

Corrections Meet Explanations: A Unified Framework for Explainable Grammatical Error Correction

Anonymous ACL submission

Abstract

Grammatical Error Correction (GEC) faces the important yet challenging issue of explainability, especially when GEC systems are developed for language learners who often struggle to understand the correction results without reasonable explanations. Extractive evidence words and grammatical error types are two crucial factors of GEC explanations. However, existing work focuses on extracting evidence words and predicting grammatical error types given a source sentence and/or a target sentence as input, ignoring the interaction between explanations and corrections. To bridge the gap, we introduce **EXGEC**, a unified explainable GEC framework that jointly perform explanation and correction tasks in a sequence-to-sequence generation manner, hypothesizing both tasks would benefit each other. Extensive experiments enable us to fully understand and establish the interaction between tasks. Especially, if models are required to jointly predict corrections and explanations, the performance of both tasks improves compared to their respective single-task baselines. Additionally, we observe that **EXPECT**, a recent explainable GEC dataset, contains considerable noise that may confuse model training and evaluation. Therefore, we rebuild **EXPECT** to eliminate the noise, leading to an objective training and evaluation pipeline ¹.

1 Introduction

Writing is a learnt skill that is particularly challenging for second-language (L2) speakers, who often struggle to create grammatical and comprehensible texts (Bryant et al., 2022). To address the problem of ungrammatical writing, GEC systems are designed to identify and correct all grammatical errors in texts. Research in the field of GEC has extended to include multi-language (Rothe

et al., 2021), multi-modality (Fang et al., 2023), document-level (Yuan and Bryant, 2021) and domain adaptation (Zhang et al., 2023).

However, the explainability of GEC is still underdeveloped due to its inherent challenges (Hanawa et al., 2021; Kaneko et al., 2022). Since neural GEC systems are typically complex black-box systems, their inner working mechanisms are opaque (Zhao et al., 2023). The lack of explainability can lead to insufficiency in an educational context, where L2-speakers may struggle to thoroughly grasp the writing skills from GEC systems without understanding why a correction is needed. Equipping corrections with explanations builds appropriate trust by elucidating the linguistic knowledge and reasoning mechanism behind model predictions in an understandable manner, assisting pedagogically end users with elementary language proficiency (Bitchener et al., 2005; Sheen, 2007). Additionally, explainability provides insight to identify unintended biases and risks for researchers and developers, acting as a debugging aid to quickly advance model performance (Ludan et al., 2023).

To help language learners better understand why GEC systems make a certain correction, Fei et al. (2023) introduce **EXPECT**, a large dataset annotated with *evidence words* and *grammatical error types*. Evidence words, which are formally called extractive rationales ², provides specific clues for corrections, helping L2-speakers understand “why to correct”. the error types in **EXPECT** cover 15 pragmatism-based categories (Skehan, 1998; Gui, 2004), facilitating L2-speakers in inferring abstract grammar rules from specific errors in an inductive reasoning manner. However, Fei et al. (2023) focus on explaining GEC given an ungrammatical source and/or a corrected sentence, ignoring the interaction between explanation and correction

¹All the source codes and data will be released after the review anonymity period.

²We use the term “evidence words” throughout the paper except Section 6, following Fei et al. (2023).

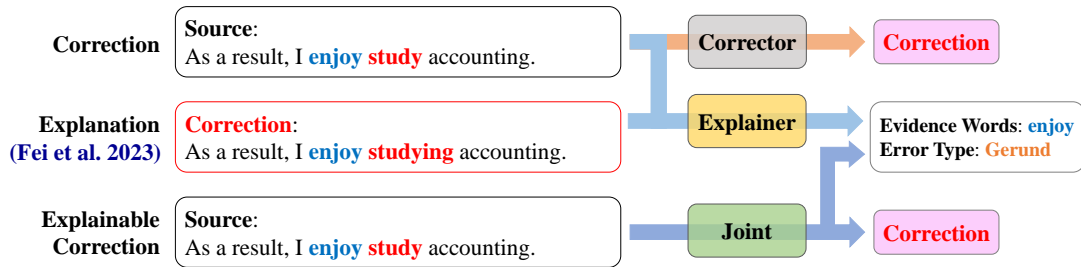


Figure 1: Comparison between correction, explanation (Fei et al., 2023) and our explainable GEC.

078 tasks, as shown in Figure 1. Previous studies have
 079 shown that training models to jointly output task
 080 predictions and explanations can improve the task
 081 performance on vision-language tasks (Majumder
 082 et al., 2022) and diversity downstream NLP tasks,
 083 including text classification (Li et al., 2022a), com-
 084 mon-sense reasoning (Veerubhotla et al., 2023), and
 085 complaint detection (Singh et al., 2023).

086 To establish the interaction between explana-
 087 tion and correction tasks, we propose **EXGEC**
 088 (**EX**plainable **G**rammatical **E**rror **C**orrection), a
 089 unified explainable GEC framework that reframes
 090 the multi-task problem as a sequence-to-sequence
 091 (Seq2Seq) generation task. With pointing mecha-
 092 nism (Vinyals et al., 2015), EXGEC can extract evi-
 093 dence words by directly generating source indexes
 094 of an ungrammatical source sentence in an auto-
 095 regressive manner. EXGEC can jointly correct un-
 096 grammatical sentences, extract evidence words and
 097 classify grammatical errors in a unified architec-
 098 ture. To the best of our knowledge, we first propose
 099 to jointly perform both correction and explanation
 100 tasks. Our findings illustrate that learning correc-
 101 tion and explanation tasks concurrently can benefit
 102 each other. Specifically, pre-explaining models
 103 achieve higher correction performance yet lower
 104 explanation performance than post-explaining mod-
 105 els. However, both models achieve better or compa-
 106 rable correction and explanation performance than
 107 their respective baselines.

108 Additionally, we observe that EXPECT is not a
 109 well-specified dataset for explainable GEC. This
 110 is due to the presence of considerable unidentified
 111 grammatical errors in EXPECT, which hinder the
 112 performance of both tasks. As a result, we rebuild
 113 EXPECT to re-correct the unidentified errors while
 114 ensuring that each sentence contains only a single
 115 unique error, as described by Fei et al. (2023). By
 116 training on rebuilt EXPECT, we significantly im-
 117 prove the performance of both tasks, demonstrating
 118 the effectiveness of our rebuild process.

2 Rebuilt EXPECT Dataset

119 In this paper, we utilize the EXPECT dataset (Fei
 120 et al., 2023). The dataset comprises a total of
 121 20,016 samples that are split into train, dev and
 122 test sets. EXPECT is annotated based on the high-
 123 quality GEC dataset, W&I+LOCNESS (Bryant
 124 et al., 2019), which is designed to represent a much
 125 wider range of English levels and abilities than pre-
 126 vious corpora. To reduce the difficulty of the model
 127 learning and evaluation, EXPECT is constructed
 128 using a special process. Specifically, for a sentence
 129 from W&I+LOCNESS with n grammatical errors,
 130 the authors repeat the sentence n times and keep a
 131 single unique error in each sentence. Considering
 132 the challenges of explainable GEC, it is reasonable
 133 and desirable as it smooths the task by classifying
 134 a grammatical error and extracting evidence
 135 words for a single unique grammatical error each
 136 time, avoiding the confusion caused by multiple
 137 interactive grammatical errors in a sentence.
 138

139 However, we argue that the official EXPECT
 140 dataset is not well-specified. Specifically, for
 141 a sentence with $n(n > 1)$ grammatical errors
 142 from W&I+LOCNESS, the authors correct a single
 143 grammatical error and leave the remaining
 144 $n - 1$ errors unidentified, as shown in Table 1.
 145 These unidentified grammatical errors may confuse
 146 models, making it uncertain which error should
 147 be corrected and explained, and leading to uncer-
 148 tainty in model training and evaluation. To address
 149 the problem, we re-correct the unidentified gram-
 150 matical errors, while leaving the single original
 151 grammatical error unchanged. The entire rebuild-
 152 ing process is automatic since we re-correct all
 153 the unidentified grammatical errors by comparing
 154 sentences from EXPECT and W&I+LOCNESS. We
 155 first retrieve the original parallel samples of
 156 W&I+LOCNESS by using the open-source toolkit
 157 TheFuzz³, and then identify and correct the un-

³<https://github.com/seatgeek/thefuzz>

| | |
|-------------------------------|---|
| W&I+LOCNESS Source | However I sometimes do a skipping to fit myself . |
| W&I+LOCNESS Target | However , I sometimes do skipping to keep myself fit . |
| EXPECT Source | However I sometimes do skipping to keep myself . |
| EXPECT Target | However I sometimes do skipping to keep myself fit . |
| Rebuilt Source | However , I sometimes do skipping to keep myself . |
| Rebuilt Target | However , I sometimes do skipping to keep myself fit . |
| W&I+LOCNESS Source | i have a dog it name 's chente , it is a golden retriver . |
| W&I+LOCNESS Target | I have a dog and its name 's Chente . It is a golden retriever . |
| EXPECT Source | i have a dog its name 's chente , it is a golden retriver . |
| EXPECT Target | i have a dog and its name 's chente , it is a golden retriver . |
| Rebuilt Source | I have a dog its name 's Chente . It is a golden retriever . |
| Rebuilt Target | I have a dog and its name 's Chente . It is a golden retriever . |

Table 1: Examples of our rebuilt EXPECT. We mark grammatical errors in blue and corrections in red.

| | Train | Dev | Test | |
|-----------------|----------------------|--------|--------|--------|
| Official | #Sent. | 15,187 | 2,413 | 2,416 |
| | #Evi. Sent. | 11,261 | 1,426 | 1,444 |
| | Perc. | 74.15% | 59.10% | 59.77% |
| | Avg. Words | 28.68 | 29.06 | 29.23 |
| | Avg. Edits | 1.03 | 1.08 | 1.07 |
| | Avg. EW/Sent. | 2.59 | 3.00 | 3.01 |
| Rebuilt | #Sent. | 15,187 | 2,413 | 2,416 |
| | #Evi. Sent. | 11,261 | 1,425 | 1,443 |
| | Perc. | 74.15% | 59.06% | 59.73% |
| | Avg. Words | 28.52 | 29.53 | 29.72 |
| | Avg. Edits | 1.03 | 1.08 | 1.07 |
| | Avg. EW/Sent. | 2.59 | 3.00 | 3.00 |

Table 2: Statistics of the official and rebuilt EXPECT datasets, including the number of sentences (#Sent.), the average number of words per sentence (Avg. Words), the average number of edits per sentence (Avg. Edits), the number and percentage of sentences with annotated evidence (#Evi. Sent. and Perc.), and the average number of evidence words per sentence (Avg. EW/Sent.).

derlying grammatical errors by leveraging GEC evaluation toolkits ERRANT (Bryant et al., 2017). It is worth noting that the evaluation for the official and rebuilt EXPECT datasets are fairly comparable since the grammatical errors and evidence words are retained during the rebuild process, except for a few extreme cases⁴. Totally, 277 (1.82%), 1,311 (54.33%), and 1,323 (54.76%) sentences in our rebuilt train/dev/test sets differ from their original sentences of official EXPECT. Detailed statistics of both EXPECT datasets are listed in Table 2.

3 Methodology

3.1 Problem Definition

The goal of this work is to perform both correction and explanation tasks jointly in a Seq2Seq-based

⁴One sample from the dev set and one sample from the test set are free from evidence words since their evidence words overlap with the unidentified grammatical errors.

generation approach. Formally, given an ungrammatical source sentence $X = \{x_0, x_1, \dots, x_n\}$, where n is the length of the source sentence, joint models are designed to learn both correction and explanation tasks. The correction task involves transforming the ungrammatical source into a grammatical target $Y = \{y_0, y_1, \dots, y_m\}$, where m is the length of the target. The explanation task consists of two sub-tasks: 1) **classifying** grammatical errors, and 2) **extracting** evidence words. The classification task requires joint models to output a grammatical error type label $c (c \in C)$, where C is the set of 15 candidate grammatical error type classes defined in EXPECT. And the extraction task requires models to extract evidence words $E(X) = \{e_0, e_1, \dots, e_k\} \subset X$ that can provide informative and complete clues for corrections.

3.2 Explainable GEC as Generation Task

To investigate the interaction between explanation and correction tasks, we propose four different training settings, as illustrated in Figure 3: 1) no explanations (*Baseline*), which is the conventional setting, 2) explanations as additional input (*Infusion*), 3) explanations as output (*Explanation*), and 4) explanations as additional output (*Self-Rationalization*). To enable all these settings in a single architecture, we propose EXGEC, a unified generative framework for explainable GEC. In the Infusion setting, we introduce a special token “<sep>” to separate the source sentence and the following explanation, which includes evidence words and an error type. In the Explanation setting, the model generates an explanation given only a source sentence. As for the Self-rationalization setting, models are required to output a correction and an explanation separated by the special token “<sep>”. The relative positions of corrections and

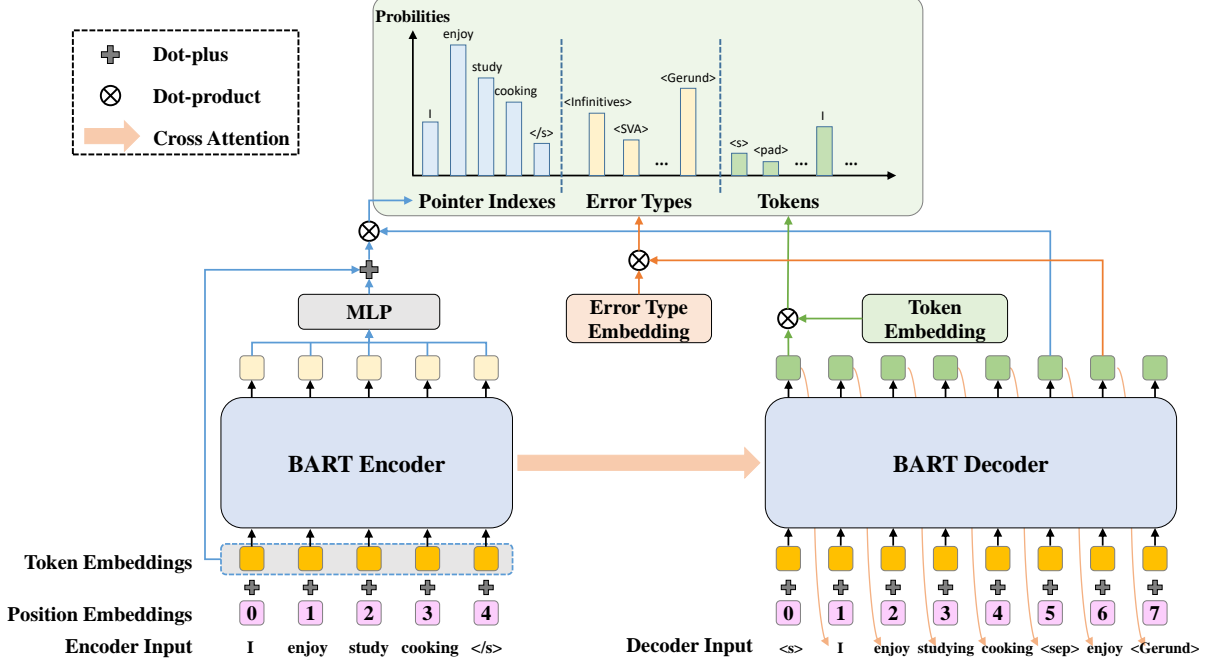


Figure 2: Overview of our Seq2Seq-based *Self-rationalization* model. The decoder can 1) output corrections from BART’s token vocabulary, 2) generate evidence words as source indexes by leveraging pointer mechanism, and 3) predict an error type from the predefined set of error type classes.

| | Input | Output |
|----------------------|---|---|
| Baseline | Source | Correction |
| Infusion | Source <sep> Evidence Words Error Type | Correction |
| Explanation | Source | Evidence Words Error Type |
| Self-rationalization | Source | Correction <sep> Evidence Words Error Type |

Figure 3: Comparison of four settings, all of which can be implemented in our proposed unified architecture.

explanations can be reversed, which allows us to understand the interaction between both tasks.

Without loss of generality, we clarify how our EXGEC tackles tasks in a unified generative framework in the *Self-rationalization* setting. Given an ungrammatical source sentence X , the encoder encodes X into hidden representation \mathbf{H} as follow:

$$\mathbf{H}^e = \text{Encoder}(X), \quad (1)$$

where $\mathbf{H}^e \in \mathbb{R}^{n \times d}$, and d is the hidden size.

At each time step t , the decoder produces the hidden state $\mathbf{h}_t^d \in \mathbb{R}^d$ based on the previous output sequence $\hat{Y}_{<t}$, which is computed as follow:

$$\mathbf{h}_t^d = \text{Decoder}(\mathbf{H}^e, \hat{Y}_{<t}). \quad (2)$$

Next, the hidden state $\mathbf{h}_t^d \in \mathbb{R}^d$ is utilized to calculate three types of logits: 1) *token logits*, which

are responsible for the correction part (Vaswani et al., 2017), 2) *pointer logits*, used to determine the probabilities of source indexes for evidence extraction, and 3) *type logits*, utilized for error type classification. Inspired by Yan et al. (2021), we calculate the probability distribution P_t as follows:

$$\mathbf{E}^e = \text{TokenEmbed}(X) \in \mathbb{R}^{n \times d}, \quad (3)$$

$$\bar{\mathbf{H}}^e = \alpha \mathbf{E}^e + (1 - \alpha) \text{MLP}(\mathbf{H}^e) \in \mathbb{R}^{n \times d}, \quad (4)$$

$$\mathbf{V}^d = \text{TokenEmbed}(V) \in \mathbb{R}^{|V| \times d}, \quad (5)$$

$$\mathbf{C}^d = \text{TypeEmbed}(C) \in \mathbb{R}^{|C| \times d}, \quad (6)$$

$$P_t = \text{softmax}([\mathbf{V}^d \otimes \mathbf{h}_t^d; \bar{\mathbf{H}}^e \otimes \mathbf{h}_t^d; \mathbf{C}^d \otimes \mathbf{h}_t^d]), \quad (7)$$

where TokenEmbed refers to the embeddings that are shared between the encoder and decoder, $\alpha \in \mathbb{R}$ is a hyper-parameter responsible for balancing the trade-off between embeddings and encoder hidden representation, V represents the token vocabulary, $[\cdot; \cdot]$ denotes the concatenation operation in the first dimension, the symbol \otimes means the dot product operation, and $P_t \in \mathbb{R}^{|V|+n+|C|}$ represents the probability distribution at the current time step t .

It is worth noting that the pointer index cannot be directly inputted to the decoder, so we introduce the `Index2Token` conversion to convert indexes into

tokens (Yan et al., 2021). Additionally, we can rearrange the generation order of corrections and explanations, which may provide helpful insight into further understanding the interaction of both tasks. In the Baseline and Infusion settings, the probability distribution is limited to the token vocabulary. However, in the Explanation setting, the probability distribution is limited to the combination of pointer indexes and error type classes.

3.3 Loss Weighting

Taking into account the heterogeneity of correction and explanation tasks, we construct the overall loss function in the form of weighted sum, which is defined as follow:

$$\begin{aligned} \mathcal{L} &= \mathcal{L}_{cor} + \lambda \cdot \mathcal{L}_{exp} \\ &= - \sum_{i=0}^m \left[\mathbb{I}(y_i \in V) \log p_i + \lambda \mathbb{I}(y_i \notin V) \log p_i \right], \end{aligned} \quad (8)$$

where λ is responsible for balancing both tasks, and \mathbb{I} is the indicator function. During the inference stage, we generate the entire target sequence in an autoregressive manner and then separate different parts from the target.

4 Experiments

4.1 Experimental Settings

Backbone model. We adopt the Seq2Seq-based pre-trained model BART-Large (Lewis et al., 2020) as our backbone model. All experiments are conducted using the open-source sequence modeling toolkit Fairseq (Ott et al., 2019), and subwords are obtained using the byte-pair-encoding (BPE) (Sennrich et al., 2016) algorithm. It is worth noting that adopting BART is non-trivial because the BPE tokenization may split a word into several BPEs, making it tricky to extract evidence words. Considering evidence words are usually short and not always contiguous, we stipulate that the pointer indexes should contain all BPEs of the evidence words. In other words, if a word is an evidence word, models in the Explanation and Self-rationalization settings are desired to output the pointer indexes of all its BPEs. If an instance has no evidence word, the target skips the prediction of pointer indexes. Additionally, we apply the Dropout-Src mechanism (Junczys-Dowmunt et al., 2018) to source-side word embeddings following

previous work (Zhang et al., 2022). Detailed hyperparameter settings are provided in Appendix A.

Training Settings. As discussed in Section 3.2, we attempt to conduct experiments on four distinct training settings leveraging a single unified framework with minimal modification. Notably, the Self-rationalization setting can be further divided into two settings based on the generation order of the correction and explanation parts: 1) *pre-explaining* models first output the explanation part and then the correction part, while 2) *post-explaining* models work in reverse order. In general, we extract evidence words first and then predict error types since we find that the generation order of evidence words and error types does not significantly affect the performance in our preliminary experiments.

Evaluation. We evaluate the model performance in three aspects. 1) Correction. Following the authors of the W&I+LOCNESS dataset (Bryant et al., 2019), we report correction performance evaluated by ERRANT (Bryant et al., 2017). 2) Extraction of evidence words. Following Fei et al. (2023), we also employ token-level evaluation metrics such as Precision, Recall, F_1 and $F_{0.5}$. However, we do *not* adopt the exact match (EM) metric since it is reported to be the least correlated with human evaluation⁵. The findings (Fei et al., 2023) show that the $F_{0.5}$ score achieves the highest correlation with human evaluation in terms of Pearson coefficient, followed by the F_1 score. 3) Classification of grammatical errors. We report label accuracy as the classification performance of grammatical error types. Unlike previous work (Fei et al., 2023), we disentangle the evaluation of extraction and classification, which might provide a clearer perspective on aspects of model performance. Specifically, we deem an evidence word as a True Positive (TP) if all of its BPEs are extracted, which is not in line with the previous evaluation (Fei et al., 2023) that considers an evidence word as a TP only if both BPEs and its error type are correctly predicted. The results are averaged over three runs with different random seeds, and the EXPECT-*dev* set serves as the validation set in all experiments.

4.2 Experiments on Rebuilt Datasets

To demonstrate the effectiveness of our rebuild process, we first respectively train post-explaining

⁵Surprisingly, we find that *do-nothing* systems achieve higher EM scores than almost all well-trained systems, but 0 F_1 and $F_{0.5}$ scores.

| System | EXPECT-dev | | | EXPECT-test | | |
|-----------------------------|----------------------------------|--|--|----------------------------------|--|--|
| | Cor. (P / R / F _{0.5}) | Exp. (P / R / F ₁ / F _{0.5} / Acc) | | Cor. (P / R / F _{0.5}) | Exp. (P / R / F ₁ / F _{0.5} / Acc) | |
| BART Baseline | 36.14 / 34.87 / 35.88 | - | | 36.33 / 35.49 / 36.16 | - | |
| BERT Explanation | - | 53.60 / 35.46 / 42.68 / 48.63 / 52.09 | | - | 51.73 / 36.34 / 42.69 / 47.69 / 50.83 | |
| BART Explanation | - | 44.43 / 32.93 / 37.82 / 41.53 / 33.36 | | - | 42.34 / 33.13 / 37.18 / 40.11 / 26.95 | |
| Infusion | | | | | | |
| + Evidence | 45.78 / 44.55 / 45.53 | - | | 46.02 / 44.13 / 45.63 | - | |
| + Type | 35.31 / 47.87 / 35.22 | - | | 36.00 / 35.37 / 35.87 | - | |
| + Evidence&Type | 44.28 / 47.55 / 44.90 | - | | 44.96 / 47.50 / 45.44 | - | |
| Self-rationalization | | | | | | |
| Pre-explaining | 38.25 / 34.18 / 37.36 | 36.01 / 35.58 / 35.79 / 35.92 / 26.56 | | 38.68 / 35.41 / 37.98 | 36.77 / 36.85 / 36.81 / 36.79 / 26.24 | |
| Post-explaining | 36.34 / 40.15 / 37.05 | 48.95 / 42.72 / 45.63 / 47.56 / 40.32 | | 36.52 / 40.41 / 37.24 | 49.43 / 44.10 / 46.61 / 48.26 / 39.86 | |

Table 3: Results of different settings for the single model. All models except ‘‘BERT Explanation’’ are initialized with pre-trained BART weights.

| Official EXPECT-dev | |
|----------------------------------|--|
| Cor. (P / R / F _{0.5}) | Exp. (P / R / F ₁ / F _{0.5} / Acc) |
| 30.94 / 35.49 / 31.75 | 45.92 / 38.42 / 41.84 / 44.19 / 37.63 |
| Rebuilt EXPECT-dev | |
| Cor (P / R / F _{0.5}) | Exp (P / R / F ₁ / F _{0.5} / Acc) |
| 36.34 / 40.15 / 37.05 | 48.95 / 42.72 / 45.65 / 47.56 / 40.32 |

Table 4: Comparison of *post-explaining* models trained on the official and rebuilt EXPECT datasets. We have similar findings on other settings, which are listed in Appendix B.1.

models on the official and our rebuilt EXPECT datasets. The results in Table 4 indicate that our rebuilt EXPECT dataset can significantly improve the performance of both correction and explanation tasks. This is because we have identified and corrected grammatical errors that were previously overlooked. As a results, we conduct the remaining experiments on the rebuilt EXPECT dataset.

4.3 Main Results

Here, we examine and analyze the interaction between the correction and explanation tasks by conducting experiments with various training settings. We first explore the Infusion setting, where we append different additional explanation information to the input source. Infusion models can be considered as oracle baselines since human-annotated explanations are usually unavailable in real applications, through which we can understand how explanations benefit the correction task. We also train a sequence labeling-based BERT model by reproducing the baseline provided in (Fei et al., 2023) under the same training and evaluation conditions as our other experiments. The results presented in Table 3 illustrate the following conclusions.

Evidence words, rather than grammatical error types, can provide invaluable information for corrections. Recent studies have highlighted

that incorporating human-annotated explanations as additional input can enhance task performance to a certain degree (Hase et al., 2020; Yao et al., 2023), and we have also observed similar results in the ‘‘Infusion’’ block of Table 3. Specifically, we notice that the additional information provided by grammatical error types does not improve correction performance. However, on the other hand, the information provided by evidence words can increase the F_{0.5} score by approximately 10 points, even though about only 60% of the samples in the dev and test sets are annotated with evidence words, demonstrating that ground truth evidence words are very helpful for the correction task.

Jointly learning correction and explanation tasks is beneficial for each task. Practically, explanations are usually unavailable during the inference stage, so Self-rationalization models are responsible for answering whether training with explanations as additional output could improve correction performance. Interestingly, experiments show that pre-explaining and post-explaining models perform differently. Specifically, pre-explaining models achieve better correction performance at the cost of decreased explanation performance compared to the ‘‘BART Explanation’’ single-task baseline, demonstrating that even noisy predicted explanations can still provide benefits towards the correction task. On the other hand, post-explaining models achieve comparable correction performance but very high explanation performance, indicating that predicted corrections are very beneficial towards the explanation task.

We also notice that the performance of grammatical error type classification for BART-based models is greatly lower than that of BERT-based models. We speculate that this may be due to the inner bias induced by the distinction between BART’s generative denoising and BERT’s masked language model

| γ | Cor. (P / R / F _{0.5}) | Exp. (P / R / F ₁ / F _{0.5} / Acc) |
|----------|----------------------------------|--|
| 0.5 | 36.16 / 35.68 / 36.06 | 57.00 / 06.87 / 12.26 / 23.18 / 19.15 |
| 0.8 | 35.47 / 36.92 / 35.74 | 51.77 / 21.63 / 30.51 / 40.49 / 23.46 |
| 1.0 | 35.10 / 36.96 / 35.46 | 48.82 / 26.55 / 34.40 / 41.81 / 25.94 |
| 1.5 | 36.12 / 36.34 / 36.16 | 50.95 / 22.01 / 30.74 / 40.34 / 24.66 |
| 2.0 | 35.93 / 35.38 / 35.82 | 52.48 / 22.29 / 31.29 / 41.29 / 28.06 |

Table 5: Results of sequence labeling-based multi-task BART baselines for varying loss weights γ on rebuilt **EXPECT-dev**.

(MLM) pre-training objectives. This is supported by the experiments in Section 5.1, which indicate that sequence labeling is not the crucial factor for grammatical error type classification.

5 Analysis

5.1 Does Sequence Labeling Help?

Motivated by recent studies in multi-task GEC frameworks (Zhao et al., 2019; Yuan et al., 2021), which combine a sequence labeling task with a sentence-level correction task, we also develop a multi-task baseline for explainable GEC, keeping the experimental setup the same as our other experiments. Specifically, we append a random-initialized tagging head after the encoder to perform the explanation task as a sequence labeling task, like BERT-based models. To predict each token’s tag, we pass the encoder’s hidden representation \mathbf{H}^e through a softmax after an affine transformation, which is computed as follow:

$$P(T | X) = \text{softmax}(W^\top \mathbf{H}^e), \quad (9)$$

where T is the resulting tagging sequence in BIO scheme. The token-level sequence labeling task is introduced to replace the role of pointer mechanism, so we conduct only the correction task at the decoder side. Similarly, we introduce loss weighting to trade-off the losses of correction generation and sequence labeling, which is defined as follow:

$$\mathcal{L} = \mathcal{L}_{cor} + \gamma \cdot \mathcal{L}_{tag} \quad (10)$$

where γ represents the trade-off factor, and we minimize the cross-entropy between predicted tokens/labels and ground truth tokens/labels.

The results of varying γ selected from the alternative set $\{0.5, 0.8, 1.0, 1.5, 2.0\}$ are shown in Table 5. Compared to Self-rationalization models, sequence labeling-based multi-task models achieve lower correction performance but mediate explanation performance between pre-explaining models and post-explaining models. Therefore, we can

conclude that our proposed EXGEC is more effective than sequence labeling-base baselines.

5.2 Position Leakage

One may suspect that the enhancement of Infusion models is due to the leakage effect of evidence words’ positions, since it is reported that a significant number of instances have at least one evidence word within the first or second-order nodes of correction words in the dependency parsing tree (Fei et al., 2023). To address this concern, we synthesize datasets with artifact explanations in two ways: 1) *random explanations*, which are randomly selected from the entire source tokens, and 2) *adjacent explanations*, which are randomly chosen from candidate source tokens located within a distance of 1~5 from the correction. Given that a substantial number of samples lack annotated evidence words, we generate an equal number of synthesized evidence words as the ground truth ones to ensure the fairness of our experiments. We train models using synthesized evidence words, but evaluation is performed with ground truth evidence words, allowing us to investigate whether the models learn to extract evidence words through this unsupervised approach. The results are presented in Table 6.

For the Infusion setting, it is no surprise that random evidence words would not improve correction performance as expected. However, we observe that adjacent synthesized evidence words do make a noticeable impact, resulting in a moderate improvement compared to random evidence words but still lower than the benefits provided by ground truth evidence words. This suggests that the leakage effect of positions does indeed exist. However, it is important to note that this effect alone is unable to fully capture all the advantages offered by ground truth evidence words.

For the pre-explaining and post-explaining settings, it seems that learning to output adjacent evidence words can improve correction performance to some extent. However, it falls short of surpassing the performance achieved by incorporating ground truth evidence words. This reaffirms the importance of joint learning for both correction and explanation tasks. On the contrary, the inclusion of random evidence words does not contribute to the improvement of correction performance. Furthermore, the models’ explanation performance reveals their inclination to disregard the influence of these random evidence words. Additionally, we observe a significant decrease in explanation per-

| System | EXPECT-dev | | | EXPECT-test | | |
|------------------------|------------------------------|--|--|------------------------------|--|--|
| | Cor. (P/R/F _{0.5}) | Exp. (P/R/F ₁ /F _{0.5} /Acc) | | Cor. (P/R/F _{0.5}) | Exp. (P/R/F ₁ /F _{0.5} /Acc) | |
| BART Baseline | 36.14 / 34.87 / 35.88 | - | | 36.33 / 35.49 / 36.16 | - | |
| Infusion | | | | | | |
| + G.T. Evidence | 45.78 / 44.55 / 45.53 | - | | 46.02 / 44.13 / 45.63 | - | |
| + Ran. Evidence | 35.88 / 33.26 / 35.33 | - | | 36.44 / 33.20 / 35.74 | - | |
| + Adj. Evidence | 38.46 / 42.81 / 39.26 | - | | 39.66 / 43.01 / 40.28 | - | |
| Pre-explaining | | | | | | |
| + G.T. Evidence | 38.25 / 34.18 / 37.36 | 36.01 / 35.58 / 35.79 / 35.92 / 26.56 | | 38.68 / 35.41 / 37.98 | 36.77 / 36.85 / 36.81 / 36.79 / 26.24 | |
| + Ran. Evidence | 36.17 / 33.72 / 35.65 | 13.60 / 00.40 / 00.77 / 01.79 / 15.83 | | 37.63 / 34.83 / 37.04 | 14.38 / 00.53 / 01.02 / 02.31 / 15.02 | |
| + Adj. Evidence | 36.53 / 38.73 / 36.95 | 26.97 / 03.37 / 06.00 / 11.23 / 17.03 | | 37.09 / 39.52 / 37.55 | 29.00 / 04.02 / 07.06 / 12.93 / 16.02 | |
| Post-explaining | | | | | | |
| + G.T. Evidence | 36.34 / 40.15 / 37.05 | 48.95 / 42.72 / 45.63 / 47.56 / 40.32 | | 36.52 / 40.41 / 37.24 | 49.43 / 44.10 / 46.61 / 48.26 / 39.86 | |
| + Ran. Evidence | 36.36 / 34.37 / 35.95 | 14.39 / 00.45 / 00.86 / 02.00 / 16.04 | | 36.86 / 34.87 / 36.44 | 07.45 / 00.16 / 00.32 / 00.74 / 15.02 | |
| + Adj. Evidence | 36.36 / 34.14 / 35.89 | 23.68 / 02.53 / 04.57 / 08.86 / 15.79 | | 37.34 / 35.18 / 36.88 | 26.74 / 03.28 / 05.84 / 11.00 / 15.48 | |

Table 6: Results of models trained on ground truth (G.T.), random (Ran.) or adjacent (Adj.) evidence words.

497 performance when learning without ground truth evi-
498 dence words, indicating the inherent challenge of
499 explaining with alignment to human preference in
500 an unsupervised way.

501 6 Related Works

502 **Explainable GEC.** Currently, most GEC sys-
503 tems are trained to correct errors without providing
504 explanations. To bridge the gap, recent studies have
505 explored several methods to facilitate the explain-
506 ability of GEC systems. One such method is the
507 feedback comment generation (FCG) task (Nagata,
508 2019; Nagata et al., 2021), which is designed to
509 automatically generate feedback comments such
510 as hints or explanatory notes for writing learning.
511 Hanawa et al. (2021) investigate three different ar-
512 chitectures for FCG and highlight the challenges
513 of the task. Another approach is Example-based
514 GEC (Kaneko et al., 2022; Vasselli and Watan-
515 abe, 2023), which improves explainability by re-
516 trieving examples similar to an input instance ac-
517 cording to pre-defined grammar rules. Kaneko
518 and Okazaki (2023) explore generating natural lan-
519 guage explanations by prompting large language
520 models (LLMs), showing the feasibility of eliciting
521 controlled and comprehensive explanations for
522 grammatical errors from LLMs. However, there
523 has been no work systematically exploring the in-
524 teraction between correction and explanation tasks.

525 **Learning with Explanations.** As an important
526 part of this work, Self-rationalization models
527 jointly generate task predictions and correspond-
528 ing explanations, aiming to improve explainabil-
529 ity or task performance of neural networks. Two
530 approaches that currently predominate the build-
531 ing of self-rationalization models are 1) extract-
532 ing highlight input tokens responsible for task pre-

533 dictions, known as extractive rationals (DeYoung
534 et al., 2020), and 2) generating natural language
535 explanations (Narang et al., 2020), which pro-
536 vide a natural interface between machine compu-
537 tation and human end-users. To improve upon the
538 task performance and trustworthiness of Seq2Seq
539 models, Lakhotia et al. (2021) develop an extrac-
540 tive fusion-in-decoder architecture in the ERASER
541 benchmark (DeYoung et al., 2020), which is a pop-
542 ular benchmark for rationale extraction across mul-
543 tiple datasets and tasks. Li et al. (2022a) propose
544 a joint text classification and rationale extraction
545 model to improve explainability and robustness.
546 Recognizing the complementarity of extractive
547 rationals and natural language explanations, Ma-
548 jumder et al. (2022) combine both ingredients in a
549 unified self-rationalization framework.

550 Powered by in-context learning (Brown et al.,
551 2020) and chain-of-thought (CoT) reasoning (Wei
552 et al., 2022; Chu et al., 2023) of LLMs, recent
553 works leverage the natural language explanations
554 generated by LLMs with chain-of-thought prompt-
555 ing (Lampinen et al., 2022; Li et al., 2023) to en-
556 hance the training of small reasoners using knowl-
557 edge distillation for task performance (Li et al.,
558 2022b; Ho et al., 2023; Hsieh et al., 2023) or faith-
559 fulness (Wang et al., 2023) improvement.

560 7 Conclusion

561 In this paper, we propose a unified generative
562 framework, EXGEC, designed to jointly perform
563 both correction and explanation tasks. EXGEC is
564 designed to be compatible with multiple training
565 settings, enabling us to understand and establish
566 the interaction between tasks. Additionally, we re-
567 build the existing noisy explainable GEC dataset,
568 EXPECT. Our experiments demonstrate the effec-
569 tiveness of our rebuild process and EXGEC.

570 Limitations

571 **Inherent nature of Seq2Seq-based models.** We
572 have noticed that our adopted backbone, BART,
573 falls short in explanation performance, including
574 extracting evidence words and classifying gram-
575 matical errors, compared to BERT-based models.
576 This can be attributed to BART’s inherent nature as
577 a sequence-to-sequence generative model. These
578 limitations may have a negative impact on correc-
579 tion performance, particularly for post-explaining
580 models that correct sentences based on previously
581 predicted explanations. In our future work, we in-
582 tend to explore a more effective approach to handle
583 and integrate both tasks.

584 **Inflexibility of structured explanations.** In the
585 era of large language models (LLMs), it has be-
586 come increasingly practical and favorable to ex-
587 press explanations as free-form natural language
588 texts. However, in this particular paper, we focus
589 our studies on structured explanations due to the
590 limited availability of free-form explanations in the
591 field of GEC. Nevertheless, we are committed to
592 advancing the development of explainable GEC
593 datasets in our future work, aiming to incorporate
594 more sophisticated and comprehensive approaches.

595 Ethics Statement

596 In this paper, we have identified significant noise
597 in the official EXPECT dataset, which has the po-
598 tential to create confusion during model training
599 and evaluation. To address this issue, we recon-
600 struct the EXPECT dataset to remove the noise,
601 resulting in an objective training and evaluation
602 pipeline. For our methods, we have exclusively uti-
603 lized source data from publicly accessible project
604 resources on legitimate websites, ensuring the ab-
605 sence of sensitive information. Furthermore, all the
606 baselines and datasets utilized in our experiments
607 are publicly available, and we have given credit to
608 the corresponding authors by citing their work.

609 References

610 John Bitchener, Stuart Young, and Denise Cameron.
611 2005. The effect of different types of corrective feed-
612 back on esl student writing. *Journal of second lan-
613 guage writing*, 14(3):191–205.

614 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
615 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
616 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
617 Askell, et al. 2020. Language models are few-shot

learners. *Advances in neural information processing
systems*, 33:1877–1901. 618 619

Christopher Bryant, Mariano Felice, Øistein E. Ander- 620
sen, and Ted Briscoe. 2019. [The BEA-2019 shared 621](#)
[task on grammatical error correction](#). In *Proceedings 622*
of the Fourteenth Workshop on Innovative Use of NLP 623
for Building Educational Applications, pages 52–75, 624
Florence, Italy. Association for Computational Lin- 625
guistics. 626

Christopher Bryant, Mariano Felice, and Ted Briscoe. 627
2017. [Automatic annotation and evaluation of error 628](#)
[types for grammatical error correction](#). In *Proceed- 629*
ings of the 55th Annual Meeting of the Association for 630
Computational Linguistics (Volume 1: Long Papers), 631
pages 793–805, Vancouver, Canada. Association for 632
Computational Linguistics. 633

Christopher Bryant, Zheng Yuan, Muhammad Reza 634
Qorib, Hannan Cao, Hwee Tou Ng, and Ted Briscoe. 635
2022. Grammatical error correction: A survey of the 636
state of the art. *arXiv preprint arXiv:2211.05166*. 637

Zheng Chu, Jingchang Chen, Qianglong Chen, Weijiang 638
Yu, Tao He, Haotian Wang, Weihua Peng, Ming Liu, 639
Bing Qin, and Ting Liu. 2023. A survey of chain of 640
thought reasoning: Advances, frontiers and future. 641
arXiv preprint arXiv:2309.15402. 642

Jay DeYoung, Sarthak Jain, Nazneen Fatema Rajani, 643
Eric Lehman, Caiming Xiong, Richard Socher, and 644
Byron C. Wallace. 2020. [ERASER: A benchmark to 645](#)
[evaluate rationalized NLP models](#). In *Proceedings 646*
of the 58th Annual Meeting of the Association for 647
Computational Linguistics, pages 4443–4458, Online. 648
Association for Computational Linguistics. 649

Tao Fang, Jinpeng Hu, Derek F. Wong, Xiang Wan, 650
Lidia S. Chao, and Tsung-Hui Chang. 2023. [Improv- 651](#)
[ing grammatical error correction with multimodal 652](#)
[feature integration](#). In *Findings of the Association 653*
for Computational Linguistics: ACL 2023, pages 654
9328–9344, Toronto, Canada. Association for Com- 655
putational Linguistics. 656

Yuejiao Fei, Leyang Cui, Sen Yang, Wai Lam, Zhen- 657
zhong Lan, and Shuming Shi. 2023. [Enhancing gram- 658](#)
[matical error correction systems with explanations](#). 659
In *Proceedings of the 61st Annual Meeting of the 660*
Association for Computational Linguistics (Volume 661
1: Long Papers), pages 7489–7501, Toronto, Canada. 662
Association for Computational Linguistics. 663

Shichun Gui. 2004. A cognitive model of corpus-based 664
analysis of chinese learners’ errors of english. *Mod- 665*
ern Foreign Languages(Quarterly), 27(2):129–139. 666

Kazuaki Hanawa, Ryo Nagata, and Kentaro Inui. 2021. 667
[Exploring methods for generating feedback com- 668](#)
[ments for writing learning](#). In *Proceedings of the 669*
2021 Conference on Empirical Methods in Natural 670
Language Processing, pages 9719–9730, Online and 671
Punta Cana, Dominican Republic. Association for 672
Computational Linguistics. 673

| | | | |
|-----|--|---|-----|
| 674 | Peter Hase, Shiyue Zhang, Harry Xie, and Mohit Bansal. | Can language models learn from explanations in context? In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 537–563, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. | 731 |
| 675 | | | 732 |
| 676 | | | 733 |
| 677 | | | 734 |
| 678 | | | 735 |
| 679 | | | |
| 680 | | | |
| 681 | Namgyu Ho, Laura Schmid, and Se-Young Yun. 2023. | Large language models are reasoning teachers . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 14852–14882, Toronto, Canada. Association for Computational Linguistics. | |
| 682 | | | |
| 683 | | | |
| 684 | | | |
| 685 | | | |
| 686 | | | |
| 687 | Cheng-Yu Hsieh, Chun-Liang Li, Chih-kuan Yeh, | | |
| 688 | Hootan Nakhost, Yasuhisa Fujii, Alex Ratner, Ranjay | | |
| 689 | Krishna, Chen-Yu Lee, and Tomas Pfister. 2023. | Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 8003–8017, Toronto, Canada. Association for Computational Linguistics. | |
| 690 | | | |
| 691 | | | |
| 692 | | | |
| 693 | | | |
| 694 | | | |
| 695 | | | |
| 696 | Marcin Junczys-Dowmunt, Roman Grundkiewicz, | | |
| 697 | Shubha Guha, and Kenneth Heafield. 2018. | Approaching neural grammatical error correction as a low-resource machine translation task . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 595–606, New Orleans, Louisiana. Association for Computational Linguistics. | |
| 698 | | | |
| 699 | | | |
| 700 | | | |
| 701 | | | |
| 702 | | | |
| 703 | | | |
| 704 | | | |
| 705 | | | |
| 706 | Masahiro Kaneko and Naoaki Okazaki. 2023. | Controlled generation with prompt insertion for natural language explanations in grammatical error correction. <i>arXiv preprint arXiv:2309.11439</i> . | |
| 707 | | | |
| 708 | | | |
| 709 | | | |
| 710 | Masahiro Kaneko, Sho Takase, Ayana Niwa, and Naoaki | | |
| 711 | Okazaki. 2022. | Interpretability for language learners using example-based grammatical error correction . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 7176–7187, Dublin, Ireland. Association for Computational Linguistics. | |
| 712 | | | |
| 713 | | | |
| 714 | | | |
| 715 | | | |
| 716 | | | |
| 717 | Diederik P Kingma and Jimmy Ba. 2014. | Adam: A method for stochastic optimization . <i>arXiv preprint arXiv:1412.6980</i> . | |
| 718 | | | |
| 719 | | | |
| 720 | Kushal Lakhota, Bhargavi Paranjape, Asish Ghoshal, | | |
| 721 | Scott Yih, Yashar Mehdad, and Srini Iyer. 2021. | FiDex: Improving sequence-to-sequence models for extractive rationale generation . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 3712–3727, Punta Cana, Dominican Republic. Association for Computational Linguistics. | |
| 722 | | | |
| 723 | | | |
| 724 | | | |
| 725 | | | |
| 726 | | | |
| 727 | | | |
| 728 | Andrew Lampinen, Ishita Dasgupta, Stephanie Chan, | | |
| 729 | Kory Mathewson, Mh Tessler, Antonia Creswell, | | |
| 730 | James McClelland, Jane Wang, and Felix Hill. 2022. | | |
| | | | 731 |
| | | | 732 |
| | | | 733 |
| | | | 734 |
| | | | 735 |
| | | | |
| | | | 736 |
| | | | 737 |
| | | | 738 |
| | | | 739 |
| | | | 740 |
| | | | 741 |
| | | | 742 |
| | | | 743 |
| | | | 744 |
| | | | |
| | | | 745 |
| | | | 746 |
| | | | 747 |
| | | | 748 |
| | | | 749 |
| | | | 750 |
| | | | |
| | | | 751 |
| | | | 752 |
| | | | 753 |
| | | | 754 |
| | | | 755 |
| | | | |
| | | | 756 |
| | | | 757 |
| | | | 758 |
| | | | 759 |
| | | | 760 |
| | | | |
| | | | 761 |
| | | | 762 |
| | | | 763 |
| | | | 764 |
| | | | 765 |
| | | | 766 |
| | | | 767 |
| | | | 768 |
| | | | |
| | | | 769 |
| | | | 770 |
| | | | 771 |
| | | | 772 |
| | | | 773 |
| | | | 774 |
| | | | |
| | | | 775 |
| | | | 776 |
| | | | 777 |
| | | | 778 |
| | | | 779 |
| | | | 780 |
| | | | 781 |
| | | | |
| | | | 782 |
| | | | 783 |
| | | | 784 |
| | | | 785 |
| | | | 786 |
| | | | 787 |
| | | | 788 |

| | | | |
|-----|--|--|--|
| 789 | Sharan Narang, Colin Raffel, Katherine Lee, Adam Roberts, Noah Fiedel, and Karishma Malkan. 2020. Wt5?! training text-to-text models to explain their predictions. <i>arXiv preprint arXiv:2004.14546</i> . | Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. 2015. Pointer networks. <i>Advances in neural information processing systems</i> , 28. | 844 845 846 |
| 793 | Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. <i>arXiv preprint arXiv:1904.01038</i> . | Peifeng Wang, Zhengyang Wang, Zheng Li, Yifan Gao, Bing Yin, and Xiang Ren. 2023. SCOTT: Self-consistent chain-of-thought distillation. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5546–5558, Toronto, Canada. Association for Computational Linguistics. | 847 848 849 850 851 852 853 |
| 797 | Sascha Rothe, Jonathan Mallinson, Eric Malmi, Sebastian Krause, and Aliaksei Severyn. 2021. A simple recipe for multilingual grammatical error correction. In <i>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)</i> , pages 702–707, Online. Association for Computational Linguistics. | Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in Neural Information Processing Systems</i> , 35:24824–24837. | 854 855 856 857 858 |
| 798 | | Hang Yan, Tao Gui, Junqi Dai, Qipeng Guo, Zheng Zhang, and Xipeng Qiu. 2021. A unified generative framework for various NER subtasks. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 5808–5822, Online. Association for Computational Linguistics. | 859 860 861 862 863 864 865 866 |
| 799 | | Bingsheng Yao, Prithviraj Sen, Lucian Popa, James Hendler, and Dakuo Wang. 2023. Are human explanations always helpful? towards objective evaluation of human natural language explanations. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 14698–14713, Toronto, Canada. Association for Computational Linguistics. | 867 868 869 870 871 872 873 874 |
| 800 | Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1715–1725, Berlin, Germany. Association for Computational Linguistics. | Zheng Yuan and Christopher Bryant. 2021. Document-level grammatical error correction. In <i>Proceedings of the 16th Workshop on Innovative Use of NLP for Building Educational Applications</i> , pages 75–84, Online. Association for Computational Linguistics. | 875 876 877 878 879 |
| 801 | | Zheng Yuan, Shiva Taslimipoor, Christopher Davis, and Christopher Bryant. 2021. Multi-class grammatical error detection for correction: A tale of two systems. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 8722–8736, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics. | 880 881 882 883 884 885 886 |
| 802 | | Yue Zhang, Bo Zhang, Haochen Jiang, Zhenghua Li, Chen Li, Fei Huang, and Min Zhang. 2023. NaSGEC: a multi-domain Chinese grammatical error correction dataset from native speaker texts. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 9935–9951, Toronto, Canada. Association for Computational Linguistics. | 887 888 889 890 891 892 893 |
| 803 | | Yue Zhang, Bo Zhang, Zhenghua Li, Zuyi Bao, Chen Li, and Min Zhang. 2022. SynGEC: Syntax-enhanced grammatical error correction with a tailored GEC-oriented parser. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 2518–2531, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics. | 894 895 896 897 898 899 900 901 |
| 804 | | | |
| 805 | Younghee Sheen. 2007. The effect of focused written corrective feedback and language aptitude on esl learners’ acquisition of articles. <i>TESOL quarterly</i> , 41(2):255–283. | | |
| 806 | | | |
| 807 | | | |
| 808 | | | |
| 809 | | | |
| 810 | | | |
| 811 | | | |
| 812 | | | |
| 813 | | | |
| 814 | | | |
| 815 | | | |
| 816 | | | |
| 817 | | | |
| 818 | | | |
| 819 | | | |
| 820 | | | |
| 821 | | | |
| 822 | | | |
| 823 | | | |
| 824 | | | |
| 825 | | | |
| 826 | | | |
| 827 | | | |
| 828 | | | |
| 829 | | | |
| 830 | | | |
| 831 | | | |
| 832 | | | |
| 833 | | | |
| 834 | | | |
| 835 | | | |
| 836 | | | |
| 837 | | | |
| 838 | | | |
| 839 | | | |
| 840 | | | |
| 841 | | | |
| 842 | | | |
| 843 | | | |

Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu, Huiqi Deng, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, and Mengnan Du. 2023. Explainability for large language models: A survey. *arXiv preprint arXiv:2309.01029*.

Wei Zhao, Liang Wang, Kewei Shen, Ruoyu Jia, and Jingming Liu. 2019. Improving grammatical error correction via pre-training a copy-augmented architecture with unlabeled data. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 156–165, Minneapolis, Minnesota. Association for Computational Linguistics.

A Experiment Hyper-Parameters

We list the main hyper-parameters in Table 7. For the training stage, we follow the same hyper-parameters as described in (Zhang et al., 2022). The total training time is about 4 hours.

| Configuration | Value |
|-----------------------|---|
| Training | |
| Backbone | BART-large (Lewis et al., 2020) |
| Devices | 1 Tesla A100 GPU (80GB) |
| Epochs | 60 |
| Batch size per GPU | 4096 tokens |
| Gradient Accumulation | 4 |
| Optimizer | Adam (Kingma and Ba, 2014) |
| | $(\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 1 \times 10^{-8})$ |
| Learning rate | 3×10^{-5} |
| Warmup updates | 500 |
| Max source length | 256 |
| Dropout | 0.3 |
| Dropout-src | 0.1 |
| α | 0.5 |
| Loss weight λ | 1.0 |
| Inference | |
| Beam size | 12 |
| Max length | 256 |

Table 7: Hyper-parameter values used in our experiments.

B Extra Analyses

B.1 Detailed Results on EXPECT Datasets

We report the detailed results on the official and our rebuilt EXPECT datasets in Table 9. All the models trained on our rebuilt EXPECT achieve better performance of both correction and explanation tasks, demonstrating the effectiveness of our rebuild process.

B.2 Impact of Loss Weighting

In this section, we investigate the trade-off of learning on both correction and explanation task by varying the loss weight λ . Considering the promising performance of post-explaining models on both

| λ | Cor. (P / R / F _{0.5}) | Exp. (P / R / F ₁ / F _{0.5} / Acc) |
|-----------|----------------------------------|--|
| 0.5 | 35.40 / 38.03 / 35.90 | 39.77 / 38.88 / 39.32 / 39.59 / 32.02 |
| 1.0 | 36.34 / 40.15 / 37.05 | 48.95 / 42.72 / 45.63 / 47.56 / 40.32 |
| 1.5 | 36.03 / 38.42 / 36.49 | 43.90 / 42.82 / 43.35 / 43.68 / 36.88 |
| 2.0 | 35.41 / 38.61 / 36.00 | 47.98 / 42.86 / 45.28 / 46.86 / 40.07 |

Table 8: Results of *post-explaining* models for varying loss weights λ on rebuilt **EXPECT-dev**.

correction and explanation tasks, we train post-explaining models with the loss weight λ alternatively selected from $\{0.5, 1.0, 1.5, 2.0\}$ and report the results on EXPECT-dev in Table 8. The results show that giving preference to either tasks harms the performance of both tasks. We speculate that the supervised explanation information during training is too weak to guide the dynamics of correction learning if λ is small. On the other hand, a large λ value might neglect correction learning, thus leading to lower explanation performance since explanation of post-explaining models are produced based on predicted corrections.

| System | Official EXPECT-dev | | Rebuilt EXPECT-dev | |
|-----------------------------|----------------------------------|--|----------------------------------|--|
| | Cor. (P / R / F _{0.5}) | Exp. (P / R / F ₁ / F _{0.5} / Acc) | Cor. (P / R / F _{0.5}) | Exp. (P / R / F ₁ / F _{0.5} / Acc) |
| BART Baseline | 30.59 / 33.72 / 31.17 | - | 36.14 / 34.87 / 35.88 | - |
| Infusion | | | | |
| + Evidence | 40.72 / 43.31 / 41.22 | - | 45.78 / 44.55 / 45.53 | - |
| + Type | 31.15 / 35.14 / 31.87 | - | 35.31 / 47.87 / 35.22 | - |
| + Evidence&Type | 40.79 / 42.50 / 41.11 | - | 44.28 / 47.55 / 44.90 | - |
| Self-rationalization | | | | |
| Pre-explaining | 32.62 / 31.29 / 32.35 | 33.75 / 44.12 / 38.25 / 35.41 / 28.22 | 38.25 / 34.18 / 37.36 | 36.01 / 35.58 / 35.79 / 35.92 / 26.56 |
| Post-explaining | 30.94 / 35.49 / 31.75 | 45.92 / 38.42 / 41.84 / 44.19 / 37.63 | 36.34 / 40.15 / 37.05 | 48.95 / 42.72 / 45.63 / 47.56 / 40.32 |

Table 9: Further comparison of models trained on the official and rebuilt EXPECT datasets.