

# EXCGEC: A Benchmark of Edit-wise Explainable Chinese Grammatical Error Correction

Anonymous ACL submission

## Abstract

Existing studies explore the explainability of Grammatical Error Correction (GEC) in a limited scenario, where they ignore the interaction between corrections and explanations. To bridge the gap, this paper introduces the task of EXplainable GEC (EXGEC), which focuses on the integral role of both correction and explanation tasks. To facilitate the task, we propose EXCGEC, a tailored benchmark for Chinese EXGEC consisting of 8,216 explanation-augmented samples featuring the design of hybrid edit-wise explanations. We benchmark several series of LLMs in multiple settings, covering post-explaining and pre-explaining. To promote the development of the task, we introduce a comprehensive suite of automatic metrics and conduct human evaluation experiments to demonstrate the human consistency of the automatic metrics for free-text explanations.<sup>1</sup>

## 1 Introduction

Despite the notable advancements in Grammatical Error Correction (GEC) (Bryant et al., 2023; Zhao et al., 2018; Bryant et al., 2019a), there still exists a lack of profound examination into the explainability of GEC (Dwivedi et al., 2023), which is critical in educational scenarios for L2 (Language second)-speakers (Wang et al., 2021) or school-age children (Li et al., 2023b). These mainstream users, who often face challenges in creating grammatically accurate and fluent texts, may be confused or even misguided if they are provided with limited access to only corrective texts. Therefore, augmenting the explainability of GEC is unquestionably beneficial for the progression of the GEC community as well as related fields, such as essay scoring (Ashiya Katuka et al., 2024; Stahl et al., 2024), intelligent tutoring systems (Montenegro-Rueda et al., 2023) and other emerging educational scenarios (Lan et al., 2024).

<sup>1</sup>All the codes and data will be released after the review.

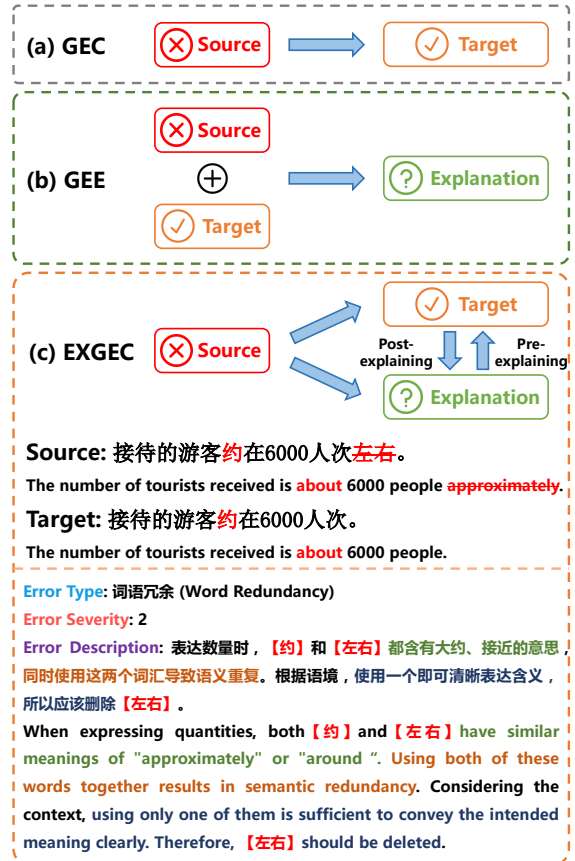


Figure 1: Task definitions of GEC, GEE and EXGEC. For the error description of EXGEC, we highlight **evidence words**, **linguistic knowledge**, **error causes**, and **revision advice** parts in different colors.

As illustrated in Figure 1, existing tasks like GEC and Grammatical Error Explanation (GEE) typically address either correction or explanation, ignoring the interaction between the two. To bridge the gap, we introduce the task of EXplainable Grammatical Error Correction (EXGEC). By integrating these two tasks, EXGEC enables systems to elucidate the linguistic knowledge and reasoning mechanism underlying predicted corrections, thereby achieving the best of both worlds. Additionally, EXGEC can function as a test bed for

051	determining the explainable abilities of large language models (LLMs) and identifying any unintended biases and risks in educational scenarios.	alignment workload for the LLMs. Our contributions in this paper are listed as follows:	103
052			104
053			
054	To facilitate EXGEC, we present <b>EXCGEC</b> , a tailored benchmark for Chinese EXGEC, featuring the design of hybrid edit-wise explanations. Each explanation, based on a particular edit, consists of three elements: 1) <i>Error types</i> , which allow learners to absorb syntax and semantic knowledge in an inductive way (Fei et al., 2023). We establish a hierarchical and pragmatic two-tier taxonomy for Chinese grammatical errors. 2) <i>Error severity levels</i> ranging from 1 ~ 5 points, which are beneficial to prioritize core corrections. 3) <i>Error descriptions</i> , presented as the form of natural language explanation (Camburu et al., 2018; He et al., 2023), provide evidence words, relevant linguistic knowledge or syntax rules, error causes, and revision advice for edits. The design provides more detailed and faithful guidance for learners, allowing them to comprehend each grammatical error committed. This is unlikely achievable for other designs such as example-based (Kaneko et al., 2022) or sentence-level explanations (Nagata et al., 2021).	(1) We introduce the EXGEC task and establish a corresponding benchmark consisting of a Chinese EXGEC dataset and a comprehensive set of automatic metrics, contributing to the stable development of the field of EXGEC.	105
055			106
056			107
057			108
058			109
059		(2) We develop EXGEC baseline models and investigate the abilities of various LLMs using our proposed benchmark.	110
060			111
061			112
062		(3) We conduct detailed analyses on our proposed dataset and baselines to gain further insights. Human evaluation experiments are also conducted to confirm the effectiveness of automatic metrics for error descriptions.	113
063			114
064			115
065			116
066			117
067			
068		<b>2 Related Work</b>	118
069		Exploration of explainable GEC has witnessed a paradigm shifting from fine-tuning (Kaneko and Okazaki, 2023) to prompting (Zhao et al., 2024), with the focus being local explanations of individual predictions. Fei et al. (2023) construct an explainable GEC dataset EXPECT, which is annotated with evidence words and error types based on the standard GEC benchmark (Bryant et al., 2019b). However, EXPECT falls short of flexibility due to the lack of natural language explanations. To fill the gap, Song et al. (2023) propose the task of grammatical error explanation. They observe that GPT-4 suffers from identifying and explaining errors with limited access to only parallel source-target pairs. To address this issue, they fine-tune an extra LLM as an edit extractor, which is trained on synthesized data. However, all these studies overlook the benefits of effectiveness and efficiency brought by multi-task learning both correction and explanation tasks, which is extensively explored in this work.	119
070			120
071			121
072			122
073			123
074			124
075			125
076	Stimulated by recent success of synthetic data generation (Shum et al., 2023; Whitehouse et al., 2023), we employ a semi-automatic dataset construction solution to enhance efficiency, while minimizing annotation costs. Initially, we synthesize the EXCGEC dataset by prompting GPT-4 (Liu et al., 2024). Then we hire native annotators to filter invalid data and provide a detailed analysis of invalid data, ensuring the quality of the dataset (Ding et al., 2024). We finally obtain 8,216 clean explanation-augmented samples for benchmarking. We also introduce automatic metrics to evaluate performance across both tasks, and conduct human evaluation experiments to ascertain the correlation between these metrics and human judgement, thus demonstrating their effectiveness.		126
077			127
078			128
079			129
080			130
081			131
082			132
083			133
084			134
085			135
086			136
087			137
088			138
089			139
090			140
091			141
092			142
093			143
094			144
095			145
096			146
097			
098		<b>3 Task Definition</b>	147
099		<b>3.1 Grammatical Error Correction</b>	148
100		GEC has been studied for decades, witnessing the shift from rule-based methods to LLM-based meth-	149
101			150
102			

ods. Formally, given an ungrammatical text (source text)  $X = \{x_1, x_2, \dots, x_T\}$ , a GEC model is required to correct  $X$  into a grammatically correct counterpart (target text)  $Y = \{y_1, y_2, \dots, y_{T'}\}$  without changing the original semantic as far as possible. Typically, GEC is usually treated as a sequence-to-sequence (Seq2Seq) task, the training objective of which is formulated as follow:

$$\mathcal{L}_{\text{GEC}} = - \sum_{t=1}^{T'} \log P(y_t | Y_{<t}, X) \quad (1)$$

### 3.2 Grammatical Error Explanation

GEE has been explored in several methodologies, including sentence-level explanation and edit-wise explanation. Since sentence-level explanation suffer from over-generalization and confusion especially when a sentence contains multiple grammatical errors, this work focuses solely on edit-wise explanations. Given a source text  $X$  and its target counterpart  $Y$ , the GEE model needs to explain each grammatical error  $e_i$  in  $X$ . Specifically, GEE is typically solved in a two-step pipeline consisting of edit extraction and edit-wise explanation. 1) **Edit extraction** produces an edit set  $E = \{e_1, e_2, \dots, e_n\}$  that represent grammatical errors in  $X$  and also clarify the transformation from ungrammatical segments of  $X$  to target segments of  $Y$ . Typically, an edit contains four key elements: source position  $sp$ , source content  $sc$ , target position  $tp$ , and target content  $tc$ . The process of edit extraction can be easily accomplished using alignment-based evaluation toolkits like ERRANT (Bryant et al., 2017; Felice et al., 2016) and CLEME (Ye et al., 2023). 2) **Edit-wise explanation** generates a set of explanations  $E' = \{e'_1, e'_2, \dots, e'_n\}$ , with each explanation  $e'_i$  corresponding to  $e_i$ , given the source and the target texts. Although the design of explanation varies across related work (Song et al., 2023; Zhao et al., 2024), the typical training objective of GEE models is presented as follows:

$$E = f(X, Y) \quad (2)$$

$$\mathcal{L}_{\text{GEE}} = - \sum_{i=1}^n \log P(e'_i | X, Y, e_i) \quad (3)$$

where  $f : (X, Y) \rightarrow E = \{(sp_i, sc_i, tp_i, tc_i)\}_{i=1}^n$  is the edit extraction function used to extract edits of  $X$  and  $Y$ , and  $n$  is the number of edits.

Existing studies (Song et al., 2023; Fei et al., 2023) focus on developing GEE models that can generate more reasonable explanations. However, an extra GEC model is compulsory to allow GEE models to generate explanations if only source texts are offered, thus resulting in an issue of low efficiency. Furthermore, there exists a gap between GEC and GEE models if they are trained on different data with domain shift.

### 3.3 Explainable Grammatical Error Correction

To get rid of the drawbacks brought by the natures of GEE, we propose the EXGEC task which aims to perform both correction and explanation tasks simultaneously. The motivation of combining these two tasks majorly falls on two aspects. First, a branch of existing studies (Wiegrefe and Marasovic, 2021; Hartmann and Sonntag, 2022; Li et al., 2022, 2024) have demonstrated training with access to human explanations can improve model performance. It is also intuitive that either of GEC and GEE tasks can mutually benefit from each other when training in a multi-task manner. Second, it is more time-saving and cost-efficient to deploy a single EXGEC model rather than two detached models in foreign language education platforms.

In this task, the only input element is an ungrammatical source text  $X$ , and the EXGEC model learns to output both the grammatically target text  $Y$  and explanations  $E'$ . Similar to GEE, EXGEC follows the edit-wise style of explanation, and it is categorized into two different settings by the order of correction and explanation tasks, with the basic scheme of multi-task learning.

**Post-explaining.** Models are trained first to generate target texts (Camburu et al., 2018), which allows the explanations to be explicitly conditioned on the target texts, thus ensuring high faithfulness of explanations towards the target texts. The training objective is as follows:

$$\mathcal{L}_{\text{post}} = - \sum_{t=1}^{T'} \log P(y_t | Y_{<t}, X) - \sum_{i=1}^n \log P(e'_i | X, Y, e_i) \quad (4)$$

The inference of post-explaining models is represented as follows:

$$\hat{Y} = \text{EXGEC}_{\text{post}}(X) \quad (5)$$

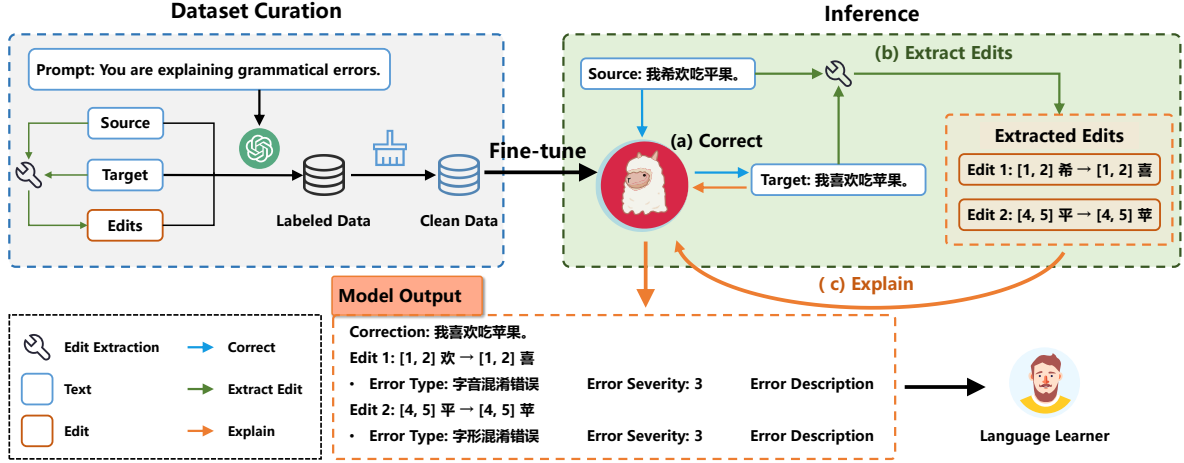


Figure 2: Overview of benchmark construction and model development. We show the inference process of the post-explaining model in particular.

$$\hat{E}' = \text{EXGEC}_{\text{post}}(X, Y, f(X, \hat{Y})) \quad (6)$$

With the target texts generated ahead, post-explaining models can output explanations conditioned on the specific edits that are extracted by aligning the source and the target texts, thus improving accuracy and faithfulness of explanations.

**Pre-explaining.** This type of models are trained conversely, whose mechanism is similar to the Chain of Thought (CoT) technique. Pre-explaining models are supposed to make full use of synthesized explanations to generate elaborated target texts. With minimal modification from Equation (4), the training objective of pre-explaining models is as follow:

$$\mathcal{L}_{\text{pre}} = - \sum_{i=1}^n \log P(e'_i | X) - \sum_{t=1}^{T'} \log P(y_t | Y_{<t}, X, E') \quad (7)$$

Notably, pre-explaining models may struggle to generate well-formed edit-wise explanations due to the inaccessibility to the edit extraction function  $f$ , which necessitate both the source and the target texts. Similarly, the inference of pre-explaining models is presented as follows:

$$\hat{E}' = \text{EXGEC}_{\text{pre}}(X) \quad (8)$$

$$\hat{Y} = \text{EXGEC}_{\text{pre}}(X, E') \quad (9)$$

## 4 EXCGEC Benchmark

To facilitate the development of EXGEC task, we construct EXCGEC, the first benchmark for ex-

plainable Chinese GEC particularly. As illustrated in Figure 2, we begin by the process of data curation, which consists of explanation design in Section 4.1, explanation synthesis and refinement in Section 4.2. Then we gain in-depth understanding of GPT-4 (Achiam et al., 2023) in EXGEC by further analyzing generated data in Section 4.3, where we summarize common failure modes in invalid instances. Finally, we introduce a series of automatic metrics for evaluating explanations in Section 4.4.

### 4.1 Explanation Design

In the pursuit of comprehensiveness and plausibility, we adopt a hybrid strategy for edit-wise explanations, where each edit is explained through three aspects, including error type labels, error severity levels, and free-text error description. 1) **Error type labels** allow language learners to comprehend and infer syntax and grammar rules in an inductive manner. In particular, we employ a two-tier hierarchical taxonomy including 5 major types and 16 minor types shown in Table 1, inspired by existing studies (Liping, 2014; Peng et al., 2021; Zhang et al., 2022). The detailed description of various error types are included in Appendix A.1 and A.2. If an edit covers multiple error types, we select the one with the highest granule. 2) **Error severity levels**, ranging from 1 to 5 points, indicate the significance of a specific grammatical error. 3) **Error descriptions** are the most crucial and flexible element. These provide keywords, pertinent linguistic knowledge, causes of errors, and revision guidance in a free-text format. We stipulate well-defined error description should meet three principles: fluency, reasonability (making sense to humans), and

Major Type	Minor Type	
<b>Punctuation-level Error</b>	标点冗余 (Punctuation Redundancy)	
	标点丢失 (Punctuation Missing)	
	标点误用 (Punctuation Misuse)	
<b>Spelling-level Error</b>	字音混淆错误 Phonetic Confusion Error	
	字形混淆错误 Glyph Confusion Error	
	词内部字符异位错误 Internal Character Misplacement Error	
	命名实体拼写错误 Named Entity Misspelling	
	<b>Word-level Error</b>	词语冗余 (Word Redundancy)
		词语丢失 (Word Missing)
词语误用 (Word Misuse)		
<b>Sentence-level Error</b>	词序不当 (Improper Word Order)	
	逻辑不通 (Illogicality)	
	句式杂糅 (Run-on Sentence)	
<b>Other Special Error</b>	照应错误 (Inconsistency Error)	
	歧义错误 (Ambiguity Error)	
	语气不协调 (Inconsistent Tone)	
<b>Other</b>		

Table 1: Hierarchical taxonomy of grammatical error types defined in our benchmark.

faithfulness (targeted to a specific edit). To ensure the reasonability and faithfulness, the error description must mostly conform to the syllogism form of deductive reasoning: [*major premise: semantic rules and related knowledge*], [*minor premise: the reason for the error in the text*], and [*explain how to correct it*]. Further, any evidence from the source  $X$  must be enclosed within special markers  $\llbracket \ \rrbracket$ . Similarly, correction content that occurs in the target sentence  $Y$  must be enclosed within  $\{ \}$ , as indicated in Figure 1.

## 4.2 Explanation Synthesizing

Annotating high-quality explanations in a large scale poses a huge challenge to our benchmark construction. Hence, we leverage GPT-4 to synthesize edit-wise explanations efficiently. To achieve this, we first select 10,000 parallel samples across 6 existing benchmarks or datasets of Chinese GEC, with the details listed in Table 2. We pick out only the samples with changed target sentences, and select the single target sentence with the most edits if a sample is annotated with multiple target sentences. Then, we prompt GPT-4 to generate edit-wise explanations following in-context learning. To ensure faithfulness of synthesized explanation, we first extract edits using the toolkit CLEME (Ye et al., 2023). Inspired by Li et al. (2022), we then employ the Rationalization Prompting (RP) strategy, where we concatenate task definition, demon-

Dataset	Sentences	Edits/Sent.	Chars/Sent.
FCGEC	41,340	1.0	53.1
YACL- <b>minimal-dev</b>	1,839	2.9	25.9
MuCGEC- <b>dev</b>	1,137	3.2	38.5
NaCGEC- <b>dev</b>	500	1.1	56.2
NLPCC- <b>test</b>	2,000	2.0	29.7
HSK	156,870	1.4	27.2
EXCGEC (FCGEC)	2,308	1.1	55.1
EXCGEC (YACL)	1,235	3.5	24.3
EXCGEC (MuCGEC- <b>dev</b> )	789	3.3	40.4
EXCGEC (NaCGEC- <b>dev</b> )	449	1.1	56.1
EXCGEC (NLPCC- <b>test</b> )	1,611	1.7	28.9
EXCGEC (HSK)	1,824	2.1	32.0
EXCGEC	8,216	2.0	38.8

Table 2: Dataset statistics of the EXCGEC benchmark.

strations, and a parallel sample  $(X, Y)$  with extracted edits  $E = \{e_1, e_2, \dots, e_n\}$  as the prompt. For each error type, we provide the definition, a suggested template of error description, and a demonstration. The prompt is listed in Appendix A.3.

## 4.3 Explanation Refinement and Analysis

Benefiting from the extensive knowledge acquired during the large-scale pre-training process, GPT-4 is able to generate fluent, reasonable and plausible explanations in most cases, meeting the requirements with specified instructions. However, GPT-4 is not guaranteed to produce all high-quality explanations due to hallucination, and the patterns of those invalid explanations are referred to as failure modes. Therefore, we hire 12 native speakers, all of whom are graduated students, to screen out invalid explanations. We finally obtain 8,216 clean samples out of 10,000 samples. We further investigate the failure modes of invalid explanations generated by GPT-4, which is provided in Appendix A.4.

## 4.4 Automatic Metrics

Recent studies leverage human evaluation for evaluation of GEE due to the lack of enough annotated samples, posing a challenge for efficient development of EXGEC systems. In this paper, we introduce a comprehensive set of automatic metrics for both correction and explanation parts.

**Correction.** We employ CLEME and ChERRANT to evaluate the correction performance. Both are edit-based metrics that output P/R/ $F_{0.5}$  scores, and they have been proven reliable metrics for GEC on CoNLL-2014 (Ye et al., 2023).

**Explanation.** Since an edit-wise explanation consists of three critical elements, we define respectively automatic metrics for them. 1) Accuracy

---

**Algorithm 1** COTE Decoding Algorithm

---

**Input:** Source text  $X$ , a post-explaining model  $\mathcal{M}$ , and the edit extraction function  $f$ .

**Output:** Target text  $\hat{Y}$ , and explanations  $\hat{E}'$ .

```
1:  $\hat{Y} \leftarrow \text{BeamSearch}(\mathcal{M}(\text{Json}(X)))$ 
2:  $\hat{E}' \leftarrow \emptyset$ 
3: if  $\hat{Y} = X$  then
4:   return  $\hat{Y}, \hat{E}'$ 
5: end if
6:  $E \leftarrow f(X, \hat{Y})$ 
7:  $\hat{E}' \leftarrow \text{Top-P}(\mathcal{M}(\text{Json}(X, Y, E)))$ 
8: return  $\hat{Y}, \hat{E}'$ 
```

---

and Macro-F1 scores are computed for error type clarification, following the conventional evaluation protocol of text clarification (Li et al., 2020). 2) We report mean absolute error (MAE) to show the deviation of hypothesis error severity levels towards ground truth ones. 3) We employ various metrics for evaluating the free-text explanation description, including BLEU, METEOR, ROUGE-1, ROUGE-2, and ROUGE-L. We leave the analysis on effectiveness of these metrics to Section 7.2.

## 5 Method

### 5.1 Training

To streamline the training process covering all the tasks mentioned in Section 3, we treat all of them as a unified Seq2Seq task. To achieve this, we linearize the data in the format of json (Gao et al., 2023). This structured approach simplifies the process of output parsing involving three types elements of edit-wise explanations, and provides a consistent and controllable view to distinguish tasks, enabling the model understand essential task elements and their relations. Therefore, we train all models using the same smooth cross entropy loss, regardless of the specific task.

### 5.2 Inference

For post-explaining EXGEC models, we design a specific **Correct-Then-Explain** decoding algorithm called **COTE**, which is presented in Algorithm 1. First, we employ the greedy beam search decoding strategy for the correction part, which is beneficial to relieve the over-correction problem that is common on LLMs (Cao et al., 2023; Loem et al., 2023; Li et al., 2023a). Then, we apply CLEME to extract edits. Notably, we merge adjacent edits with distance less than 2 characters to avoid fragmented edits. Finally, we leverage the Top-p decoding strategy for generating explanations,

encouraging diversified natural language explanations. It is worth noting that COTE is not accessible to pre-explaining models since the edit extraction tool necessitates both a source text and a target text.

## 6 Experiments

### 6.1 Experimental Settings

**Backbones.** We benchmark three series of LLMs, including Llama-3 (Touvron et al., 2023), Qwen-1.5 (Bai et al., 2023), and DeepSeek (Bi et al., 2024). For each series of LLMs, we experiment with their base and chat (or instruct) versions to investigate whether further alignment training benefits the task. All results are based on EXCGEC-test. Training details are reported in Appendix B.1.

**Evaluation.** We report experiment results using the metrics introduced in Section 4.4, calculated using open-source toolkits including *NLTK* (Bird and Loper, 2004), *rouge* (Lin, 2004), and *scikit-learn* (Pedregosa et al., 2011). Particularly, we observe many hypothesis edits are not covered in references, making it impossible to evaluate the subsequent explanations for these edits. To address this, we introduce two extra indicators, namely *Hit* and *Miss* rates. A hypothesis edit overlapping with a reference edit is designated as a hit edit, while a reference edit without any match with hypothesis edits is deemed a miss edit. The hit rate is defined as the ratio of hit edits to all hypothesis edits, and the miss rate as the ratio of miss edits to all reference edits. Only hit edits are used to determine the evaluation outcomes for explanations.

### 6.2 Results of Multi-task Models

The preliminary results from both the post- and pre-explaining models are presented in Table 3, from which we can make some conclusions.

**Post-explaining models consistently outperform pre-explaining models.** In relation to the correction aspect, all post-explaining models obtain higher  $F_{0.5}$  scores than pre-explaining models, regardless of the applied backbones. A similar pattern is observed in the explanation part, where all the pre-explaining models invariably underperform their post-explaining counterparts. This suggests that a complexity for LLMs in initially explaining grammatical errors. And once pre-explaining models generate flawed explanations, the ensuing distraction impedes their ability to accurately correct the source text.

Model	Correction $\uparrow$			Explanation						
	CLEME (P / R / F <sub>0.5</sub> )	ChERRANT (P / R / F <sub>0.5</sub> )	Hit $\uparrow$	Miss $\downarrow$	Acc $\uparrow$	F1 $\uparrow$	MAE $\downarrow$	BLEU $\uparrow$	METEOR $\uparrow$	ROUGE-(1 / 2 / L) $\uparrow$
Qwen1.5-7B-base	26.00 / 26.54 / <b>26.10</b>	33.87 / <b>20.16</b> / 29.81	67.29	56.81	60.99	<b>29.82</b>	0.80	15.22	<b>39.05</b>	49.74 / 23.28 / 34.32
Qwen1.5-7B-chat	<b>28.31</b> / 21.21 / <b>26.54</b>	<b>36.74</b> / 17.26 / <b>29.98</b>	68.94	64.83	<b>61.98</b>	29.62	<b>0.75</b>	<b>15.49</b>	38.88	<b>50.32</b> / <b>24.25</b> / <b>35.24</b>
Llama3-8B-base	20.92 / 23.60 / 21.40	28.81 / 17.78 / 25.63	61.54	58.38	58.39	25.12	0.91	14.54	37.84	49.53 / 23.19 / 34.58
Llama3-8B-instruct	21.33 / 26.05 / 22.14	29.00 / 19.40 / 26.39	61.40	<b>55.71</b>	59.16	25.63	0.88	14.70	36.89	49.41 / 23.54 / 34.87
DeepSeek-7B-base	26.21 / 7.00 / 16.92	36.00 / 7.04 / 19.75	<b>69.92</b>	85.39	60.64	26.47	0.79	15.07	38.05	50.19 / 24.10 / 34.90
DeepSeek-7B-chat	25.46 / 18.51 / 23.68	34.02 / 15.75 / 27.62	67.52	66.64	58.11	24.45	0.84	13.94	36.97	48.66 / 22.70 / 34.23
Qwen1.5-7B-chat	<b>13.76</b> / <b>13.42</b> / <b>13.69</b>	<b>19.27</b> / <b>9.93</b> / <b>16.22</b>	<b>29.49</b>	80.24	23.35	8.22	<b>1.17</b>	<b>7.75</b>	<b>27.67</b>	<b>40.47</b> / <b>15.00</b> / <b>28.20</b>
Llama3-8B-instruct	7.12 / 11.17 / 7.68	10.86 / 8.57 / 10.31	23.88	<b>73.06</b>	<b>24.31</b>	<b>8.78</b>	1.21	5.78	23.07	37.57 / 13.47 / 27.19
DeepSeek-7B-chat	9.93 / 8.26 / 9.55	14.28 / 7.07 / 11.86	24.72	78.67	19.12	5.84	1.29	5.91	23.95	37.59 / 13.11 / 26.78

Table 3: Main results of multi-task learning models. Results of post-explaining models are listed in the *top* block, while those of pre-explaining models are in the *bottom* block.

Model	Correction $\uparrow$			Explanation						
	CLEME (P / R / F <sub>0.5</sub> )	ChERRANT (P / R / F <sub>0.5</sub> )	Hit $\uparrow$	Miss $\downarrow$	Acc $\uparrow$	F1 $\uparrow$	MAE $\downarrow$	BLEU $\uparrow$	METEOR $\uparrow$	ROUGE-(1 / 2 / L)
Qwen1.5-7B-base	— / — / —	— / — / —	98.42	6.14	85.00	<b>43.32</b>	0.73	19.70	43.18	53.48 / 27.79 / 38.12
Qwen1.5-7B-chat	62.59 / 87.35 / 66.35	67.58 / 69.53 / 67.96	<b>99.93</b>	0.43	81.53	39.56	0.73	17.88	41.40	51.73 / 28.81 / 36.51
Llama3-8B-base	— / — / —	— / — / —	99.42	2.27	83.27	40.51	0.89	20.52	43.37	54.32 / 29.05 / 39.49
Llama3-8B-instruct	<b>69.10</b> / <b>90.90</b> / <b>72.58</b>	<b>73.75</b> / <b>74.37</b> / <b>73.87</b>	99.63	1.67	<b>85.99</b>	41.84	0.78	20.73	42.98	<b>54.60</b> / <b>29.64</b> / <b>40.04</b>
DeepSeek-7B-base	— / — / —	— / — / —	<b>99.93</b>	3.54	85.06	40.19	<b>0.71</b>	<b>20.78</b>	<b>43.48</b>	54.07 / 29.18 / 39.57
DeepSeek-7B-chat	41.12 / 79.02 / 45.48	48.35 / 53.20 / 49.25	<b>99.93</b>	<b>0.40</b>	81.17	35.93	0.74	19.57	42.32	53.12 / 28.03 / 38.59

Table 4: Ground truth results of multi-task learning models. We report the explanation performance (**right** block) of *post-explaining* models conditioned on source texts and ground truth target texts. Contrarily, we report the correction performance (**left** block) of *pre-explaining* models conditioned on source sentences and ground truth explanations.

Model	Correction $\uparrow$			Explanation						
	CLEME (P / R / F <sub>0.5</sub> )	ChERRANT (P / R / F <sub>0.5</sub> )	Hit $\uparrow$	Miss $\downarrow$	Acc $\uparrow$	F1 $\uparrow$	MAE $\downarrow$	BLEU $\uparrow$	METEOR $\uparrow$	ROUGE-(1 / 2 / L)
Post-explaining	28.31 / 21.21 / 26.54	36.74 / 17.26 / 29.98	68.94	64.83	61.98	29.62	<b>0.75</b>	15.49	38.88	50.32 / 24.25 / 35.24
Pre-explaining	13.76 / 13.42 / 13.69	19.27 / 9.93 / 16.22	29.49	80.24	23.35	8.22	1.17	7.75	27.67	40.47 / 15.00 / 28.20
Pipeline	<b>32.45</b> / <b>23.93</b> / <b>30.29</b>	<b>40.50</b> / <b>19.58</b> / <b>33.37</b>	<b>88.53</b>	<b>57.42</b>	<b>74.80</b>	<b>34.84</b>	<b>0.75</b>	<b>16.44</b>	<b>39.76</b>	<b>50.62</b> / <b>24.56</b> / <b>35.71</b>

Table 5: Comparison of the multi-task solutions and the GEC-GEE pipeline solution based on Qwen1.5-7B-chat.

**Chat models outperform base models.** For post-explaining models, we observe all chat or instruct models gain slightly higher F<sub>0.5</sub> correction scores, and they also marginally outperform their base version counterparts in the explanation task. It indicates that additional alignment training (Wang et al., 2023) can benefit the EXGEC task.

### 6.3 Ground Truth Results

In order to study the isolated performance of multi-task models, we provide part ground truth information in advance during the inference stage. Specifically, we provide ground truth target texts for post-explaining and report their performance of explanation. Conversely, we offer ground truth explanations for pre-explaining and report their performance of correction. This experimental setting allows for revealing the specialized performance, eliminating the distraction of previously generated contents. The results are presented in Table 4.

**For the task of explanation, two base models slightly outperform chat models.** Specially, the base version models of Qwen and DeepSeek exhibit a minor increase in performance over their chat/instruct counterparts on classifying error types and providing error descriptions. However, this is not true for Llama3, where the Llama3-instruct model obtain the highest Acc, METEOR and ROUGE scores. Also noteworthy is the significantly lower miss rates of chat/instruct models compared to base models, indicating a tendency for the latter to overlook explanations, even when ground truth target texts are available. These findings contradict the joint results in Table 3. We speculate the reason is base models may be more susceptible to low-quality self-generated corrections.

**Ground truth explanations tremendously improve correction performance.** Since the explanations include explicit clues for corrections such as evidence words and revision advice, it is effort-

Model	Correction $\uparrow$		Explanation							
	CLEME (P / R / F <sub>0.5</sub> )	ChERRANT (P / R / F <sub>0.5</sub> )	Hit $\uparrow$	Miss $\downarrow$	Acc $\uparrow$	F1 $\uparrow$	MAE $\downarrow$	BLEU $\uparrow$	METEOR $\uparrow$	ROUGE- (1 / 2 / L)
Beam search	28.31 / 21.21 / 26.54	36.74 / 17.26 / 29.98	99.22	19.05	83.93	44.48	0.71	22.71	44.28	55.55 / 32.26 / 42.34
Top-p	19.45 / 27.05 / 20.61	24.83 / 19.14 / 23.44	99.93	0.40	81.53	39.56	0.74	17.88	41.40	51.73 / 25.81 / 36.51

Table 6: Comparison of the post-explaining model with different token-wise decoding strategies. Note that the explanation performance is conditioned on ground truth target texts in order to exclude unrelated interference.

	Hit $\uparrow$	Miss $\downarrow$	Acc $\uparrow$	F1 $\uparrow$	MAE $\downarrow$	BLEU $\uparrow$	METE $\uparrow$	ROUGE- (1/2/L) $\uparrow$
w COTE	99.93	0.43	81.53	39.56	0.74	17.88	41.40	51.73 / 25.81 / 36.51
w/o COTE	49.64	54.01	42.51	17.77	0.93	11.53	33.81	46.35 / 19.34 / 31.28

Table 7: Ablation results of COTE from the same Qwen1.5-7B-chat post-explaining model.

	Pearson	Spearson
Human v.s. BLEU	0.9222	0.6571
Human v.s. METEOR	0.9280	0.7714
Human v.s. ROUGE-1	0.9464	0.8286
Human v.s. ROUGE-2	0.9175	0.4857
Human v.s. ROUGE-L	0.9352	0.6571
A <sub>1</sub> v.s. A <sub>2</sub>	0.9874	0.9429

Table 8: Correlations between human judgements and metrics for error descriptions.

less for pre-explaining models to correct the source.

## 6.4 Comparison with Pipeline

We compare the results of multi-task models and GEC-GEE pipeline with COTE in Table 5. It indicates that the pipeline can improve both correction and explanation performance, highlighting the challenges of learning a multi-task model for EXGEC.

## 7 Analysis

### 7.1 Ablation Results

We conduct ablation studies on Qwen1.5-7B-chat to provide in-depth insights into post-explaining models. We also study the effect of model sizes in Appendix B.2 and provide a case study for different LLMs in Appendix B.3.

**Effect of COTE.** The impact of COTE introduced in Section 5.2 is examined in this section. We provide the post-explaining model with ground truth target texts, which allows us to focus on the explanation performance. The results presented in Table 7 reveal a huge performance drop if we do not leverage COTE, especially the hit and miss rates. This demonstrates the effectiveness of COTE.

**Effect of token-wise decoding strategies.** By default, we employ beam search decoding for cor-

rections and top-p decoding for explanations. In this section, we explore the reverse setting, and the results are reported in Table 6. When switching from beam search to top-p for correction, we observe a huge performance drop in precision and F<sub>0.5</sub> and increase in recall, which means top-p encourages LLMs to over-correct (Cao et al., 2023). On the other hand, leveraging beam search improves explanation performance, suggesting the potential benefits of a greedy decoding algorithm for the task. However, we notice that beam search also increases the miss rate. We speculate that beam search may discard some low-likelihood explanations.

### 7.2 Human Evaluation for Error Descriptions

Despite the efficiency of automatic metrics in evaluating error descriptions, their accuracy remains to be confirmed. Therefore, this section attempts to demonstrate the suitability of different metrics by comparing their corrections with human judgements. We report the correlations between two human annotators and the ones between average human ratings and metric scores in Table 8. We observe the inter-annotator correlations are close to 1, meaning it is relatively easy to determine the quality of error descriptions for human. Among various metrics, ROUGE-1 achieve the highest correlations, followed by METEOR. All the introduced metrics show moderate or high correlations, indicating that it is advisable to employ them as proxies for human evaluation. We provide the detailed annotation guidance and rating rules in Appendix B.4.

## 8 Conclusion

We propose and formulate the task of EXGEC, overcoming the limitation of previous studies that fail to establish the interaction of both correction and explanation tasks. To develop the task, we propose the EXCGEC benchmark, based on which we develop baseline models in multiple settings. Extensive experiments and analyses reveal several challenges of the task, and we hope this paper can serve as a starting point for future exploration.



## 555 Limitations

556 **Inferior performance of multi-task models.** In  
557 our experiments, we observe the pipeline solution  
558 outperform the multi-task solutions, regardless of  
559 correction or explanation tasks. This suggests that  
560 the multi-task models struggle to reap positive  
561 benefits from the interaction of both tasks. We  
562 leave the exploration of effective multi-task learn-  
563 ing EXGEC models to the future work.

564 **Limitations of synthesizing datasets.** LLM-  
565 augmented datasets may include some unintended  
566 biases towards or inaccuracies, resulting in skewed  
567 or unfair outcomes in applications. Second, it is  
568 necessary to manually filter out invalid data in order  
569 to ensure the quality of datasets. But it is indeed a  
570 advisable method to construct datasets using LLMs,  
571 considering its efficiency.

572 **Adaptation to other languages.** The general de-  
573 sign of our proposed edit-wise explanations can  
574 be easily adapted to other languages. However,  
575 the detailed design may not be suitable to other  
576 languages. For example, the two-tier hierarchical  
577 taxonomy of error types is tailored for Chinese.

## 578 Ethics Statement

579 Our proposed benchmark is built upon existing  
580 datasets, backbones and metrics, all of which are  
581 publicly available. We have cited the corresponding  
582 authors or projects of them, and confirm that they  
583 are consistent with their intended use.

584 Additionally, we conduct human evaluation ex-  
585 periments to ensure the quality of the dataset and  
586 find out the correlations between metrics and hu-  
587 man judgements. To achieve this, we hire 12 native  
588 speakers, all of whom are graduated students. Each  
589 annotator could complete the entire annotation pro-  
590 cess within approximately 6~8 working hours. All  
591 annotators were paid for their work, with an aver-  
592 age salary of approximately \$5 per hour.

## 593 References

594 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama  
595 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
596 Diogo Almeida, Janko Altenschmidt, Sam Altman,  
597 Shyamal Anadkat, et al. 2023. Gpt-4 technical report.  
598 *arXiv preprint arXiv:2303.08774*.

599 Gloria Ashiya Katuka, Alexander Gain, and Yen-Yun  
600 Yu. 2024. Investigating automatic scoring and feed-  
601 back using large language models. *arXiv e-prints*,  
602 pages arXiv-2405.

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,  
Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei  
Huang, et al. 2023. Qwen technical report. *arXiv  
preprint arXiv:2309.16609*. 603  
604  
605  
606

Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen,  
Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong,  
Qiushi Du, Zhe Fu, et al. 2024. Deepseek llm: Scal-  
ing open-source language models with longtermism.  
*arXiv preprint arXiv:2401.02954*. 607  
608  
609  
610  
611

Steven Bird and Edward Loper. 2004. **NLTK: The natu-  
ral language toolkit**. In *Proceedings of the ACL In-  
teractive Poster and Demonstration Sessions*, pages  
214–217, Barcelona, Spain. Association for Compu-  
tational Linguistics. 612  
613  
614  
615  
616

Christopher Bryant, Mariano Felice, Øistein E Ander-  
sen, and Ted Briscoe. 2019a. The bea-2019 shared  
task on grammatical error correction. In *Proceedings  
of the fourteenth workshop on innovative use of NLP  
for building educational applications*, pages 52–75. 617  
618  
619  
620  
621

Christopher Bryant, Mariano Felice, Øistein E. Ander-  
sen, and Ted Briscoe. 2019b. **The BEA-2019 shared  
task on grammatical error correction**. In *Proceedings  
of the Fourteenth Workshop on Innovative Use of NLP  
for Building Educational Applications*, pages 52–75,  
Florence, Italy. Association for Computational Lin-  
guistics. 622  
623  
624  
625  
626  
627  
628

Christopher Bryant, Mariano Felice, and Ted Briscoe.  
2017. **Automatic annotation and evaluation of error  
types for grammatical error correction**. In *Proceed-  
ings of the 55th Annual Meeting of the Association for  
Computational Linguistics (Volume 1: Long Papers)*,  
pages 793–805, Vancouver, Canada. Association for  
Computational Linguistics. 629  
630  
631  
632  
633  
634  
635

Christopher Bryant, Zheng Yuan, Muhammad Reza  
Qorib, Hannan Cao, Hwee Tou Ng, and Ted Briscoe.  
2023. Grammatical error correction: A survey of  
the state of the art. *Computational Linguistics*,  
49(3):643–701. 636  
637  
638  
639  
640

Oana-Maria Camburu, Tim Rocktäschel, Thomas  
Lukasiewicz, and Phil Blunsom. 2018. e-snli: Natu-  
ral language inference with natural language expla-  
nations. *Advances in Neural Information Processing  
Systems*, 31. 641  
642  
643  
644  
645

Qi Cao, Takeshi Kojima, Yutaka Matsuo, and Yusuke  
Iwasawa. 2023. Unnatural error correction: Gpt-4  
can almost perfectly handle unnatural scrambled text.  
In *Proceedings of the 2023 Conference on Empiri-  
cal Methods in Natural Language Processing*, pages  
8898–8913. 646  
647  
648  
649  
650  
651

Bosheng Ding, Chengwei Qin, Ruochen Zhao, Tianze  
Luo, Xinze Li, Guizhen Chen, Wenhan Xia, Jun-  
jie Hu, Anh Tuan Luu, and Shafiq Joty. 2024.  
Data augmentation using llms: Data perspectives,  
learning paradigms and challenges. *arXiv preprint  
arXiv:2403.02990*. 652  
653  
654  
655  
656  
657





874	• 标点丢失 ( <b>Punctuation Missing</b> ). Punctu-	words in Chinese, such as person names, or-	921
875	ation missing mainly refers to the omission	ganization names, place names, and all other	922
876	of punctuation that should have existed in the	entities identified by terminologies. These	923
877	middle and end of a sentence. For the explana-	words are also very prone to spelling errors.	924
878	tion of punctuation missing errors, first point		
879	out the evidence words and missing punctua-	<b>Word-level Error.</b> Word-level errors often refer	925
880	tion symbols, and then explain the role of the	to misuse of individual words or idioms in a sen-	926
881	added punctuation in this context.	tence, but the syntactic structure of the sentence is	927
		correct. This type of error belongs to the most com-	928
882	• 标点误用 ( <b>Punctuation Misuse</b> ). Misuse of	mon category in Chinese text errors and can usually	929
883	punctuation is very common in daily Chinese	be subdivided into the following three types:	930
884	writing. For cases of punctuation misuse, first		
885	briefly explain the roles of misused punctua-	• 词语冗余 ( <b>Word Redundancy</b> ). The simul-	931
886	tion and correct punctuation, and then explain	taneous appearance of words with the same or	932
887	sufficient reasons for correction.	similar meanings in a sentence can cause se-	933
		semantic repetition and sentence redundancy,	934
888	<b>Spelling-level Error.</b> Spelling-level errors refer	which is known as word redundancy. Re-	935
889	to people who, due to carelessness or lack of knowl-	peated words often appear adjacent to each	936
890	edge, write incorrect characters or words during the	other, so it is important to pay attention to	937
891	writing process. The type is so common that Chi-	whether the meanings of adjacent words are	938
892	nese Spelling Check (CSC), as a standard NLP task	exactly the same. If they are the same, it may	939
893	specialized in spelling-level errors, attract the at-	lead to the problem of word redundancy.	940
894	tention from many researchers. Inspired by these		
895	studies, we categorize spelling-level errors further	• 词语丢失 ( <b>Word Missing</b> ). In modern Chi-	941
896	into 4 sub-classes.	nese, sentences generally have six major com-	942
		ponents, namely subject, predicate, object, at-	943
897	• 字音混淆错误 ( <b>Phonetic Confusion Error</b> ).	tributive, adverbial, complement, etc. A sen-	944
898	Phonetic confusion errors are caused by mis-	tence must express a complete meaning, and	945
899	using the Chinese characters with the same or	its structure must also be complete. The so-	946
900	similar pinyin. The vast majority of Chinese	called complete structure does not mean that a	947
901	Internet users are using pinyin input method,	sentence must have the usual six components,	948
902	so many Chinese spelling-level errors on In-	but rather that the sentence should be com-	949
903	ternet fall in this type.	posed of the necessary components to express	950
		the complete meaning. If the necessary sen-	951
904	• 字形混淆错误 ( <b>Glyph Confusion Error</b> ).	tence components are missing, it will cause	952
905	In addition to pinyin input method, some	the phenomenon of word missing.	953
906	users apply Wubi input method or other glyph-		
907	based input methods. In this case, they are	• 词语误用 ( <b>Word Misuse</b> ). Word Misuse in-	954
908	prone to spelling errors due to confusion of	dicates improper use of words in the text. The	955
909	fonts or strokes.	main cause of this error is the author's insuffi-	956
		cient understanding of the meaning and part	957
910	• 词内部字符异位错误 ( <b>Internal Character</b>	of speech of a certain word.	958
911	<b>Misplacement Error</b> ). Internal character mis-		
912	placement error refers to expressing a multi-	<b>Sentence-level Error.</b> This type mainly involves	959
913	character word in disorder of characters. The	sentence-level issues, not just individual words or	960
914	type seldom happens for native speakers, but	characters. Sentence-level errors are often caused	961
915	sometimes in texts written by L2-speakers.	by violating common syntactic structures, or not	962
916	For example, the spelling-level error “共公”	following objective reasoning.	963
917	falls in this type and should be corrected to		
918	“公共”.	• 词序不当 ( <b>Improper Word Order</b> ). Proper	964
		word order is essential to express exact mean-	965
919	• 命名实体拼写错误 ( <b>Named Entity Mis-</b>	ing in Chinese. Writing texts without accurate	966
920	<b>spelling</b> ). There are numerous named entity	word order results in the type of improper	967
		word order. If a sentence is not combined	968

according to the intended meaning, it may lead to confusion in the sentence structure, resulting in an imbalance in the relationship between sentence components and affecting the expression of sentence meaning.

- **逻辑不通 (Illogicality)**. Illogicality refers to a sentence that conforms to grammatical norms but does not conform to logical reasoning. Illogicality can be caused by many reasons such as improper logical order, causal confusion, and reversal of subject and object.
- **句式杂糅 (Run-on Sentence)**. Run-on Sentence in Chinese usually refers to the use of two formats or sentences with similar or identical meanings in one sentence. People originally used one format when writing sentences, but due to interference from other factors such as sentence content, they may unconsciously switch to another format, resulting in a mixture of the two formats.

**Other Special Error.** Besides the above grammatical error types, other several types can not easily fit in the mentioned major types. So we classify them to other special errors.

- **照应错误 (Inconsistency Error)**. Inconsistency errors are ones involved in the mistaken referential relationship between two words, and explaining this grammatical error requires knowledge of the referential relationship between each word.
- **歧义错误 (Ambiguity Error)**. Ambiguity errors happen when a word or a sentence can be understood as having multiple meanings.
- **语气不协调 (Inconsistent Tone)**. Inconsistent tone refers to the inconsistency of tone between the preceding and following sentences.

Additionally, we define the grammatical error type *Other* as ones that do not fit in any of the above error types. These errors are usually involved in rather significant modification and sometimes change the original semantics.

## A.2 Examples of Error Types

We list the examples of error types in Figure 4, 5, 6.

## A.3 Prompt of Generating Explanations

The prompt we use to generate explanations is shown in Figure 8. We also provide an English version in Figure 9.

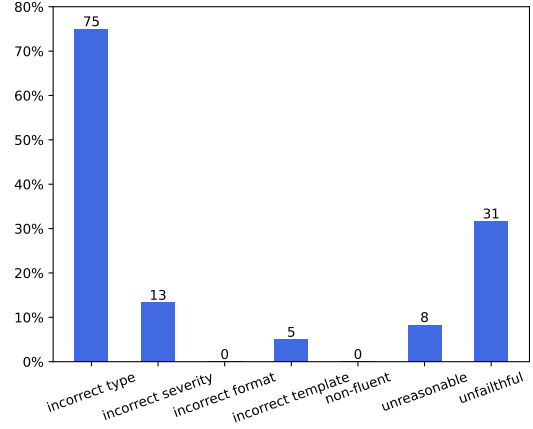


Figure 3: Distribution of 7 kinds of LLM errors.

Configuration	Value
<b>Fine-tuning</b>	
Devices	2 Tesla A100 GPU (80GB)
Epochs	5
Finetuning type	Lora
Train batch size per GPU	2
Eval batch size per GPU	1
Gradient accumulation steps	16
Optimizer	AdamW
Learning rate	$(\beta_1 = 0.9, \beta_2 = 0.98, \epsilon = 1 \times 10^{-6})$
Learning rate schedule	$5 \times 10^{-5}$
Warmup steps	cosine decay
Eval steps	20
Cutoff length	200
Preprocessing workers number	1024
Numerical precision	16
Weight decay	fp16
	0.05
<b>Inference</b>	
Beam size	5
Top-p	0.8
Max new tokens	2048
Temperature	0.7

Table 9: Hyper-parameter values used in our experiments.

## A.4 Detailed Description of LLM Failure Modes

We categorize the failure modes in our case into seven major reasons: incorrect type, incorrect severity, incorrect format, incorrect template, non-fluency, unreasonability, and unfaithfulness. One expert annotator is asked to classified the sampled 100 invalid explanations, where an explanation may be categorized into multiple failure modes. The annotation results, illustrated in Figure 3, reveal that GPT-4 tend to mis-classify grammatical errors and providing unfaithful error descriptions. On the other hand, GPT-4 is capable to a large extend to offer well formed, fluent, and reasonable explanations, demonstrating the effectiveness of LLM annotation on this task.

The definitions of seven failure modes of expla-

Model	Correction $\uparrow$			Explanation						
	CLEME (P / R / F <sub>0.5</sub> )	ChERRANT (P / R / F <sub>0.5</sub> )	Hit $\uparrow$	Miss $\downarrow$	Acc $\uparrow$	F1 $\uparrow$	MAE $\downarrow$	BLEU $\uparrow$	METEOR $\uparrow$	ROUGE- (1 / 2 / L)
Qwen1.5-1.8B-chat	21.11 / 19.28 / 20.72	28.91 / 15.70 / 24.74	59.94	65.14	55.80	23.27	0.89	10.19	34.35	48.66 / 22.70 / 34.23
Qwen1.5-4B-chat	22.49 / 20.84 / 22.14	30.57 / 16.85 / 26.29	62.91	<b>62.70</b>	57.16	25.31	0.85	11.61	35.91	46.83 / 19.59 / 30.86
Qwen1.5-7B-chat	<b>28.31 / 21.21 / 26.54</b>	<b>36.74 / 17.26 / 29.98</b>	<b>68.94</b>	64.83	<b>61.98</b>	<b>29.62</b>	<b>0.75</b>	<b>15.49</b>	<b>38.88</b>	<b>50.32 / 24.25 / 35.24</b>

Table 10: Comparison of post-explaining models with various model sizes.

nations are as follows:

- **Incorrect type:** the error type is incorrect.
- **Incorrect format:** the evidence content and the correction content are not highlighted by special markers [ ] or { }.
- **Incorrect template:** the error description does not follow the syllogism form of deductive reasoning.
- **Non-fluency:** the error description is non-fluent or unreadable.
- **Unreasonability:** the error description contains obvious mistakes about linguistics, thus making it unacceptable for human.
- **Unfaithfulness:** the error description is not targeted to the given edit.

## B Experimental Details and Extra Results

### B.1 Implementation Details.

We train all models for 5 epochs and select the best model validated on EXCGEC-dev and report its performance on EXCGEC-test. The detailed training hyperparameter values of the all models in our experiments are shown in Table 9.

### B.2 Effect of Model Sizes

Table 10 indicates the varying performance across model sizes ranging from 1.8B to 7B. We observe consistent performance enhancement with increasing model sizes.

### B.3 Case Study

We provide a case study in Table 7.

### B.4 Details of Human Rating

Specifically, we hire 2 native Chinese speakers to rate the explanations generated by 6 post-explaining models in Table 3 conditioned on ground truth target texts. The rating scores range from 0 to 100, and each annotator concurrently rate 6 explanations for each sample. We randomly

select 100 samples for annotation. We provide annotators with general scoring suggestions:

- **100 points:** Explain and describe fluently (fluency), introduce relevant semantic knowledge to enhance persuasiveness (rationality), and explain that it is aimed at the current editor (loyalty). All aspects are impeccable, and there is almost no better explanation or description than this.
- **80~100 points:** The explanation and description are expressed fluently, satisfy fidelity, and have a certain degree of rationality, but there are certain degrees of flaws.
- **60~80 points:** The explanation and description are expressed fluently, but the fidelity or rationality is not good enough, but it is somewhat helpful for correcting the grammar error in understanding.
- **30~60 points:** The explanation and description are expressed fluently, but the rationality is poor, and it is not very helpful for correcting the grammar error in understanding.
- **0~30 points:** The explanation and description are expressed fluently, but the fidelity is poor, and the object of explanation is not the current editor. There is no help in correcting the grammar error for understanding.
- **0~30 points:** The explanation and description are vague and cannot be understood. There is no help in correcting the grammar error for understanding.

```

1 标点级别错误
# 标点冗余
{
  "input": "所以一些人说：‘读书一点用处都没有。’"
  "output": "所以一些人说：‘读书一点用处都没有。’",
  "explanations": [
    {
      "error_severity": 1,
      "error_type": "标点冗余",
      "error_description": "【：】直接用于【说】等总说性或提示性词语后面，提起下文，没有必要在【：】前插入逗号。应删去【说】之后的冒号。"
    }
  ]
},
# 标点丢失
{
  "input": "人为了生存不管是干净的空气还是污染的空气都要呼吸。"
  "output": "人为了生存，不管是干净的空气还是污染的空气，都要呼吸。",
  "explanations": [
    {
      "error_severity": 1,
      "error_type": "标点丢失",
      "error_description": "【人为了生存】和【不管……】是两个分句，复句内各分句之间应使用逗号表示停顿。应在【为了生存】后添加逗号。"
    }
  ]
}
# 标点误用
{
  "input": "那我们一定要参加这个活动吗。"
  "output": "那我们一定要参加这个活动吗？",
  "explanations": [
    {
      "error_severity": 1,
      "error_type": "标点误用",
      "error_description": "‘句号’主要表示句子的陈述语气，而‘问号’主要表示句子的疑问语气。【吗】意味着该句是一个疑问句，故应【吗】后的句号改为问号。"
    }
  ]
}
2 拼写级别错误
# 字音混淆错误
{
  "input": "我们舒舍有四个人。"
  "output": "我们宿舍有四个人。",
  "explanations": [
    {
      "error_severity": 1,
      "error_type": "字音混淆错误",
      "error_description": "【宿舍}指学校或用人单位等提供给学生和职工的房屋，对应句子中的【有四个人】。{宿舍}和【舒舍}发音相近，导致了此处的拼写错误。应将【舒舍}改为{宿舍}。"
    }
  ]
}
# 字形混淆错误
{
  "input": "这座关溢非常雄伟。"
  "output": "这座关隘非常雄伟。",
  "explanations": [
    {
      "error_severity": 1,
      "error_type": "字形混淆错误",
      "error_description": "【关隘}指险要的关口，在交通要道设立的防务设施，又称关卡。{关隘}和【溢}字形相近，导致了此处的拼写错误。应将【关溢}改为{关隘}。"
    }
  ]
}
# 词内部字符异位
{
  "input": "我非常爱吃阴冬功。"
  "output": "我非常爱吃冬阴功。",
  "explanations": [
    {
      "error_severity": 2,
      "error_type": "词内部字符异位",
      "error_description": "【冬阴功}是泰国和老挝的一道富有特色的酸辣口味汤品，书写者错误地将该词写成{阴冬功}。应将【阴冬功}改为{冬阴功}。"
    }
  ]
}

```

Figure 4: Examples of error types.

```

# 命名实体拼写错误
{
  "input": "我们都是海南詹州人。"
  "output": "我们都是海南儋州人。",
  "explanations": [
    {
      "error_severity": 2,
      "error_type": "命名实体拼写错误",
      "error_description": "中国【海南】不存在【詹州】这一地名，但存在字形相近的{儋州}。【詹】与{儋}字形相近，导致了此处的拼写错误。应将【詹州】改为{儋州}。"
    }
  ]
},

3 词语级别错误
# 词语冗余
{
  "input": "终于看到了大熊猫，儿子显得特别兴奋极了。"
  "output": "终于看到了大熊猫，儿子显得特别兴奋。",
  "explanations": [
    {
      "error_severity": 3,
      "error_type": "词语冗余",
      "error_description": "【特别】与【极了】都是修饰【兴奋】的程度副词，两者重复。应删去【特别】与【极了】其中一个。"
    }
  ]
},

# 词语丢失
{
  "input": "最终经过他的不懈努力，成为了一个地位很高的长官。"
  "output": "最终经过不懈努力，他成为了一个地位很高的长官。",
  "explanations": [
    {
      "error_severity": 3,
      "error_type": "词语丢失",
      "error_description": "状语从句【经过不懈努力】和主句的谓语【成为】具有共同主语【他】，此处将【他】放在【经过】的后面导致句子缺失主语。可以把【他】放在【成为】之前，也可以把【他】提到【经过】的前面，充当状语从句和主句的共同主语。"
    }
  ]
},

# 词语误用
{
  "input": "这样一个年过八旬的老奶奶在她即将逝去的生命中仍然绽放着希望的光辉。"
  "output": "这样一个年过八旬的老奶奶在她即将逝去的生命中仍然散发着希望的光辉。",
  "explanations": [
    {
      "error_severity": 3,
      "error_type": "词语误用",
      "error_description": "谓语句【绽放】和宾语【光辉】搭配不当，【绽放】一般用于形容花开时由花蕾花瓣紧闭展开的样子。应将【绽放】改为{散发}。"
    }
  ]
},

4 句法级别错误
# 词序不当
{
  "input": "改革开放后，中国的经济增长速度加快明显起来。"
  "output": "改革开放后，中国的经济增长速度明显加快起来。",
  "explanations": [
    {
      "error_severity": 4,
      "error_type": "词序不当",
      "error_description": "状语【明显】用于修饰谓语句【加快】，一般放在谓语句之前。应将【明显】提到【加快】前面。""
    }
  ]
},

# 逻辑不通
{
  "input": "我们要注意多多提高总结自己。"
  "output": "我们要注意多多总结提高自己。",
  "explanations": [
    {
      "error_severity": 3,
      "error_type": "逻辑不通",
      "error_description": "按照动作的发生顺序，应该先【总结】，再【提高】。应将逻辑顺序不当的【提高总结】改为{总结提高}。"
    }
  ]
},

```

Figure 5: Examples of error types.



```

# 句式杂糅
{
  "input": "形成沼泽的原因是水体沼泽化的结果。"
  "output": "形成沼泽是水体沼泽化的结果。",
  "explanations": [
    {
      "error_severity": 4,
      "error_type": "句式杂糅",
      "error_description": "【原因是.....】和【是.....的结果】都是表示原因的句式，将两个意思相同或相近的格式放在一个句子里使用导致格式杂糅。应删去其中一个。"
    }
  ]
}

5 其它特殊错误
# 照应错误
{
  "input": "蜗牛感到自己很悲哀，不只是因为它生命的长短。"
  "output": "蜗牛感到自己很悲哀，不只是因为它生命很短。",
  "explanations": [
    {
      "error_severity": 4,
      "error_type": "照应错误",
      "error_description": "生命的【短】可以照应【悲哀】，而【长】则不能照应。应将【生命的长短】改为【生命很短】。"
    }
  ]
},

# 歧义错误
{
  "input": "他没有听从导师的建议，对论文的观点加以修改，影响了论文的水平。"
  "output": "他没有根据导师的建议对论文的观点加以修改，影响了论文的水平。",
  "explanations": [
    {
      "error_severity": 4,
      "error_type": "歧义错误",
      "error_description": "句子有歧义，一方面可理解为导师建议他对论文的观点加以修改，但是他没有听从；另一方面可理解为导师认为论文不需要修改。本句可改为{他没有根据导师的建议对论文的观点加以修改，影响了论文的水平}。"
    }
  ]
}

# 语气不协调
{
  "input": "生活在这个万物滋长的社会中，我不禁慨叹社会的安详宁静仅仅是因为制度的限制吗？"
  "output": "生活在这个万物滋长的社会中，我不禁疑惑社会的安详宁静仅仅是因为制度的限制吗？",
  "explanations": [
    {
      "error_severity": 3,
      "error_type": "语气不协调",
      "error_description": "动词【慨叹】指感慨叹息，一般不用于疑问句。应将【慨叹】改为常用于疑问句的【疑惑】。"
    }
  ]
}

```

Figure 6: Examples of error types.

```

{
  "input": "我不想反叛宋老师的期待。",
  "output": "我不想辜负宋老师的期待。",

  # qwen1.5-7B-Chat
  "explanations": [
    {
      "error_type": "词语误用",
      "error_severity": 3,
      "error_description": "【反叛】通常指违反现有规则、制度或法律的行为，而根据上下文，'辜负期待'更符合句意，即没有达到或辜负期望。应将【反叛】改为{辜负}，以表达不辜负老师期望的意思。"
    }
  ],

  # qwen1.5-7B-Base
  "explanations": [
    {
      "error_type": "词语误用",
      "error_severity": 2,
      "error_description": "动词【反叛】一般指公开或秘密地反对政府或现有社会秩序的行为，而在这里应当使用表示辜负或使失望的动词。因此，应将【反叛】替换为{辜负}，以更准确地表达不想让宋老师失望的意思。"
    }
  ],

  # Llama3-Instruct
  "explanations": [
    {
      "error_type": "词语误用",
      "error_severity": 3,
      "error_description": "【反叛】通常指对某个制度、组织或权威的公开抵制或反对行为，而在这里应表达的是不想让宋老师失望的意思。正确的表达是【辜负】，表示不想让别人失望或不值得信赖。应将【反叛】改为{辜负}。"
    }
  ],

  # Llama3-Base
  "explanations": [
    {
      "error_type": "词语误用",
      "error_severity": 2,
      "error_description": "【反叛】一词通常指对某人或某个集体的背叛或叛变，而在此句中，作者表达的是对老师的期待的不负责任的含义，应该使用{辜负}。【辜负】一词含义是指不尽承诺或期望的义务，适用于此处语境。应将【反叛】替换为{辜负}。"
    }
  ],

  # DeepSeek-Chat
  "explanations": [
    {
      "error_type": "词语误用",
      "error_severity": 3,
      "error_description": "在此上下文中，【反叛】一词通常用于表达与权威或传统对立的行为，而【辜负】则用于表达未能达到他人的期望或未能履行承诺。根据句意，应表达的是未能满足老师期望的情感，所以应使用{辜负}。"
    }
  ],

  # DeepSeek-Base
  "explanations": [
    {
      "error_type": "词语误用",
      "error_severity": 2,
      "error_description": "在此上下文中，【反叛】通常用于表达对权威或传统观念的反抗或挑战，而【辜负】则表示没有达到或未能实现他人的期望。根据句意，{反叛}一词并不恰当，因为这里表达的是对宋老师期待的尊重和感激。应将【反叛】改为{辜负}，以更准确地表达对宋老师期待的尊重和感激之情。"
    }
  ],

  # GPT-4
  "explanations": [
    {
      "error_type": "词语误用",
      "error_severity": 3,
      "error_description": "动词【反叛】通常用来描述违背或背叛权威、组织或原则等的行为，而在这个语境中，作者意图表达的是不想让宋老师失望。因此，应用表达失望或希望落空的词汇{辜负}替换不当用词【反叛】，使句子意义正确传达。"
    }
  ]
}

```

Figure 7: A case study of all the LLMs involved in our experiments.

你是一个优秀的语法纠错解释模型，你需要在仔细对比错误句和纠正句的前提下，对错误句中的标点错误、拼写错误、词语错误和句法错误等提供流畅、合理且忠实的解释，解释包括语法错误类型、错误程度和错误描述。流畅性要求解释本身没有语法错误且表达流畅；合理性要求对语法错误的解释是能被人们接受的；忠实性要求对句子中所有语法错误都有对应解释，且解释能对应正确句的纠正方式。

每个语法错误由一个编辑改动 (edit) 来表示，为了提升解释的合理性和忠实性，你必须遵守以下原则：

- 1) 必须对每个给定的语法错误进行解释，禁止私自修改编辑中的错误内容 (src\_content) 和纠正内容 (tgt\_content)。
- 2) 必须对每个语法错误分别给出相应的错误类型 (error\_type)、错误程度 (error\_severity) 和错误描述 (error\_description)。
- 3) 如果一处编辑改动存在多个语法错误，选择优先级最高的语法错误进行解释，优先级顺序：句法级别错误>词语级别错误>拼写级别错误>标点级别错误。
- 4) 错误类型禁止自主捏造，只能来自下列错误类型：

- 标点冗余、标点丢失、标点误用
- 字音混淆错误、字形混淆错误、词内部字符异位错误、命名实体拼写错误
- 词语冗余、词语丢失、词语误用
- 词序不当、逻辑不通、句式杂糅
- 照应错误、歧义错误、语气不协调
- 其他错误

中的一个。语法错误类型将在下文给出定义和示例。当无法确定具体的错误类型时，统一分类为“其他错误”。

5) 错误程度的打分范围为1-5分，下面是每种分数在语法、语义层面上的详细描述和例句：

- 1分（无关紧要的错误）：可能是一些常规的打字错误或者一些影响很小的误用词语。例如：“他擅长数学和英语”应为“他擅长数学和英文”。
- 2分（轻度语法错误）：可能引起表达混淆，但并不会影响完整的理解。例如：“我喜欢狗和猫播放电子游戏”应为“我喜欢玩电子游戏，还喜欢狗和猫”。
- 3分（中度语法错误）：可能会导致句子部分不流畅，使读者需要重新阅读以理解含义。例如：“我走家去了”应为“我走去家了”。
- 4分（严重语法错误）：不仅会对理解产生困扰，还可能完全改变句子的意思。例如：“我想借用你的手机扮演职业摄影师”应为“我想借用你的手机拍摄一些专业的照片”。
- 5分（极度严重的语法错误）：可能导致句子无法理解。例如：“他举妈妈，我去购物车”应为“他举着妈妈的购物车，我就去了”。

6) 错误描述必须符合演绎推理的三段论形式：[大前提：语法规则和相关知识] [小前提：当前文本的错误原因] [阐述如何纠正]

7) 错误描述需要提供充分且全面的纠正证据词，并使用以下符号强调错误描述中的证据词和纠正方式：

- 证据词必须是出现在错误句中的文本段，并且前后使用【】包围。
- 纠正方式必须是出现在纠正句中的文本段，并且前后使用{}包围。

注意：下列大多数示例仅包含一个语法错误，但是正式输入数据通常包含多个语法错误，你必须对每个语法错误都分别给出相应的解释。输出必须严格符合 json 格式。

1 标点级别错误，即涉及标点符号的语法错误。

1.1 标点冗余：指在不必要的地方插入了标点。对于标点冗余错误，首先阐述所涉及标点符号的作用，然后解释标点冗余的原因。

解释标点冗余的建议模板为：[解释冗余标点和相关证据词的基本用法] [解释标点冗余的原因] 应删去[冗余标点]

标点冗余输入示例：

```
{
  "error_sentence": "所以一些人说：，"读书一点用处都没有。",
  "correct_sentence": "所以一些人说："读书一点用处都没有。",
  "edit": [
    {
      "src_interval": [6, 7],
      "tgt_interval": [6, 6],
      "src_content": "，",
      "tgt_content": ""
    }
  ]
}
```

标点冗余输出示例：

```
{
  "edits": [
    {
      "src_interval": [6, 7],
      "tgt_interval": [6, 6],
      "src_content": "，",
      "tgt_content": "",
      "error_severity": 1,
      "error_type": "标点冗余",
      "error_description": "【：】直接用于【说】等总说性或提示性词语后面，提起下文，没有必要在【：】前插入逗号。应删去【说】之后的逗号。"
    }
  ]
}
```

1.2 标点丢失：主要指的是在句中、句末漏写了本应存在的标点。对于标点丢失错误的解释，首先要点明证据词和缺失的标点符号，然后阐述所加标点在此处的作用。

解释标点丢失的建议模板为：[解释丢失标点和相关证据词的基本用法] [解释标点丢失的原因] 应在[证据词]前/后添加[丢失标点]

.....

Figure 8: The prompt used for explanation generation. For each error type, We provide the definition, a suggested template of error description, and a demonstration for GPT-4.

You are an excellent grammar error correction explanation model. Your task is to provide fluent, reasonable, and faithful explanations for punctuation errors, spelling errors, word errors, and syntactic errors in erroneous sentences by carefully comparing the erroneous sentences with the corrected sentences. The explanations should include the type of grammatical error, the severity of the error, and a description of the error. Fluency requires that the explanation itself has no grammatical errors and is expressed fluently; reasonableness requires that the explanation of the grammatical error is acceptable to people; faithfulness requires that all grammatical errors in the sentence have corresponding explanations, and the explanations should correspond to the correction methods of the correct sentence.

Each grammatical error is represented by an edit. To improve the reasonableness and faithfulness of the explanations, you must follow these principles:

1. Each given grammatical error must be explained, and the error content and correction content in the edits must not be modified.
2. Each grammatical error must be given a corresponding error type, error severity, and error description.
3. If an edit contains multiple grammatical errors, choose the grammatical error with the highest priority to explain. The priority order is: syntactic-level errors > word-level errors > spelling-level errors > punctuation-level errors.
4. Error types must not be fabricated; they can only come from the following error types:
  - Punctuation Redundancy, Punctuation Missing, Punctuation Misuse
  - Phonetic Confusion Error, Glyph Confusion Error, Internal Character Misplacement Error, Named Entity Misspelling
  - Word Redundancy, Word Missing, Word Misuse
  - Improper Word Order, Illogicality, Run-on Sentence
  - Inconsistency Error, Ambiguity Error, Inconsistent Tone
  - Other errors

The definitions and examples of grammatical error types will be provided later. When it is impossible to determine the specific error type, classify it as "Other errors".
5. The scoring range for error severity is 1-5 points. Here is a detailed description and examples of each score at the grammatical and semantic levels:
  - 1 point (trivial error): It may be some routine typing errors or minor word misuse that has little impact. Example: "他擅长数学和英语" should be "他擅长数学和英文".
  - 2 points (minor grammatical error): It may cause confusion in expression but does not affect the overall understanding. Example: "我喜欢狗和猫播放电子游戏" should be "我喜欢玩电子游戏, 还喜欢狗和猫".
  - 3 points (moderate grammatical error): It may cause parts of the sentence to be incoherent, requiring the reader to reread to understand the meaning. Example: "我走家去了" should be "我走去家了".
  - 4 points (serious grammatical error): It not only causes confusion in understanding but may also completely change the meaning of the sentence. Example: "我想借用你的手机扮演职业摄影师" should be "我想借用你的手机拍摄一些专业的照片".
  - 5 points (extremely serious grammatical error): It may make the sentence incomprehensible. Example: "他举妈妈, 我去购物车" should be "他举着妈妈的购物车, 我就去了".
6. The error description must follow the deductive reasoning form of a syllogism: [Major premise: semantic rules and related knowledge] [Minor premise: the reason for the current text error] [Explain how to correct it].
7. The error description needs to provide sufficient and comprehensive correction evidence words and use the following symbols to emphasize the evidence words and correction methods:
  - Evidence words must be text segments appearing in the erroneous sentence, surrounded by **【】**.
  - Correction methods must be text segments appearing in the corrected sentence, surrounded by **{}**.

Note: Most examples below contain only one grammatical error, but formal input data usually contains multiple grammatical errors, and you must provide corresponding explanations for each grammatical error. The output must strictly follow the JSON format.

1. Punctuation-level errors: These involve grammatical errors related to punctuation marks.

1.1 Punctuation Redundancy : Refers to inserting punctuation marks unnecessarily. For redundant punctuation errors, first explain the function of the involved punctuation mark, then explain the reason for the redundant punctuation.

Suggested template for explaining redundant punctuation: [Explain the basic usage of redundant punctuation and related evidence words] [Explain the reason for the redundant punctuation] Delete [redundant punctuation] Example of input with redundant punctuation:

```
{
  "error_sentence": "所以一些人说, : "读书一点用处都没有。",
  "correct_sentence": "所以一些人说: "读书一点用处都没有。",
  "edit": [
    {
      "src_interval": [6, 7],
      "tgt_interval": [6, 6],
      "src_content": ", ",
      "tgt_content": ""
    }
  ]
}
```

Example of output for redundant punctuation:

```
{
  "edits": [
    {
      "src_interval": [6, 7],
      "tgt_interval": [6, 6],
      "src_content": ", ",
      "tgt_content": "",
      "error_severity": 1,
      "error_type": "标点冗余",
      "error_description": "【】直接用于【说】等总说性或提示性词语后面, 提起下文, 没有必要在【:】前插入逗号, 应删去【说】之后的逗号。"
    }
  ]
}
```

1.2 Punctuation Missing : Mainly refers to missing punctuation marks that should be present in the sentence, either within or at the end of the sentence. For explaining missing punctuation errors, first identify the evidence words and the missing punctuation mark, then explain the function of the punctuation mark in that context.

Suggested template for explaining missing punctuation: [Explain the basic usage of the missing punctuation and related evidence words] [Explain the reason for the missing punctuation] Add the missing punctuation before/after [evidence words]

Figure 9: The English prompt used for explanation generation.