# Improving Explainability of Sentence-level Metrics via Edit-level Attribution for Grammatical Error Correction

Anonymous ACL submission

#### Abstract

Various evaluation metrics have been proposed 001 for Grammatical Error Correction (GEC), but many, particularly reference-free metrics, lack 004 explainability. This lack of explainability hinders researchers from analyzing the strengths and weaknesses of GEC models and limits the ability to provide detailed feedback for users. 800 To address this issue, we propose attributing sentence-level scores to individual edits, providing insight into how specific corrections con-011 tribute to the overall performance. For the attribution method, we use Shapley values, from 012 cooperative game theory, to compute the contribution of each edit. Experiments with ex-014 isting sentence-level metrics demonstrate high consistency across different edit granularities and show approximately 70% alignment with 018 human evaluations. In addition, we analyze biases in the metrics based on the attribution results, revealing trends such as the tendency to ignore orthographic edits. Our implementation is available at LINK<sup>1</sup>.

## 1 Introduction

027

034

Grammatical error correction (GEC) is the task of automatically correcting grammatical or superficial errors in an input sentence. Automatic evaluation metrics play a key role in improving GEC performance, but their effectiveness depends on their level of explainability. For example, metrics that evaluate at the edit level are more explainable than sentence-level metrics, as they allow us to identify which specific edits are effective and which are not, even when a GEC system makes multiple edits. Such explainable metrics enable researchers to analyze the strengths and weaknesses of GEC models, providing valuable insights into how models can be improved. Furthermore, in education applications, explainable metrics can provide language learners with detailed feedback on their writing, supporting their learning more effectively.



(b) Our proposed method improves explainability.

Figure 1: Overview of the proposed method with an example using three edits. Figure (a) shows the low-explainability of existing metrics that only estimate the sentence-level score, but Figure (b) shows that the edit-level attribution solves this issue.

041

042

043

044

045

046

047

049

051

052

054

057

060

In GEC, explainable reference-based metrics, such as ERRANT (Felice et al., 2016; Bryant et al., 2017) are limited because references cannot account for all valid corrections. Preparing test data with comprehensive references is often impractical, especially when targeting domains such as medical or academic writing that differ from existing datasets. To address this issue, reference-free metrics have been proposed to evaluate corrected sentences without relying on references (Choshen and Abend, 2018; Yoshimura et al., 2020; Islam and Magnani, 2021; Maeda et al., 2022). Although these reference-free metrics achieve high correlation with human evaluations, many are designed to assign scores at the sentence level, limiting their explainability on individual edits. This lack of granularity makes it difficult to analyze how specific edits contribute to the overall sentence score. For example, as shown in Figure 1, a metric evaluates a corrected sentence created by applying the three

<sup>&</sup>lt;sup>1</sup>Here will be replaced with an actual link in camera-ready.

063

101

102

103

104

105

106

107

109

edits. As shown in Figure 1a, the sentence-level metric assigns an overall score of 0.75, but it does not indicate whether all edits are valid, or if both valid and invalid edits have been applied.

To improve the explainability of metrics with low or no explanation, we propose attributing sentencelevel scores to individual edits as illustrated in Figure 1b. In the proposed method, the total contribution of all edits is calculated as the difference between the scores of the input sentence and the corrected sentence. This difference is then attributed to the individual edits. For example, in Figure 1b, a difference of -0.05 is distributed among three edits with contributions of 0.2, 0.1, and -0.35. The attribution results are intrepreted using the sign and magnitude of these scores: the sign indicates whether an edit is the valid or invalid, while the 077 magnitude represents the degree of its influence on the final sentence-level score. We employ Shapley values (Shapley et al., 1953) from cooperative game theory to fairly distribute the total score among the edits. By considering various combinations edits, Shapley values allow us to precisely attribute each edit's contribution to the overall sentence score, offering insights into their individual impact. Unlike previous feature attribution methods (Lundberg and Lee, 2017; Sundararajan et al., 2017), the proposed method is novel in attributing the difference between the input sentence and the corrected sentence.

In the experiments, we apply the proposed method to two popular reference-free metrics, SOME (Yoshimura et al., 2020) and IM-PARA (Maeda et al., 2022), as well as a fluency metric based on GPT-2 (Radford et al., 2019) perplexity. The results show that the proposed attribution method produces consistent scores across different granularities of edits and that edits with larger absolute attribution scores align more closely with human evaluations. We introduce Shapley sampling values (Strumbelj and Kononenko, 2010) to mitigate the time-complexity issues of calculating Shapley values. Additionally, we demonstrate that the proposed method can explain metric decisions at both the sentence and corpus levels, categorized by error types. These analyses reveal the types of edits that metrics give more weight to, as well as provide insights into the strengths and weaknesses of GEC systems.

## 2 Background

Edits in GEC. The GEC task aims to correct grammatical errors in a source sentence S and output a corrected sentence H. The differences between S and H are often represented as N edits  $e = \{e_i\}_{i=1}^N$  to enable evaluation (Dahlmeier and Ng, 2012; Bryant et al., 2017; Gong et al., 2022; Ye et al., 2023), ensembling (Tarnavskyi et al., 2022), and post-processing (Sorokin, 2022) at the edit level. These edits can be automatically extracted using edit extraction methods (Felice et al., 2016; Bryant et al., 2017; Belkebir and Habash, 2021; Korre et al., 2021; Uz and Eryiğit, 2023). Each edit typically includes a word-level span in S and its corresponding correction, although it may also include an error type (Bryant et al., 2017). The error type categorizes each edit, indicating the partof-speech or grammatical aspect it relates to, which helps to analyze the strengths and weaknesses of the GEC systems.

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

127

128

129

130

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

157

Sentence-level Metrics. A sentence-level metric M computes the score of the corrected sentence given the source sentence, denoted as  $M(H|S) \in \mathbb{R}$ . The source sentence is used to assess meaning preservation, as GEC requires correcting errors while maintaining the original meaning of the source sentence. This formulation has been adopted by several reference-free metrics (Yoshimura et al., 2020; Islam and Magnani, 2021; Maeda et al., 2022; Kobayashi et al., 2024a). Sentence-level metrics aim to rank GEC systems in alignment with humans judgments, as evidenced by the fact that the meta-evaluation is performed using the correlation between metric-generated rankings or scores and those of humans. However, these metrics are limited to sentence-level scoring and cannot explain how individual edits contribute to the final score.

## 3 Method

Our attribution method assumes that the overall contribution of edits is the difference in scores before and after correction. We distribute the difference  $\Delta M(H|S) = M(H|S) - M(S|S)$  across each edit  $e = \{e_i\}_{i=1}^N$ , where M(S|S) is the score of the source sentence treated as its own corrected sentence.

The goal of our attribution method is to compute the contribution for each edit denoted as  $\{\phi_i(M) \in$ 

197

198 199

200

201

204

205

206

207

209

210

211

212 213

214

215

216

217

218

219

220

221

222

224

225

226

227

228

230

231

232

233

234

235

236

 $\mathbb{R}$ <sub>*i*=1</sub>, so that the following equation is satisfied:

158

159

178

179

181

182

184

186

187

188

$$\Delta M(H|S) = \sum_{i=1}^{N} \phi_i(M).$$
 (1)

We refer to  $\phi_i(M)$  as attribution scores. A positive 160 score  $(\phi_i(M) > 0)$  indicates an edit that improves 161 the metric  $M(\cdot)$ , while a negative score ( $\phi_i(M) <$ 162 0) indicates an edit that worsens it. The absolute value  $|\phi_i(M)|$  represents the degree of the edit's 164 impact. 165

Shapley. For the attribution method, we intro-166 duce Shapley values (Shapley et al., 1953) from 167 cooperative game theory. In cooperative game theory, multiple players work together towards a common goal and share the total benefit based on their 170 contributions. Shapley values distribute this benefit 171 among players fairly, ensuring that those players 172 who contributes more receive a larger share. For 173 our purpose, we regard  $\Delta M(H|S)$  as the total ben-174 efit, edits e as the players, and  $\phi_i(M)$  as the Shap-175 176 ley values. The Shapley value  $\phi_i(M)$  for a given metric  $M(\cdot)$  is calculated as follows: 177

$$\phi_{i}(M) = \sum_{e' \subseteq e \setminus \{e_{i}\}} \frac{|e'|!(N - |e'| - 1)!}{N!} \\ (\Delta M(S_{e' \cup \{e_{i}\}}|S) - \Delta M(S_{e'}|S)),$$
(2)

where  $S_e$  denotes the source sentence after applying the edit set e. Equation 2 calculates the weighted sum of the differences in evaluation scores when including and excluding the edit  $e_i$ . For example, us- $= \{e_1, e_2, e_3\}$ ing Figure 1 with e= $\{[A \rightarrow The], [job \rightarrow work], [is \rightarrow was]\}, one of$ the terms in the calculation for  $\phi_1(M)$  with e' = $\{e_2\}$  is

$$\frac{1}{6} \left( \Delta M(S_{\{e_1, e_2\}} | S) - \Delta M(S_{\{e_2\}} | S) \right)$$
  
=  $\frac{1}{6} \left( \Delta M(\text{The work is performed by him.} | S) - \Delta M(\text{A work is performed by him.} | S) \right).$   
(3)

Here, bold words indicate the edit being attributed, and underlined words show other edits. The terms 190 for  $e' = \{\phi\}, \{e_3\}$ , and  $\{e_2, e_3\}$  are computed in a 191 similar way. Shapley values consider various com-192 binations of edits, ensuring accurately attribution of the *i*-th edit's contribution. By design, Shapley 194

values naturally satisfy Equation 1 due to their effectiveness (Shapley et al., 1953). However, the computational complexity is  $\mathcal{O}(2^N)$ .

Shapley Sampling Values. To improve computational efficiency, we introduce Shapley sampling values (Strumbelj and Kononenko, 2010), an approximation of Shapley values. Equation 2 can be rewritten as:

$$\phi_i(M) = \frac{1}{N!} \sum_{\boldsymbol{o} \in \pi(\boldsymbol{e})} (\Delta M(S, S_{\operatorname{Pre}^i(\boldsymbol{o}) \cup \{e_i\}})) - \Delta M(S, S_{\operatorname{Pre}^i(\boldsymbol{o})}))$$
(4)

where  $\pi(e)$  is the set of all possible orders of edits, and  $Pre^{i}(o)$  is the set of edits preceding  $e_i$  in permutation o. In the example from Equation 3,  $Pre^{1}(o) = \{\phi\}$ when  $o = [e_1, e_2, e_3]$ , and  $Pre^1(o)$ =  $\{e_2, e_3\} = \{[\text{job} \rightarrow \text{work}], [\text{is} \rightarrow \text{was}]\}$  when  $o = [e_3, e_2, e_1]$ . To approximate Shapley values, we uniformly sample T permulations without replacement from  $\pi(e)$ , denoted as  $\pi(e) = \{o_1, \ldots, o_T\}$ . Shapley sampling values are then calculated using  $\pi(e)$  instead of  $\pi(e)$  in Equation 4. This approximation reduces the computational cost from  $\mathcal{O}(2^N)$  to  $\mathcal{O}(TN)$ .

Normalized Shapley Values The calculated attribution scores are not directly comparable across different sentence-level scores. For instance, an attribution score of 0.2 has a different relative impact when distributing a sentence-level score of 1.0 versus -0.05. To enable meaningful comparison, we apply L1 normalization to the attribution scores:

$$\phi_i^{\text{norm}}(M) = \frac{\phi_i(M)}{\sum_{i=1}^N |\phi_i(M)|}.$$
 (5)

This normalization, applied as a post-processing step, adjusts only the magnitude of the scores while preserving their original signs. Since the normalized scores represent the ratio of each edit's contribution, they are assumed to be comparable even when the sentence-level scores differ.

#### **Evaluation of Attribution** 4

We evaluate the proposed attribution method from two perspectives: faithfulness and explainability (Wang et al., 2024). Faithfulness measures how well the attribution results reflect the model's internal decision, while explainability assesses the

extent to which the results are understandable to humans. To demonstrate the effectiveness of the proposed method across various domains, we conduct
experiments using diverse datasets, GEC systems, and metrics.

## 4.1 Experimental Settings

## 4.1.1 Datasets

242

245

246

247

248

249

251

254

259

260

261

271

272

274

275

278

279

We use the CoNLL-2014 test set (Ng et al., 2014) and the JFLEG validation set (Heilman et al., 2014; Napoles et al., 2017). CoNLL-2014 is a benchmark for minimal edits, focusing on correcting errors while preserving the original structure of the input as much as possible. In contrast, JFLEG is a benchmark for fluency edits, allowing more extensive rewrites to produce fluent and natural sentences.

### 4.1.2 GEC Systems

We evaluate our attribution method on various GEC systems, including two tagging-based models (the official RoBERTa-based GECToR (Omelianchuk et al., 2020) and GECToR-2024 (Omelianchuk et al., 2024)), two encoder-decoder models (BART (Lewis et al., 2020) and T5 (Rothe et al., 2021)), and a causal language model (GTP-40 mini). This allows us to assess the explainability of attributions scores across different GEC architectures. For GPT-40 mini, we used a two-shot setting following Coyne et al. (2023), with examples randomly sampled once from the W&I+LOCNESS validation set (Yannakoudakis et al., 2018) and used for all input sentences. Note that we use only the corrected sentences containing 10 or fewer edits  $(N \leq 10)$  due to the computational complexity of Shapley values. According to Figure 2, which shows the cumulative sentence ratio by the number of edits, our experiments cover at least more than 97% of the sentences in both datasets.

4.1.3 **Reference-free Metrics** 

**SOME (Yoshimura et al., 2020)** trains a BERTbased regression model optimized directly on human evaluation results. We used the official pretrained model weights<sup>2</sup> and used the default coefficients for the weighted average of grammaticality, fluency, and meaning preservation scores, from the official script<sup>3</sup>.



Figure 2: Cumulative sentences ratio regarding the number of edits. The red line indicates the position where the number of edits is 10.

281

282

283

284

285

287

290

291

292

293

294

295

297

298

299

301

302

303

304

305

306

307

308

310

311

312

**IMPARA (Maeda et al., 2022)** estimates evaluation scores through similarity estimation and quality estimation. We use BERT (bert-base-cased) as the similarity estimator and train our own model for the quality estimator, as the official pre-trained weights are not available. Our quality estimator was trained following the same settings described in Maeda et al. (2022), achieving a correlation with the human ranking comparable to their reported results.

**GPT-2 Perplexity (PPL).** Our proposed method can be applied to metrics that evaluate only the quality of the corrected sentence<sup>4</sup>. To test this, we use GPT-2 (Radford et al., 2019) perplexity, with negative perplexity scores to ensure that higher values correspond to better quality.

#### 4.2 Baseline Attribution Methods

To evaluate the effectiveness of Shapley values, we employ simpler variants, i.e., ADD and Sub, as baseline attribution methods.

Add. This method observes the change in the score when each edit is applied individually to the source sentence. An edit that increases the score is considered valid for the metric. This approach corresponds to using only  $e' = \{\phi\}$  in Equation 2, with the attribution scores normalized by  $\frac{\Delta M(H|S)}{\sum_{i=1}^{N} \phi_i(M)}$  so that it satisfies Equation 1.

**Sub.** This method observes the change in the score when each edit is removed individually from the corrected sentence. An edit that decreases the score upon removal is considered valid for the metric. This approach corresponds to using only

<sup>&</sup>lt;sup>2</sup>https://github.com/kokeman/SOME

 $<sup>^{3}0.55*</sup>$ grammaticality + 0.43 \* fluency + 0.02 \* meaning preservation.

<sup>&</sup>lt;sup>4</sup>In this case, the sentence-level score is  $\Delta M(S, H) = M(H) - M(S)$ 



Figure 3: The results of consistency-based evaluation. Each row shows the different datasets and each column shows different metrics. "Mag." means the magnitude. Colors show the attribution scores.

 $e' = e \setminus \{e_i\}$  in Equation 2, with the attribution scores normalized by  $\frac{\Delta M(H|S)}{\sum_{i=1}^{N} \phi_i(M)}$  so that it satisfies Equation 1.

## 4.3 Consistency Evaluation

313

314

315

316

317

319

321

327

331

333

334

338

339

341

To evaluate faithfulness, we test how well the attribution scores represent the judgments of the metrics through consistency evaluation. Specifically, we first calculate the attribution scores for individual edits and then group edits with the same sign, treating them as a single edit. Next, we calculate the attribution score for the grouped edits. We hypothesize that the attribution score for a grouped edit should equal the sum of the individual attribution scores of the edits comprising the group. If this condition holds, the attribution method consistently calculates the contributions of edits, making its results reliable for practical use. We use an agreement ration to measure the consistency of the signs and use Pearson and Spearman correlations to assess the consistency of the magnitudes.

For example, in Figure 1, we group two positivity-attributed edits,  $[A \rightarrow The]$  and  $[job \rightarrow work]$ , into a single edit and compute attribution scores for the grouped edit and the remaining edit,  $[is \rightarrow was]$ . Ideally, the attribution score for the grouped edit should be 0.2 + 0.1 = 0.3, which can be verified by sign agreement and closeness to 0.3.

Figure 3 presents the results for each metrics. Our proposed Shapley method shows higher consistency than the baseline attribution methods across various domains and metrics. While the Sub metric also demonstrates high consistency, its Spearman's rank correlation occasionally drops for certain metrics, such as IMPARA. Low rank correlation can misrepresent the relative importance of edits, posing a serious issue for explainability. These results suggest that the attribution method is reliable across different edit granularities, such as edits extracted by ERRANT (Felice et al., 2016; Bryant et al., 2017) or chunks created by merging multiple edits (Ye et al., 2023). This flexibility enables a wide range of applications for the proposed method. 342

343

344

345

346

347

349

350

351

352

353

355

356

357

359

360

361

362

363

364

365

366

367

368

369

371

## 4.4 Human Evaluation

To evaluate explainability, we assess the agreement between attribution scores and human evaluation results using references. Ideally, a positively attributed edit should align with a correct edit in the reference-based evaluation, while a negativity attributed edit should correspond to an incorrect one. Furthermore, edits with larger absolute attribution scores are expected to show higher agreement with human evaluations.

In this experiment, we annotate two types of labels for each edit: one based on the sign of the attribution score and another based on referencebased evaluation. We then calculate the matching ratio between these labels at the corpus level. For the evaluation, we use the two official references for CoNLL-2014, and four official references

for JFLEG validation set. The assessment is per-372 formed on mixed outputs from five GEC systems. 373 To ensure the analysis focuses on meaningful cases, 374 we include only sentences with two or more edits. When assigning labels for reference-based evaluation with multiple references, we select the reference that results in the highest agreement with the 378 attribution scores. To further examine the relationship between the magnitude of attribution scores and agreement rates, we follow standard attribu-381 tion evaluation practices (Petsiuk, 2018; Fong and Vedaldi, 2017) by applying a threshold to the absolute values of the scores. We use only edits with normalized absolute attribution scores below the threshold for accuracy calculations. The threshold starts at 0.1 and increases in steps of 0.1 until it reaches 1.0, where all edits are included.

> Figure 4 presents the results for the CoNLL-2014 and JFLEG datasets. Overall, the results show that including edits with larger absolute attribution scores improves the agreement with human evaluation, indicating that the magnitude of these scores is meaningful. When comparing attribution methods, Shapley rarely achieves the worst agreement. For instance, in JFLEG, the SOME metric shows the order Add > Shapley > Sub, while the IMPARA metric shows Sub > Shapley > Add. Either Add or Sub often results in the worst agreement, whereas Shapley demonstrates more stable performance across different metrics and domains.

394

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

When comparing metrics, particularly in the results for JFLEG (Figure 4b), the agreement rates consistently rank in the order of PPL, SOME, and IMPARA. This trend may reflect the characteristics of these reference-free metrics in relation to reference-based evaluation. In fact, when we compute the correlation with ERRANT <sup>5</sup> using standard sentence-level meta-evaluation (Kobayashi et al., 2024b), the rankings follow the same order: of PPL (0.550), SOME (0.529), and IMPARA (0.516), with Kendall rank correlation coefficients of 0.100, 0.058, and 0.033, respectively. These results suggest that metrics more closely aligned with reference-based evaluation can be attributed more accurately, improving the reliability of our attribution method. On the other hand, for CoNLL-2014, the sentence-level correlation shows the order of PPL (0.522), IMPARA (0.479), and SOME (0.477). However, the agreement in Figure 4a does not fol-





Figure 4: Human evaluation results for CoNLL-2014 and JFLEG. Colors indicate metrics and line styles indicate attribution methods.

low this trend. This indicates that the proposed method aligns well with human judgement in case of fluency edits. Conversely, minimal edits may require further studies, but primarily depend on the development of better reference-free metrics.

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

#### 4.5 Efficiency of Shapley Values

One limitation of Shapley values is their high computational cost. Figure 5 shows the relation between the number of edits and the computation time per sentence in seconds on a single RTX 3090. The computation time increases rapidly when the number of edits exceeds 11. For this reason, we assume that sentences with more than 11 edits are impractical to attribute within a reasonable time. According to Figure 2, the affects approximately 3% of the sentences in GEC output. Similarly, tasks involving a higher number of edits, such as text simplification, could face even greater challenges.

As discussed in Section 3, we address this is-

Original (S) Correction (H)	-	Further more	,	by with	these this	evidence evidence	,	u you	will agree will agree	
Metrics $(M)$	$\Delta M(\cdot)$		Shapley values $\phi_i(M)$							
SOME	0.298	-	0.068	0.064	0.033	-	0.038	0.066	-	0.030
IMPARA	-0.027	-	0.068	0.029	0.124	-	0.145	-0.361	-	-0.033
PPL	1266.3	-	250.7	103.8	216.0	-	67.4	366.6	-	261.5
		Normalized Shapley values								
SOME		-	0.229	0.215	0.111	-	0.126	0.220	-	0.099
IMPARA		-	0.090	0.039	0.163	-	0.191	-0.475	-	-0.043
PPL		-	0.198	0.082	0.171	-	0.053	0.290	-	0.207

Table 1: An example of the proposed method's results using actual sentence.



Figure 5: The relationship between the number of edits and computation time per sentence. The solid lines are average time and ranges are standard deviation.

442

443

444

445

446 447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

sue by employing Shapley sampling values and evaluate their ability to approximate exact Shapley values by measuring the average absolute differences between them. For system-independent experiments, we use a dataset combining all GEC model corrections on the JFLEG validation set. We set T = 64 and restrict sentences to  $10 \le N \le 15$ 

Table 2 reports the errors and computation times for each metric. With Shapley sampling values, the computation time per sentence can be reduced to as little as one second. To assess the impact of errors, we also show the distribution of absolute original Shapley values. While SOME and PPL show errors below the average, IMPARA exhibits higher errors. This discrepancy with IMPARA can lead to misinterpretations of attribution scores. For example, the frequency of changes in the relative contributions of different edits is likely to increase, undermining reliability. IMPARA's higher error rate may be due to its smaller variance in evaluated values, making

Metric	Error	Time	Shapley values dist.
SOME	0.014	1.00	$\begin{array}{c} 0.019 \pm 0.020 \\ 0.052 \pm 0.071 \\ 34.549 \pm 59.472 \end{array}$
IMPARA	0.066	0.99	
PPL	17.515	0.20	

Table 2: The average error and average computation time (seconds) when using Shapley sampling values. It also shows the distribution of the absolute original Shapley values (the average  $\pm$  the standard deviation).

it less effective at quantifying impact with a limited number of calculations.

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

## 5 Applications of Attribution Scores

We demonstrate practical applications of attribution scores for users. All results in this section are based on Shapley values for the attribution method.

#### 5.1 Case Study

Attribution scores can be used to identify which edits improve or worsen the sentence-level score. Table 1 provides an example, showing attribution scores and their normalized version. The original sentence and its corrections are chunked according to edit spans, omitting scores for non-edited chunks which are all zeros. One observation is that the sentence-level score of IMPARA declines primarily due to the edit  $[u \rightarrow you]$ , as identified by the attribution score. In contrast, SOME and PPL prefer this edit. This analysis demonstrates how attribution scores can reveal weaknesses in metrics as seen in Table 1.

Normalized Shapley values enable comparison of attribution scores across metrics. For example, while SOME and IMPARA assign the same Shapley value to the edit  $[\phi \rightarrow ]$ , their normalized scores reveal differing impacts. This feature is particularly useful for comparing metrics with different value ranges, such as SOME and PPL.

<sup>&</sup>lt;sup>6</sup>When T = 64 and  $10 \le N$ , the computation cost of Shapley sampling values is consistently lower than that of exact Shapley values, as  $2^x > 64x$  holds for x > 9.20...



Figure 6: The heatmap indicating the average of normalized Shapley values per error type. The deeper color indicates higher values.

However, the metrics themselves may exhibit biases that affect attribution scores. To investigate these biases, we calculate the average normalized Shapley values for each error type (Bryant et al., 2017). As in Section 4.5, we combine the corrected sentences from five GEC systems for the JFLEG validation set to mitigate biases specific to individual GEC models. Figure 6 shows a heatmap of average normalized attribution scores for error types with a frequency greater than 30. The results indicate that different metrics emphasize different error types. For instance, orthography (ORTH) edits, such as case changes and whitespace adjustments, tend to be downplayed. Metric biases must be considered when interpreting attribution scores. It is important to not that the attribution scores reflect the internal decisions of the metric and may not align with the true correctness of edits. We leave addressing these biases to future work.

### 5.2 Precision per Error Type

While the case study focused on local, sentencelevel evaluation, the proposed method can be extended to corpus-level analysis. Typically, metrics with low explainability provide only a single numerical score at the corpus level. By applying the proposed method, we can decompose this score is into performance across different error types. Specifically, we treat edits with positive attribution scores as True Positives, and those with negative attribution scores as False Positives, enabling the calculation of precision for each error type. To handle attribution scores across multiple sentences, we use normalized Shapley values:

$$Precision = \frac{\phi_+^{\text{norm}}(M)}{\phi_+^{\text{norm}}(M) + |\phi_-^{\text{norm}}(M)|}, \quad (6)$$

where  $\phi_{+}^{\text{norm}}(M)$  and  $\phi_{-}^{\text{norm}}(M)$  represent the sum of positive and negative normalized attribution scores at the corpus-level, respectively. This is



Figure 7: The heatmap indicating the precision for each GEC systems. We used JFLEG validation set as a dataset and SOME as a metric.

similar to PT-M2 (Gong et al., 2022) which proposed an edit-level weighted evaluation. However, our method is designed to enhance the corpus-level explainability of metrics rather than to improve agreement with human evaluations. 526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

564

Figure 7 shows the precision for each error type using the JFLEG validation set and SOME as the evaluation metric. The parentheses in the y-axis labels indicate the corpus-level scores, with each row of the heatmap explaining these score in terms of error types. The results reveal that better edits in adverbs (ADV) or orthography (ORTH) contribute most to the highest corpus-level score achieved by GPT-40 mini. On the other hand, despite achieving the highest corpus-level score among the five systems, GPT-40 mini's precisions are not particularly high. Notably, T5 appears to perform better in terms of precision, as indicated by more darkcolored cells. This discrepancy may stem from an overcorrection issue, leading to a low-precision, high-recall trend in performance (Fang et al., 2023; Omelianchuk et al., 2024). While this trend is intuitive because the valid edits in the reference-based evaluation are limited to the references, we also observe a similar trend even for reference-free evaluation metrics.

## 6 Conclusion

This paper proposes a method to improve the explainability of existing low-explainable GEC metrics by attributing sentence-level scores to individual edits. Specifically, we employed Shapley values to perform attribution while accounting for various contexts in which edits are applied. Quantitative evaluations showed that the attribution scores align with metric's judgement achieve approximately 70% agreement with human evaluations. Additionally, we demonstrated how attribution scores can be used at both the sentence and corpus levels. Finally, we discussed the biases of existing metrics.

510

512

513

514

515

516

517

518

519

489

490

491

492

493

## Limitations

565

571

572

574

592

593

597

598

599

Treating False Negative Corrections. As men-566 tioned in Section 5.2, the proposed method is lim-567 ited to analyzing corrections made by the GEC 568 system, i.e. True Positives (TP) and False Positives (FP), and does not address False Negatives (FN). While we assume that the effect of FN corrections is canceled out by  $\Delta M(H|S) = M(H|S) -$ M(S|S), it may still affect the computation of attribution scores. A more detailed investigation into this issue is left for future work.

Treating dependent edits Edits might exhibit dependencies. For example, the correction [model 577 's prediction -> prediction of the model] can be split into two dependent edits: [model 's ->  $\phi$ ] and  $[\phi \rightarrow of \text{ the model}]$ . While analyzing these edits together may better capture their contribution, the proposed method evaluates each edit independently. 582 We assume that Shapley values partially capture 583 such dependent edits by considering various pat-584 terns of applying edits. However, understanding 585 dependencies fully requires error correction data annotated for edit dependencies or tools to automatically identify them. Developing such resources is 588 left as future work. 589

**Real Human Evaluation** Unlike Section 4.4, which uses a reference-based evaluation framework, we could also conduct direct human evaluation. However, we prioritize reference-based evaluation for its scalability when applying the method to new metrics or datasets. It is important to note that the primary goal of this study is not to derive attribution scores that align with human evaluation, but to explain the decision-making process of metrics at the edit level. Verifying alignment with human evaluations is a secondary finding. If the goal were to achieve consistency with human evaluation, training a dedicated model would be a more appropriate approach.

**Rectifying Metric Biases** The case study results (Section 5.1) revealed that metrics exhibit 605 biases towards specific error types. While one could attempt to mitigate such biases, we believe that sentence-level metrics benefit from implicitly weighting edits, making these biases beneficial for evaluation. However, biases related to social fac-610 tors such as gender or nationality, should be ad-611 dressed. A deeper investigation into metric biases is beyond the scope of this work, but remains an im-613

portant area for future research. Our work provides	614
a strong foundation for exploring these biases	615
Acknowledgments	616
References	617
Riadh Belkebir and Nizar Habash. 2021. Automatic	618
error type annotation for Arabic. In <i>Proceedings of</i>	619
<i>the 25th Conference on Computational Natural Lan-</i>	620
<i>guage Learning</i> , pages 596–606, Online. Association	621
for Computational Linguistics.	622
Christopher Bryant, Mariano Felice, and Ted Briscoe.	623
2017. Automatic annotation and evaluation of error	624
types for grammatical error correction. In <i>Proceed-</i>	625
<i>ings of the 55th Annual Meeting of the Association for</i>	626
<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	627
pages 793–805, Vancouver, Canada. Association for	628
Computational Linguistics.	629
Leshem Choshen and Omri Abend. 2018. Reference-	630
less measure of faithfulness for grammatical error	631
correction. In Proceedings of the 2018 Conference of	632
the North American Chapter of the Association for	633
Computational Linguistics: Human Language Tech-	634
nologies, Volume 2 (Short Papers), pages 124–129,	635
New Orleans, Louisiana. Association for Computa-	636
tional Linguistics.	636
Steven Coyne, Keisuke Sakaguchi, Diana Galvan-Sosa,	638
Michael Zock, and Kentaro Inui. 2023. Analyzing	639
the performance of gpt-3.5 and gpt-4 in grammatical	640
error correction. <i>Preprint</i> , arXiv:2303.14342.	641
Daniel Dahlmeier and Hwee Tou Ng. 2012. Better	642
evaluation for grammatical error correction. In <i>Pro-</i>	643
<i>ceedings of the 2012 Conference of the North Amer-</i>	644
<i>ican Chapter of the Association for Computational</i>	645
<i>Linguistics: Human Language Technologies</i> , pages	646
568–572, Montréal, Canada. Association for Compu-	647
tational Linguistics.	648
Tao Fang, Shu Yang, Kaixin Lan, Derek F. Wong, Jin-	649
peng Hu, Lidia S. Chao, and Yue Zhang. 2023. Is	650
chatgpt a highly fluent grammatical error correc-	651
tion system? a comprehensive evaluation. <i>Preprint</i> ,	652
arXiv:2304.01746.	653
Mariano Felice, Christopher Bryant, and Ted Briscoe.	654
2016. Automatic extraction of learner errors in ESL	655
sentences using linguistically enhanced alignments.	656
In <i>Proceedings of COLING 2016, the 26th Inter-</i>	657
<i>national Conference on Computational Linguistics:</i>	658
<i>Technical Papers</i> , pages 825–835, Osaka, Japan. The	659
COLING 2016 Organizing Committee.	660
Ruth C Fong and Andrea Vedaldi. 2017. Interpretable explanations of black boxes by meaningful pertur-	661 662

663 664

665

666

667

- Ruth explanations of black boxes by meaningful perturbation. In Proceedings of the IEEE international conference on computer vision, pages 3429-3437.
- Peiyuan Gong, Xuebo Liu, Heyan Huang, and Min Zhang. 2022. Revisiting grammatical error correction evaluation and beyond. In Proceedings of the

780

781

2022 Conference on Empirical Methods in Natural Language Processing, pages 6891–6902, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Michael Heilman, Aoife Cahill, Nitin Madnani, Melissa Lopez, Matthew Mulholland, and Joel Tetreault.
2014. Predicting grammaticality on an ordinal scale. In Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pages 174–180, Baltimore, Maryland. Association for Computational Linguistics.

673

674

675

676

677

679

680

681

684

685

688

694

697

701

703

705

706

708

709

710

711

714

715

716

718

719

720

721

722

725

- 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.
  - 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.
  - Masamune Kobayashi, Masato Mita, and Mamoru Komachi. 2024b. Revisiting meta-evaluation for grammatical error correction. *Preprint*, arXiv:2403.02674.
  - Katerina Korre, Marita Chatzipanagiotou, and John Pavlopoulos. 2021. ELERRANT: Automatic grammatical error type classification for Greek. In Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021), pages 708–717, Held Online. INCOMA Ltd.
  - 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.
  - Scott M Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. *Advances in neural information processing systems*, 30.
  - 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, and Joel Tetreault. 2017. JFLEG: A fluency corpus and benchmark for grammatical error correction. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics:*

*Volume 2, Short Papers*, pages 229–234, Valencia, Spain. Association for Computational Linguistics.

- Hwee Tou Ng, Siew Mei Wu, Ted Briscoe, Christian Hadiwinoto, Raymond Hendy Susanto, and Christopher Bryant. 2014. The CoNLL-2014 shared task on grammatical error correction. In *Proceedings of the Eighteenth Conference on Computational Natural Language Learning: Shared Task*, pages 1–14, Baltimore, Maryland. Association for Computational Linguistics.
- 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  $\rightarrow$ Online. Association for Computational Linguistics.
- Kostiantyn Omelianchuk, Andrii Liubonko, Oleksandr Skurzhanskyi, Artem Chernodub, Oleksandr Korniienko, and Igor Samokhin. 2024. Pillars of grammatical error correction: Comprehensive inspection of contemporary approaches in the era of large language models. In *Proceedings of the 19th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2024)*, pages 17–33, Mexico City, Mexico. Association for Computational Linguistics.
- V Petsiuk. 2018. Rise: Randomized input sampling for explanation of black-box models. *arXiv preprint arXiv:1806.07421*.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- 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.
- Lloyd S Shapley et al. 1953. A value for n-person games.
- Alexey Sorokin. 2022. Improved grammatical error correction by ranking elementary edits. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 11416–11429, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Erik Strumbelj and Igor Kononenko. 2010. An efficient explanation of individual classifications using game theory. *J. Mach. Learn. Res.*, 11:1–18.
- Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks. In *International conference on machine learning*, pages 3319– 3328. PMLR.

Maksym Tarnavskyi, Artem Chernodub, and Kostiantyn Omelianchuk. 2022. Ensembling and knowledge distilling of large sequence taggers for grammatical error correction. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics* (Volume 1: Long Papers), pages 3842–3852, Dublin, Ireland. Association for Computational Linguistics.

782

783

785

791

793 794

795

796

797

802

807

810

811

812 813

814

815

816 817

818

819

- Harun Uz and Gülşen Eryiğit. 2023. Towards automatic grammatical error type classification for Turkish. In Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics: Student Research Workshop, pages 134– 142, Dubrovnik, Croatia. Association for Computational Linguistics.
- Yongjie Wang, Tong Zhang, Xu Guo, and Zhiqi Shen. 2024. Gradient based feature attribution in explainable ai: A technical review. *Preprint*, arXiv:2403.10415.
- Helen Yannakoudakis, Øistein E Andersen, Ardeshir Geranpayeh, Ted Briscoe, and Diane Nicholls. 2018.
  Developing an automated writing placement system for esl learners. *Applied Measurement in Education*, 31(3):251–267.
- Jingheng Ye, Yinghui Li, Qingyu Zhou, Yangning Li, Shirong Ma, Hai-Tao Zheng, and Ying Shen. 2023. CLEME: Debiasing multi-reference evaluation for grammatical error correction. In *Proceedings of the* 2023 Conference on Empirical Methods in Natural Language Processing, pages 6174–6189, Singapore. 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.