method from averaging or summing to using a rating algorithm, such as TrueSkill. We also recommend that evaluation metrics should be developed that allow for accurate sentence-wise comparisons. This is evidenced by the fact that IMAPARA achieves a higher correlation cofficients than SOME in Table 1. In fact, IMAPARA is trained to assess the pairwise comparison results, whereas SOME is trained to evaluate sentences absolutely.

<sup>&</sup>lt;sup>1</sup>Using more number of references may solve this issue.

<sup>&</sup>lt;sup>2</sup>In the writing of this paper, we partially used an AI assistant to improve text.

### 284 Limitations

296

297

300

301

302

305

307

310

311

312

313

314

315

319

320

321

322

323

324

326

327

329

332 333

Use for Purposes Other Than System Ranking The proposed method is designed for system ranking and cannot be used for other types of evaluation, such as analyzing the strengths and weaknesses of a specific system. For instance, when analyzing whether a model excels in precision or recall, it is more useful to accumulate the number of edits at the corpus level, as done in existing evaluation methods.

**Reproducing the Outputs of Compared GEC Systems** Since the proposed ranking method requires inputting all GEC outputs being compared, it is necessary to reproduce their models. This point is different from existing absolute evaluation methods, where previously reported scores can be cited. While this may seem burdensome for researchers, it can also be seen as an important step toward promoting the publication of reproducible research results.

# Ethical Considerations

When the metric contains social biases, the proposed method cannot eliminate that bias and may reflect that bias in the rankings. However, we argue that this problem should be resolved as a metric problem.

#### References

- Abhijeet Awasthi, Sunita Sarawagi, Rasna Goyal, Sabyasachi Ghosh, and Vihari Piratla. 2019. Parallel iterative edit models for local sequence transduction. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 4260–4270, Hong Kong, China. Association for Computational Linguistics.
- Christopher Bryant, Mariano Felice, and Ted Briscoe. 2017. Automatic annotation and evaluation of error types for grammatical error correction. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 793–805, Vancouver, Canada. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages

4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

- Mariano Felice, Christopher Bryant, and Ted Briscoe. 2016. Automatic extraction of learner errors in ESL sentences using linguistically enhanced alignments. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 825–835, Osaka, Japan. The COLING 2016 Organizing Committee.
- Peiyuan Gong, Xuebo Liu, Heyan Huang, and Min Zhang. 2022. Revisiting grammatical error correction evaluation and beyond. In *Proceedings of the* 2022 Conference on Empirical Methods in Natural Language Processing, pages 6891–6902, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Roman Grundkiewicz, Marcin Junczys-Dowmunt, and Edward Gillian. 2015. Human evaluation of grammatical error correction systems. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 461–470, Lisbon, Portugal. Association for Computational Linguistics.
- Ralf Herbrich, Tom Minka, and Thore Graepel. 2006. Trueskill<sup>TM</sup>: A bayesian skill rating system. In *Advances in Neural Information Processing Systems*, volume 19. MIT Press.
- Md Asadul Islam and Enrico Magnani. 2021. Is this the end of the gold standard? a straightforward referenceless grammatical error correction metric. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3009–3015, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Masahiro Kaneko and Naoaki Okazaki. 2023. Reducing sequence length by predicting edit spans with large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 10017–10029, Singapore. Association for Computational Linguistics.
- Satoru Katsumata and Mamoru Komachi. 2020. Stronger baselines for grammatical error correction using a pretrained encoder-decoder model. In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, pages 827–832, Suzhou, China. Association for Computational Linguistics.
- Masamune Kobayashi, Masato Mita, and Mamoru Komachi. 2024a. Large language models are state-ofthe-art evaluator for grammatical error correction. In *Proceedings of the 19th Workshop on Innovative Use of NLP for Building Educational Applications (BEA* 2024), pages 68–77, Mexico City, Mexico. Association for Computational Linguistics.

334

335

336

337

338

339

341

342

343

345

346

347

348

349

350

351

352

354

355

356

357

358

359

360

361 362

363

364

365

366

367

369

370

371

372

373

374

375

376

377

378

379

380

384

386

387

388

- 38 39
- 39<sup>.</sup>
- 392

398

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424 425

426

427

428

429

430

431

432

433

434 435

436

437

438

439

440

441

442 443

444

Masamune Kobayashi, Masato Mita, and Mamoru Komachi. 2024b. Revisiting meta-evaluation for grammatical error correction. *Transactions of the Association for Computational Linguistics*, 12:837–855.

- Shota Koyama, Ryo Nagata, Hiroya Takamura, and Naoaki Okazaki. 2024. n-gram F-score for evaluating grammatical error correction. In *Proceedings of the 17th International Natural Language Generation Conference*, pages 303–313, Tokyo, Japan. Association for Computational Linguistics.
- Mengsay Loem, Masahiro Kaneko, Sho Takase, and Naoaki Okazaki. 2023. Exploring effectiveness of GPT-3 in grammatical error correction: A study on performance and controllability in prompt-based methods. In *Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)*, pages 205–219, Toronto, Canada. Association for Computational Linguistics.
  - Koki Maeda, Masahiro Kaneko, and Naoaki Okazaki. 2022. IMPARA: Impact-based metric for GEC using parallel data. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 3578–3588, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
  - Courtney Napoles, Keisuke Sakaguchi, Matt Post, and Joel Tetreault. 2015. Ground truth for grammatical error correction metrics. In Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers), pages 588–593, Beijing, China. Association for Computational Linguistics.
  - Courtney Napoles, Keisuke Sakaguchi, Matt Post, and Joel Tetreault. 2016. Gleu without tuning. *Preprint*, arXiv:1605.02592.
  - Kostiantyn Omelianchuk, Vitaliy Atrasevych, Artem Chernodub, and Oleksandr Skurzhanskyi. 2020. GECToR – grammatical error correction: Tag, not rewrite. In Proceedings of the Fifteenth Workshop on Innovative Use of NLP for Building Educational Applications, pages 163–170, Seattle, WA, USA → Online. Association for Computational Linguistics.
  - 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.
- Keisuke Sakaguchi, Matt Post, and Benjamin Van Durme. 2014. Efficient elicitation of annotations for human evaluation of machine translation. In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, pages 1–11, Baltimore, Maryland, USA. Association for Computational Linguistics.

Ryoma Yoshimura, Masahiro Kaneko, Tomoyuki Kajiwara, and Mamoru Komachi. 2020. SOME: Reference-less sub-metrics optimized for manual evaluations of grammatical error correction. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6516–6522, Barcelona, Spain (Online). International Committee on Computational Linguistics. 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

492

493

# A Detailed Experimental Settings for Evaluation Metrics

For ERRANT, we used the Python module errant=3.0.0 and evaluated it with the Spanbased Correction setting. The reference edits were manually provided, but these were reextracted and used by ERRANT. PT-ERRANT performed weighting based on BERTScore with the bert-base-uncased model, and the weights were calculated based on the F1 score. At this point, rescaling was performed using the baseline, and no adjustments were made using idf. Similar to ERRANT, reference edits were re-extracted by ER-RANT. GLEU used up to a maximum of 4-grams with 500 iterations. The seed values for sampling reference sentences followed the official implementation settings. For GREEN, up to 4-grams were used, and the evaluation metric  $F_{2,0}$  was applied. SOME used the official pre-trained model with weights for grammaticality, fluency, and meaning preservation set to 0.43, 0.55, and 0.02, respectively, following the official implementation<sup>3</sup>. For IMPARA, since the quality estimation model was not publicly available, we conducted a reproduction implementation and experiment. Following the original paper (Maeda et al., 2022), we generated 4,096 training pairs using the CoNLL-2013 corpus as the seed corpus and split them into an 8:1:1 ratio and regards training set, development set, and test set, respectively. We fine-tuned bert-base-cased on the training data. The architecture followed the BertForSequenceClassification model in the transformers library, and the representation corresponding to the CLS token was linearly transformed into a real-valued output. For inference, we also used bert-base-cased for the similarity estimation model, and the threshold for the estimated value was set to 0.9. Scribendi Score used the GPT-2 language model<sup>4</sup> and performed inference with a threshold of 0.8 for the maximum value of Levenshtein distance ratio and token sort ratio.

<sup>&</sup>lt;sup>3</sup>https://github.com/kokeman/SOME

<sup>&</sup>lt;sup>4</sup>https://huggingface.co/openai-community/gpt2