

Rationalize and Align: Enhancing Writing Assistance with Rationale via Self-Training for Improved Alignment

Hannan Cao Hai Ye Hwee Tou Ng

Department of Computer Science, National University of Singapore

caoh@u.nus.edu, e0684045@u.nus.edu

nght@comp.nus.edu.sg

Abstract

A Writing Assistant (WA) is a system that offers writing suggestions based on user instructions. Existing WAs are typically built by training large language models (LLMs) on domain-specific instruction data through supervised fine-tuning (SFT) only. However, SFT optimizes models to match a single reference, failing to capture the inherent flexibility of text editing, where multiple valid revisions exist. Therefore, solely relying on SFT limits WA performance. To address this limitation, we propose the **Rationalize and Align** framework, which enhances the WA performance with rationale (*i.e.*, linguistic explanations) and alignment. Our framework automatically generates the rationale and preference data for writing tasks via distillation and self-training, eliminating the need for human annotation. These data are then leveraged to refine WA using a novel preference optimization method. Empirical results show that our framework significantly improves WA performance. Our WA outperforms both open-source state-of-the-art WAs and the closed-source GPT-4o by 3.9 and 7.1 points on average, respectively, across eight well-established writing-related test sets.¹

1 Introduction

Large language models (LLMs) demonstrate strong proficiency in generating fluent and coherent text revisions, assisting users with various writing tasks (Achiam et al., 2023). However, general-purpose LLMs, such as GPT-4, often underperform compared to task-specific models (Cao et al., 2023b; Raheja et al., 2023). Consequently, research has focused on tailoring LLMs for writing-related tasks, such as grammatical error correction (GEC) (Bryant et al., 2023), text simplification (Baez and Saggion, 2023), and style transfer (Luo et al., 2023). Several studies (Raheja et al., 2023; Schick et al., 2023) have proposed writing assistants (WAs) that

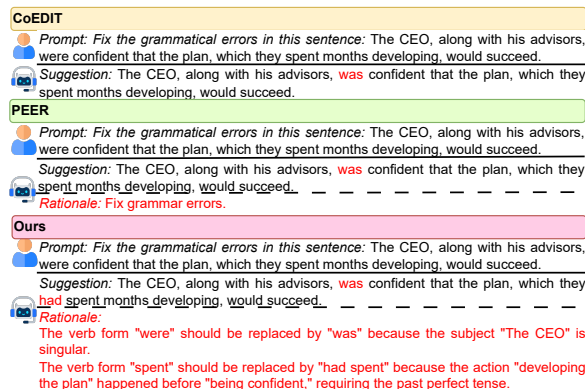


Figure 1: Comparison of our WA and related systems' output: CoEDIT (Raheja et al., 2023) and PEER (Schick et al., 2023).

generate text suggestions based on user instructions.

State-of-the-art (SOTA) WAs (Raheja et al., 2023; Zhang et al., 2023) are built by directly applying supervised fine-tuning (SFT) on labeled data. However, text editing inherently allows multiple valid revisions for a given input. Since SFT optimizes models to match a single reference revision, it fails to capture this flexibility. In contrast, task-specific evaluation metrics (*e.g.*, SARI, ROUGE) may account for such variation (Paulus et al., 2018). Consequently, relying solely on SFT may result in suboptimal performance. Furthermore, existing open-source WAs cannot provide rationales (*i.e.*, linguistic explanations) for their suggestions (Fig. 1), limiting their applicability.

To bridge this gap, we introduce the **Rationalize and Align** framework, which enhances WA with rationales and aligns WA towards suggestions with higher overall quality (*e.g.*, fluency, coherence). However, the absence of rationale and preference data presents a significant challenge, and acquiring human annotations is both time-consuming and expensive. To address this, we propose to automatically generate these data via *distillation* and *self-training*. Since existing WAs are not able to

¹<https://github.com/nusnlp/RationAlign>

provide rationales for their suggestions, we *distill* rationales (Fig. 2(a)) from GPT-4o into the labeled data, creating rationale-enhanced labeled datasets. These datasets are then used to develop a rationale-enhanced WA through SFT.

We then generate preference data to further enhance WA. However, text editing tasks lack appropriate reward models. Moreover, widely used evaluation metrics in writing tasks are reference-based metrics, restricting the generation of preference data to labeled datasets (*i.e.*, containing labeled suggestions for user instructions). However, such labeled data is scarce and costly to obtain. To address this, we propose a *self-training* alignment method that leverages unlabeled data to align WA. This method consists of two phases: (1) *reward modeling* (Fig. 2 (b)), which builds a reward model to assess suggestion quality, and (2) *preference optimization* (Fig. 2 (c)), which creates preference data and aligns WA towards higher-quality suggestions.

During *reward modeling*, we utilize labeled data (*i.e.*, data used for SFT) and evaluation metrics to generate preference data, which is then used to build the reward model. In the subsequent *preference optimization* phase, we leverage the constructed reward model and rationale-enhanced WA to create new preference data from unlabeled data for WA. This preference data, along with its corresponding reward assigned by the reward model, is then used to align WA via a novel margin-based preference optimization loss. This loss extends DPO (Rafailov et al., 2023) by replacing the reference model with the reward margin between preferred and non-preferred data, thereby capturing their quality differences more effectively.

Our experiments demonstrate that the Rationalize & Align framework significantly enhances our WA’s performance. Across eight widely-used writing-related tasks, our WA significantly outperforms both open-source SOTA WA and GPT-4o, by 3.9 and 7.1 points on average, respectively. The contributions of our paper are as follows:

- We present the first open-source explainable WA that provides rationales for its suggestions.
- We introduce the Rationalize & Align framework to enhance WA. Our analysis demonstrates that each component contributes effectively.
- To the best of our knowledge, this is the first study to assess the impact of incorporating rationales into WA. Extensive experiments confirm that rationales enhance performance and increase confidence in the generated suggestions.

2 Related Work

Alignment. DPO (Rafailov et al., 2023) serves as a prominent approach in offline preference optimization. However, its reliance on a reference model and the absence of a reward model restrict its ability to sample preference pairs from the optimal policy (Meng et al., 2024). Recent advancements have aimed to eliminate the need for a reference model. CPO (Xu et al., 2024) entirely removes the reference model, whereas SimPO (Meng et al., 2024) replaces it with a length regularization term and integrates a pre-defined margin. Similar to SimPO, ODPO (Amini et al., 2024) integrates an offset into DPO while retaining the reference model. Unlike these approaches, our approach leverages the reward values generated by our reward model and replaces the reference model with a dynamically adjusted reward difference, tailored according to the quality of different preference data.

Enhancing LLMs with Rationale. Incorporating rationales into LLMs effectively enhances their faithfulness and performance across various tasks, including reasoning (Krishna et al., 2023; Wei et al., 2022), question answering (Kassner et al., 2023), and continual relation extraction (Xiong et al., 2023). Unlike previous studies, ours focuses on leveraging writing-related rationales to improve the performance of LLMs in writing tasks.

Writing Assistants. Previous WAs (Raheja et al., 2023; Zhang et al., 2023) provide suggestions based on user instructions yet lack sufficient explanatory support. In contrast, the PEER system (Schick et al., 2023), trained using the edit history of Wikipedia, offers high-level rationales. However, these rationales (Fig. 1) tend to be too abstract to effectively assist users in understanding or validating the suggestions. Our system overcomes these limitations by providing writing suggestions and detailed linguistic explanations for each edit.

Distillation and Self-Training. Recent studies (Cao et al., 2023b; Zhao et al., 2024; Kadlčík and Štefánik, 2024) demonstrate the effectiveness of distillation and self-training across diverse NLP tasks. Specifically, Cao et al. (2023b) employ self-training to build a strong unsupervised GEC system; Kadlčík and Štefánik (2024) leverage self-training to improve LLM reasoning; and Zhao et al. (2024) combine distillation with self-training to boost LLMs’ parsing performance.

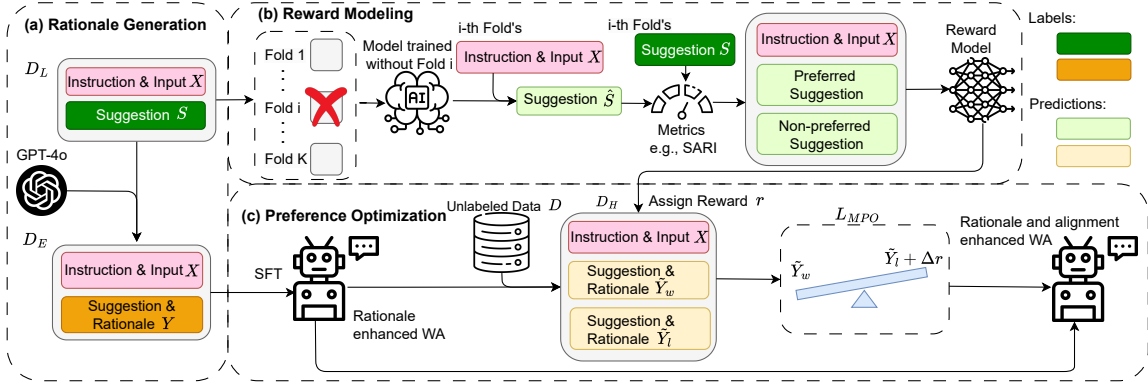


Figure 2: Our WA involves the following steps: (a) generating rationales for labeled data by distilling the output from GPT-4o, (b) constructing the reward model with \mathcal{D}_L and evaluation metrics, and (c) building a rationale and alignment enhanced WA with self-generated preference data through the preference optimization process.

3 Rationalize & Align

WAs aim to generate a proper writing suggestion s based on the user’s instruction and the input text. Existing WAs (Raheja et al., 2023; Zhang et al., 2023) focus solely on providing writing suggestions, which are trained to learn the mapping from x to s , where x represents the combination of instruction and input text. Here, we want to enhance the WA by providing a rationale (or explanation) e for its s . Therefore, our goal is to learn a mapping from x to y , where y represents the combination of suggestion s and rationale e . Fig. 2 shows the method to build our WA, comprising three key components: rationale generation, reward modeling, and preference optimization.

3.1 Rationale Generation

To enhance a WA with rationales in writing tasks, it is essential to provide linguistic explanations for its suggestions. However, current open-source WAs lack this capability, and no existing datasets have rationale annotations. To address this, we propose distilling rationales from the advanced closed-source model GPT-4o² into open-source models. Song et al. (2024) demonstrated that GPT-4 generates high-quality rationales (e) (92.2% accuracy in English) when given the input text (in), suggestion (s), and necessary edits ($edits$) that transform in into s , formulated as $e = \text{GPT-4}(in, s, edits)$.

Following their approach, we use GPT-4o to generate rationales using these three components. Specifically, we leverage the MaxMatch toolkit (Dahlmeier and Ng, 2012) to identify the edits required to transform the input into the suggestion. These *edits*, along with the *input* and *suggestion*,

are incorporated into prompts (Fig. 3) that instruct GPT-4o to generate the *rationale*. This process yields a rationale-enhanced dataset $\mathcal{D}_E = \{x, s, e\}$ from the labeled dataset $\mathcal{D}_L = \{x, s\}$.

3.2 Self-Training Alignment

Aligning LLMs requires preference data, which is unavailable for writing-related tasks. To address this, we propose a self-training alignment framework that utilizes widely available unlabeled data to perform alignment. The framework comprises two key steps: *reward modeling* and *preference optimization*. In the *reward modeling* step, we use labeled data (\mathcal{D}_L) and task-specific evaluation metrics to generate preference data for training the reward model. During *preference optimization*, we first construct an initial WA on \mathcal{D}_E . We then utilize this initial WA and the reward model to generate high-quality preference data, which are subsequently used to further refine the WA.

You are given a pair of English sentences along with a list of atomic edits. For each edit, the first word identifies content in the source sentence that is less appropriate, while the second word suggests a better phrase in the target sentence. [Task Instruction] Please generate a succinct explanation for each edit using the following template:

The word X should be deleted/inserted/replaced by Y because ...

###Source sentence:
[Input Text]

###Target sentence:
[Suggestion]

###Edits:
[dit ontent]

###Explanation:

Figure 3: Prompt to generate rationale. [Task Instruction] for each task is shown in Table 10.

²<https://chat.chatbotapp.ai/?model=gpt-4o>

$$\mathcal{L}_{DPO} = -\mathbb{E}_{(x, \tilde{y}_w, \tilde{y}_l) \sim \mathcal{D}_H} \left[\log \sigma \left(\beta \log \frac{\theta_W(\tilde{y}_w | x)}{\theta_{\text{SFT}}(\tilde{y}_w | x)} - \beta \log \frac{\theta_W(\tilde{y}_l | x)}{\theta_{\text{SFT}}(\tilde{y}_l | x)} \right) \right] \quad (1a)$$

$$= -\mathbb{E}_{(x, \tilde{y}_w, \tilde{y}_l) \sim \mathcal{D}_H} \left\{ \log \sigma \left[\beta (\log \theta_W(\tilde{y}_w | x) - \log \theta_W(\tilde{y}_l | x) - \underbrace{(\log \theta_{\text{SFT}}(\tilde{y}_w | x) - \log \theta_{\text{SFT}}(\tilde{y}_l | x))}_{\text{margin}}) \right] \right\} \quad (1b)$$

$$\mathcal{L}_M = -\mathbb{E}_{(x, \tilde{y}_w, \tilde{y}_l) \sim \mathcal{D}_H} \left\{ \log \sigma \left[\beta (\log \theta_W(\tilde{y}_w | x) - \log \theta_W(\tilde{y}_l | x) - \underbrace{\gamma (r_w - r_l)}_{\text{margin}}) \right] \right\} \quad (1c)$$

3.2.1 Reward Modeling

Our objective of reward modeling is to develop a robust reward model, r_ϕ , that can accurately evaluate the quality of any suggestion s given input x without relying on human-labeled suggestions. We achieve this by generating multiple pseudo-suggestions through sampling and then using automatic metrics to construct preference pairs, which are then used to build the reward model.

Specifically, we first partition the entire \mathcal{D}_L into K folds. For each fold, a model θ_k is trained on $K - 1$ folds to learn the mapping from x to s . Subsequently, we use θ_k to generate N responses \hat{s} for the reserved i -th fold through sampling. We then create preference pair (\hat{s}_p, \hat{s}_n) for each x by randomly selecting two predictions from these responses and evaluating them against the true label s using automatic metrics. The response with the higher score is labeled as the preferred output \hat{s}_p , and the one with the lower score as the non-preferred output \hat{s}_n . Following Wang et al. (2024), we build the reward model with the objective function: $\mathcal{L}_r = -\log(\sigma(r_\phi(\hat{s}_p) - r_\phi(\hat{s}_n)))$, where σ is the sigmoid function.

3.2.2 Preference Optimization

Labeled data is often scarce, while unlabeled data is abundant. Therefore, we propose a preference optimization method that generates preference data from unlabeled data. The preference data is generated through sampling, utilizing both a reward model and an initial WA model. We build this initial WA model, θ_W , by fine-tuning LLM on \mathcal{D}_E through SFT. To differentiate this initial WA model from our final WA model, we refer to this initial model as the SFT model and the final model as the WA model. This SFT model is optimized to generate the suggestion s followed by its rationale e , with the cross-entropy loss: $\mathcal{L}_{CE}(y) = -\log \theta_W(y|x)$.

Preference Data Generation. Utilizing the SFT model, we sample multiple responses and employ the reward model to identify the preferred and non-preferred responses, constructing preference pairs

from unlabeled data. Specifically, for each input x , we sample M responses with θ_W , where each response contains a writing suggestion \tilde{s} and its corresponding rationale. The reward model then assigns a score to each \tilde{s} , designating the highest-scoring suggestion as the winning suggestion \tilde{s}_w and the lowest-scoring as the losing suggestion \tilde{s}_l . Each suggestion is paired with its rationale to form the winning response \tilde{y}_w and the losing response \tilde{y}_l , resulting in a preference pair $(\tilde{y}_w, \tilde{y}_l)$ for each x . The reward for \tilde{s}_w and \tilde{s}_l are denoted as r_w and r_l , respectively.

However, suggestions generated through sampling exhibit varying quality (Beigi et al., 2024). To ensure that θ_W learns only from high-quality suggestions, we filter out lower-quality ones based on the reward r_w . Specifically, we rank the preference pairs by r_w and retain only the top 20% of the preference pairs with the highest reward. We denote this subset of high-quality preference pairs and their corresponding inputs as $\mathcal{D}_H = \{x, \tilde{y}_w, \tilde{y}_l\}$. Since no established methods exist for evaluating rationale quality, we rank suggestions and rationales together based on the reward value of the suggestions, hypothesizing that higher-quality suggestions correlate with better rationales.

Optimization. Eq. (1a) and Eq. (1b) present the DPO (Rafailov et al., 2023) loss. As shown in Eq. (1b), DPO ensures that the log probability of the winning response, \tilde{y}_w , exceeds that of the losing response, \tilde{y}_l , by a margin equal to their log probability difference under the SFT model, θ_{SFT} . Within our framework, since the reward model assesses suggestion quality, the quality gap between the winning and losing responses can be captured by the reward difference, $r_w - r_l$. With an effective r_ϕ (Fig.4), we hypothesize that the reward difference more accurately captures this margin³, motivating us to replace the log probability difference computed by θ_{SFT} with the reward difference. Consequently, we define our margin loss, \mathcal{L}_M , in Eq. (1c), where γ controls the weight assigned to

³Fig. 5 demonstrates the importance of this margin.

	System	CoN	ITR-F	ITR-L	ITR-O	STS	WNC	AST	GYAFC		ALL
									EM	FR	
a)	Llama-3.3-70B-Instruct	55.6	46.5	31.4	31.0	34.6	31.8	46.4	59.5 / 98.1	56.2 / 98.4	43.7
	ChatGPT	53.3	50.9	31.5	31.0	39.9	36.3	47.0	57.7 / 99.6	60.4 / 99.5	45.3
	GPT-4	59.9	51.6	32.6	32.3	42.2	40.8	46.3	60.2 / 99.6	62.4 / 99.5	47.6
	GPT-4o	59.4	51.1	32.4	32.4	42.4	41.1	47.4	62.8 / 99.2	63.7 / 99.1	48.1
	o1-preview	59.7	51.4	32.6	32.7	42.6	41.7	47.2	62.9 / 98.3	63.1 / 98.4	48.2
	Deepseek-R1	60.1	50.9	32.3	32.6	42.3	41.1	47.6	62.8 / 98.1	63.1 / 98.5	48.1
b)	Llama-3.3-70B-Instruct (R)	58.4	49.4	35.0	31.8	37.7	41.2	44.7	63.6 / 97.8	65.3 / 98.2	47.5
	ChatGPT (R)	56.1	51.1	30.3	28.7	40.6	36.6	45.0	63.1 / 98.7	63.5 / 98.9	46.1
	GPT-4 (R)	60.4	50.1	33.3	32.8	41.2	40.7	47.6	63.2 / 98.8	63.5 / 99.1	48.1
	GPT-4o (R)	60.8	51.4	32.6	32.2	43.3	40.9	46.1	64.4 / 98.4	64.6 / 98.6	48.5
	o1-preview (R)	60.9	50.6	33.3	32.9	41.2	41.1	47.5	63.0 / 98.4	64.1 / 98.9	48.3
	Deepseek-R1 (R)	60.4	51.5	32.3	32.5	43.0	40.6	45.8	64.5 / 99.1	64.3 / 99.2	48.4
c)	PEER-EDIT-11B	N.A.	52.1	32.5	32.7	28.2	54.5	29.5	N.A.	N.A.	N.A.
	Writing-Alpaca (7B)	55.9	52.8	39.4	37.1	44.6	64.4	44.7	N.A.	N.A.	N.A.
	CoEDIT-xxl (11B)	57.1*	51.6	31.8	31.5	42.9*	71.0	41.7	66.0 / 98.7*	68.7 / 97.9*	51.7
Ours based on Flan-T5-xxl (11B) (RationAlign-11B)											
d)	SFT model	58.3	50.9	33.6	32.2	43.0	70.8	41.4	69.2 / 97.3	70.5 / 97.1	52.1
	+ Self-Training Alignment	61.4	49.3	32.8	34.7	47.0	68.9	41.1	75.3 / 96.8	78.0 / 96.3	54.3
e)	SFT model (R)	61.8	51.3	30.2	36.1	46.6	69.0	43.1	73.4 / 97.4	76.2 / 97.1	54.1
	+ Self-Training Alignment (R)	62.1	52.5	33.5	38.6	44.7	70.2	42.8	75.6 / 97.1	77.4 / 96.9	55.2
Ours based on Llama 3.1 8B (RationAlign-8B)											
f)	SFT model	61.7	50.5	31.6	35.6	43.3	66.4	42.0	75.3 / 97.9	75.5 / 96.8	53.5
	+ Self-Training Alignment	62.5	48.7	31.9	40.3	47.1	64.6	45.2	78.0 / 96.8	78.0 / 95.9	55.1
g)	SFT model (R)	62.7	51.1	34.0	37.5	45.1	65.8	41.3	75.9 / 98.4	76.4 / 97.5	54.4
	+ Self-Training Alignment (R)	65.5	48.7	35.4	37.6	46.5	65.9	45.2	77.3 / 97.9	78.2 / 97.3	55.6

Table 1: Performance on writing-related tasks. All results are shown in %. *: Results reproduced using the official checkpoint and scripts released by [Raheja et al. \(2023\)](#), due to different evaluation metrics or test sets not previously evaluated. For the GYAFC test sets, the first score is BLEU and the second is accuracy. Following [Raheja et al. \(2023\)](#); [Zhang et al. \(2023\)](#), we show the averaged result under the ALL column, and we only consider the BLEU score for the GYAFC test sets when taking the average. **a)**: zero-shot performance of LLMs. **b)**: zero-shot performance of LLMs when also prompted to generate rationales (or explanations) for their writing suggestions. **c)**: SOTA WAs. **d) & f)**: RationAlign trained without rationale. **e) & g)**: RationAlign trained with rationale.

the reward difference. To better align WA with the distribution of winning responses ([Xu et al., 2024](#)), we integrate \mathcal{L}_{CE} , computed for \tilde{y}_w , with \mathcal{L}_M and propose the margin-based preference optimization loss, \mathcal{L}_{MPO} . We then optimize θ_W with \mathcal{L}_{MPO} :

$$\mathcal{L}_{MPO} = \lambda \mathcal{L}_M + \mathcal{L}_{CE}(\tilde{y}_w) \quad (2)$$

where λ is a hyper-parameter, controlling the weight assigned to \mathcal{L}_M . SimPO ([Meng et al., 2024](#)) also incorporates a margin, but our approach differs in motivation. Unlike SimPO, which uses a fixed margin as a hyper-parameter, our margin is determined by the reward difference, allowing it to adapt to varying input qualities. Furthermore, SimPO is not as effective as ours (Fig. 5).

4 Experimental Results

4.1 Data and Model Configuration

Data. Following previous works ([Raheja et al., 2023](#); [Zhang et al., 2023](#)), we evaluate our system’s performance in generating suggestions across eight writing-related tasks. These tasks closely align with EDITEVAL ([Dwivedi-Yu et al., 2022](#)): GEC,

fluency, clarity, coherence, paraphrasing, neutralization, simplification, and formality style transfer (FST). For GEC, we use the CoNLL-2014 (CoN) test set ([Ng et al., 2014](#)). Fluency, clarity, and coherence are assessed using the ITERATER test sets ([Du et al., 2022](#)), corresponding to ITR-F, ITR-L, and ITR-O, respectively. Paraphrasing is evaluated with the STSB (STS) test set ([Dwivedi-Yu et al., 2022](#)), while neutralization is assessed using the WNC test set ([Pryzant et al., 2020](#)). Simplification is evaluated using the ASSET (AST) test set ([Alva-Manchego et al., 2020](#)). FST is evaluated on the GYAFC test set ([Rao and Tetreault, 2018](#)), covering both the Entertainment & Music (EM) and Family & Relationships (FR) domains.

Metrics. Following prior research ([Raheja et al., 2023](#); [Schick et al., 2023](#); [Zhang et al., 2023](#)), we primarily utilize the SARI metric ([Xu et al., 2016](#)) for evaluation, except for the GEC and FST tasks. For GEC, we use the MaxMatch scorer ([Dahlmeier and Ng, 2012](#)), as done in previous studies ([Cao et al., 2023a](#)). For FST, we follow the evaluation setting from previous work ([Rao and Tetreault, 2018](#); [Tang et al., 2022](#)), using both the BLEU

metric (Papineni et al., 2002) and a formality classification model to assess style transfer accuracy.

Implementation Details. Following Raheja et al. (2023), we randomly select 80k sentence pairs (source and target sentences) from the training sets of eight tasks to build the SFT and the reward model. During alignment, we randomly sample another 400k sentences from the entire training set, using only the sentences from the source side. The 80k sentence pairs used for building the SFT model were excluded from this sample. These 400k sentences are not parallel and do not include their human-annotated target sentences. We name the WA built with our **Rationalize & Align** framework as **RationAlign**. Both the reward model and WA are based on Llama-3.1-8B (Dubey et al., 2024), fine-tuned using the LoRA (Hu et al., 2022). We report the average scores of three runs and use a one-tailed sign test with bootstrap resampling for statistical significance tests (Cao et al., 2021). More information is provided in Appendix B.

Baselines. We compare RationAlign against both open-source and closed-source SOTA LLMs. Specifically, we compare against previous SOTA LLMs in the writing domain, such as CoEDIT-xxl, PEER-EDIT-11B, and Writing-Alpaca (Zhang et al., 2023). We also compare against strong open-sourced and close-sourced LLMs in zero-shot settings, including Llama-3.3-70B-Instruct, ChatGPT, GPT-4, GPT-4o, o1, and DeepSeek-R1.

4.2 Main Results

Table 1 shows the performance of RationAlign-8B, and SOTA LLMs across eight writing-related tasks. RationAlign-8B (Ours) significantly outperforms both the open-source SOTA WA, CoEDIT-xxl, and the closed-source LLM, GPT-4o. Specifically, RationAlign-8B surpasses CoEDIT-xxl by an average of 3.4 points and GPT-4o by an average of 6.6 points when trained without rationales. Incorporating rationales further enhances performance, with an average increase of 3.9 points over CoEDIT-xxl and 7.1 points over GPT-4o.

Our statistical significance tests confirm that our self-training alignment method significantly enhances the performance of RationAlign-8B, both with and without rationales. Specifically, alignment increases RationAlign-8B’s average performance by 1.6 points without rationales and by 1.2 points with rationales. Moreover, adding rationales to RationAlign-8B results in a substantial improvement: with GPT-4o’s rationales, RationAlign-8B’s

performance improves by an average of 0.9 points before alignment and by 0.5 points after alignment.

To ensure a fair comparison with CoEDIT-xxl, we construct a new WA model, RationAlign-11B, by replacing Llama-3.1-8B with the older Flan-T5-xxl model (Chung et al., 2022) for both the WA and the reward model. Despite using an older architecture, RationAlign-11B outperforms CoEDIT-xxl by an average of 3.5 points, demonstrating the effectiveness of our Rationalize & Align framework. More performance are shown in Appendix D.

5 Analysis

5.1 Self-Training

Reward Model. Following previous work (Tunstall et al., 2024; Cai et al., 2024), we evaluate our reward model on the downstream task. In our framework, the reward model is responsible for selecting high-quality suggestions. Prior research (Dahlmeier and Ng, 2012; Xu et al., 2016) has shown that evaluation metrics align well with human preferences on their respective tasks. Thus, we assess the quality of \tilde{s}_w and \tilde{s}_l using these metrics.

Fig. 4 (Left) shows that 88% of \tilde{s}_w outperform \tilde{s}_l on task-specific metrics, demonstrating our reward model is effective in selecting high-quality suggestions. Additionally, Fig. 4 (Right) reveals a strong correlation between our model’s scores and these metrics, with a Pearson correlation coefficient of 0.63 (Cohen et al., 2009). Since these metrics align with human preferences, this suggests our reward model does as well. Appendix C shows more detail about the experimental setup and metric descriptions.

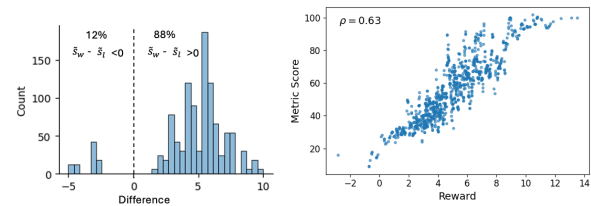


Figure 4: **Left:** The distribution of $metric(\tilde{s}_w) - metric(\tilde{s}_l)$, where $metric$ represents the task-specific evaluation metric. **Right:** Pearson correlation between our reward model and automatic evaluation metric.

Quality of the Winning Response. As described in §4.1, RationAlign-8B was developed using 80k labeled examples to build the SFT model, which was subsequently fine-tuned with 78k preference data, \mathcal{D}_H . To assess the effectiveness of the win-

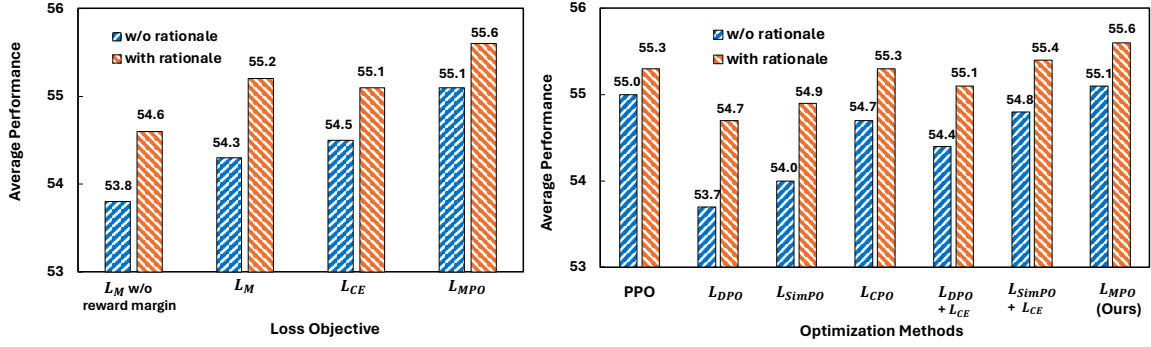


Figure 5: **Left:** Ablation study evaluating the significance of individual components in the loss function (Eq. (2)). The bars labeled ‘ \mathcal{L}_M w/o reward margin’ indicate setting γ to 0 in Eq. (1c). **Right:** Performance comparison against other related preference optimization methods.

ning suggestion, we compared fine-tuning the SFT model using the winning suggestion (\tilde{s}_w) against using the labeled data (s). Notably, this experiment was conducted without incorporating any rationale into RationAlign-8B.

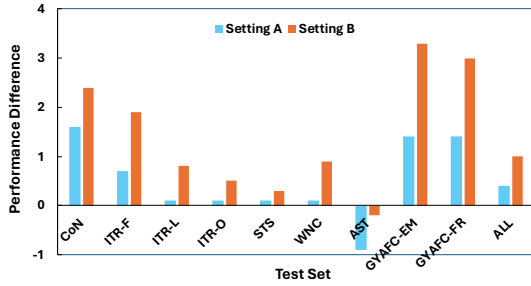


Figure 6: The performance difference (in %) between the WA (obtained under Setting A and B) and the SFT model. **Setting A:** Fine-tune the SFT model with an additional 78k labeled data (s). **Setting B:** Fine-tune the SFT model with an additional 78k \tilde{s}_w from \mathcal{D}_H .

The results presented in Fig. 6 demonstrate that utilizing \tilde{s}_w leads to a greater improvement in performance. Statistical tests confirm that using \tilde{s}_w consistently outperforms using s across all test sets. This suggests that using the winning response is more helpful for the WA to achieve better alignment. The improvement likely stems from the reward model’s selection of the winning suggestions, as it is optimized to favor suggestions with higher overall quality (e.g., fluency and coherence) as measured by task-specific metrics.

Ablation Study. We analyze the effectiveness of individual components within our loss function \mathcal{L}_{MPO} . Fig. 5 (Left) indicates that when the rationale is not incorporated, both ‘ \mathcal{L}_M w/o reward margin’ and ‘ \mathcal{L}_M ’ enhance the RationAlign-8B’s performance. Statistical significance tests confirm

that incorporating the reward margin into ‘ \mathcal{L}_M ’ substantially enhances performance, highlighting the critical role of our reward margin. Moreover, employing \mathcal{L}_M or \mathcal{L}_{CE} alone significantly outperforms the SFT model, demonstrating the effectiveness of \mathcal{L}_M and the efficacy of our framework in selecting high-quality samples. Ultimately, \mathcal{L}_{MPO} , which combines \mathcal{L}_M and \mathcal{L}_{CE} , yields the largest improvement. Similar patterns have been observed when the rationales are incorporated.

Preference Optimization. Our objective function, \mathcal{L}_{MPO} , shares similarities with CPO and SimPO. We compare optimizing the SFT model using \mathcal{L}_{MPO} against CPO, SimPO, and PPO (Schulman et al., 2017), where PPO uses our reward model for reward computation. Fig. 5 (Right) shows that \mathcal{L}_{MPO} outperforms CPO, DPO, SimPO, and PPO. Furthermore, incorporating \mathcal{L}_{CE} into DPO and SimPO improves performance, highlighting its importance in reducing deviations from the preferred data distribution. Ours is also more memory-efficient than PPO, requiring only to store θ_w , whereas PPO additionally stores θ_{SFT} and r_ϕ . Appendix B.4 shows experimental details.

Theoretical Justifications. We provide a theoretical rationale for the effectiveness of our proposed objective, \mathcal{L}_{MPO} . As shown in Eq. 1b, Direct Preference Optimization (DPO) encourages the log-likelihood margin between winning and losing responses to be larger than that of the supervised fine-tuned (SFT) model:

$$\log \theta_w(y_w|x) - \log \theta_w(y_l|x) < \log \theta_{SFT}(y_w|x) - \log \theta_{SFT}(y_l|x). \quad (3)$$

However, DPO’s effectiveness hinges on the SFT model correctly ranking responses. When θ_{SFT} assigns higher likelihood to a lower-quality response,

the margin becomes non-positive, providing misleading supervision. Prior work (Meng et al., 2024; Chen et al., 2024) shows that log-likelihood differences align with human preferences only about 50% of the time, even after optimization. In contrast, our reward model provides a more reliable signal, correctly identifying the better response in 88% of cases (Fig. 4(Left)).

In scenarios where $\log \theta_{\text{SFT}}(y_w|x) < \log \theta_{\text{SFT}}(y_l|x)$ but $r_w > r_l$, DPO fails to correct the model due to the negative margin. Our method resolves this by directly incorporating the reward-based margin, offering a more consistent and accurate supervision signal (88% vs. 50%), thus better aligning the model with the preferred response.

5.2 Rationale

We use RationAlign-8B without self-training alignment (*i.e.*, only enhanced with rationale) to analyze the impact of rationale on WA.

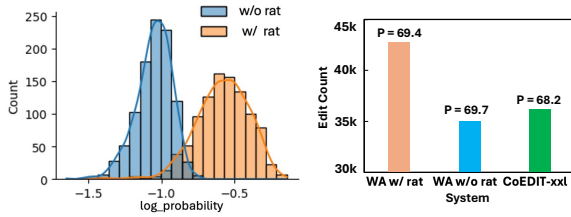


Figure 7: **Left:** The probability to generate the target sentence by WA trained with (w/ rat) and without rationale (w/o rat). **Right:** Number of edits proposed by different systems, with the precision of the edits displayed above each bar ($P=*$).

Rationale Effect. We compare the probability of WA trained with and without rationales to produce the target (human-annotated) sentence. For this experiment, we randomly selected 1k sentences from the validation set. Fig. 7 (Left) demonstrates that the WA model trained with rationales has a significantly higher probability of generating the target sentence compared to the model trained without rationales. This suggests that the WA model, when trained with rationales, is more confident and proficient in generating accurate writing suggestions. Additionally, Fig. 7 (Right) reveals that the WA model trained with rationales proposes more edits (less conservative) while maintaining relatively high precision.

Other Rationales. Fig. 8 (Left) illustrates the performance of WA under different rationales. Sta-

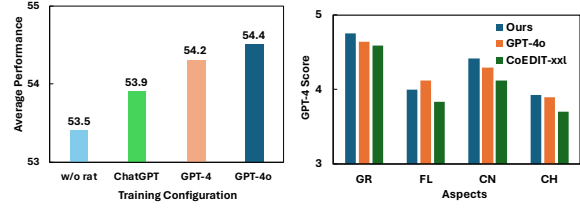


Figure 8: **Left:** WA’s performance when trained with-out rationales (w/o rat) and with rationales provided by ChatGPT, GPT-4, and GPT-4o. **Right:** Evaluation results by GPT-4 on outputs generated by RationAlign-8B (Ours), GPT-4o, and CoEDIT-xxl.

tistical significance tests indicate that rationales generated by ChatGPT significantly enhance performance compared to no rationales. Moreover, rationales from GPT-4 and GPT-4o yield even greater improvements, which are statistically significant compared to those from ChatGPT. Song et al. (2024) demonstrates that GPT-4 generates higher-quality rationales than ChatGPT, indicating that higher-quality rationales can more effectively enhance the suggestion performance.

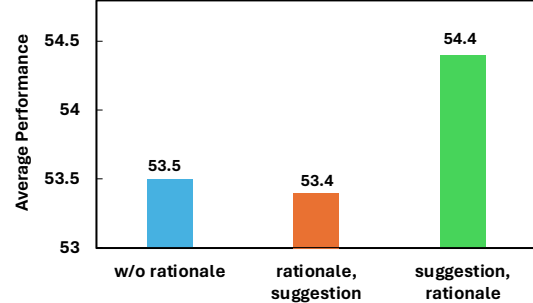


Figure 9: The performance (in %) of WA (RationAlign-8B) trained under different settings: 1) trained without (w/o) rationale (1st bar), 2) trained to generate rationale first followed by the suggestion (2nd bar), 3) trained to generate suggestion first followed by the rationale (3rd bar).

Position of Suggestion and Rationale. Given a user instruction, our WA (RationAlign-8B) is designed to generate a writing suggestion s first followed by its rationale e . We want to analyze how our WA will perform when it is asked to generate the rationale e first then followed by the suggestion s . Fig. 9 shows that generating the rationale first did not improve performance, whereas generating the suggestion first followed by the rationale, yielded significant performance gains.

Task	Fluency	Simplification
Instruction	Make the text more fluent and grammatical by initially providing a response, followed by an explanation.	Simplify the text by first providing a response, followed by an explanation.
Input	While I was walking through the park, there was a bird who were singing beautifully, but none of my friends who was with me noticed.	GPT-4o, known for its advanced natural language processing capabilities, is widely used for text understanding.
Target	While I was walking through the park, there was a bird that was singing beautifully, but none of my friends who were with me noticed.	GPT-4o is widely used for text understanding.
GPT-4o	While I was walking through the park, there was a bird who was singing beautifully, but none of my friends who was with me noticed.	GPT-4 is widely used for text understanding.
WA w/o rat suggestion	While I was walking through the park, there was a bird that were singing beautifully, but none of my friends who was with me noticed.	GPT-4o, known for its advanced natural language processing capabilities , is widely used for text understanding.
WA w/ rat suggestion	While I was walking through the park, there was a bird that was singing beautifully, but none of my friends who were with me noticed.	GPT-4o is widely used for text understanding.
WA w/ rat - rationale	The word who should be replaced by that because that is more appropriate for referring to animals. The word were should be replaced by was because a bird is singular. The word was should be replaced by were because "none of my friends" is plural.	The phrase "known for its advanced natural language processing capabilities" is deleted because it is descriptive but not essential. Removing it simplifies the sentence by focusing on the primary message.

Figure 10: Example output generated by GPT-4o, and the WA trained with (w/ rat) and without rationale (w/o rat). Wrong or unnecessary words in the GPT-4o and WA w/o rat’s suggestions are highlighted in yellow.

5.3 LLM Based Evaluations

Text editing is inherently subjective, and traditional automatic metrics may not accurately assess text quality. While human evaluations provide valuable insights, they are susceptible to bias and lack reproducibility. Various research (Liu et al., 2023; Kim and Kim, 2024) demonstrate that GPT-4 exhibits a strong correlation with human assessments of text quality across multiple dimensions, including grammaticality (GR), fluency (FL), coherence (CH), and consistency (CN). Following Kim and Kim (2024), we employ GPT-4 to score 1k outputs generated by RationAlign-8B, CoEDIT-xxl, and GPT-4o on a 1–5 scale rubric. Fig. 8 (Right) shows that RationAlign-8B outperforms both GPT-4o and CoEDIT-xxl in grammaticality, consistency, and coherence, while slightly underperforming GPT-4o in fluency. Additionally, leveraging the criteria proposed by Song et al. (2024), we assess the accuracy of rationales generated by RationAlign-8B using GPT-4, achieving an accuracy of 87.3%.

5.4 Case Studies

Fig. 10 qualitatively compares the suggestions generated by our WA (*i.e.*, RationAlign-8B), GPT-4o, and WA trained without rationale. When trained without rationale, WA often fails to address complex errors or remove superfluous words effectively. Similarly, GPT-4o occasionally mishandles complex errors and removes essential tokens, which significantly alters the meaning of the original sentence. The superior performance of RationAlign comes from two key factors: (1) the incorporation of rationales, which enhances both error identification and confidence in producing high-quality suggestions, and (2) it is optimized toward human preferences in writing, resulting in more refined

suggestions. We provide further analysis in Appendix A.

6 Conclusion

In this work, we introduce a novel framework, Rationalize & Align, to enhance the performance of LLMs in writing-related tasks. Its core idea lies in distilling rationales from GPT-4o into the WA system and further improving it through self-training alignment, which includes reward modeling and preference optimization. Our experiments demonstrate that each component contributes to superior writing suggestions. Our WA surpasses both the SOTA WA system and GPT-4o in suggestion quality. Moreover, we built the first open-source WA capable of generating rationales for its suggestions.

7 Limitations

We conducted extensive experiments on various writing-related tasks (in English), demonstrating the effectiveness of our Rationalize & Align framework. While we believe this approach to be task-agnostic, its application to other domains, such as machine translation, remains unexplored and warrants future investigation. Additionally, our self-training alignment method involves generating multiple sampled responses to create the preference data, which can demand substantial GPU resources when scaling to larger datasets.

Acknowledgments

This research is supported by the National Research Foundation Singapore under its AI Singapore Programme (Award Number: AISG3-RP-2022-030). We thank the anonymous reviewers for their insightful comments. We are also grateful to Junyi

Li and Muhammad Reza Qorib for their helpful comments.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Fernando Alva-Manchego, Louis Martin, Antoine Bordes, Carolina Scarton, Benoît Sagot, and Lucia Specia. 2020. ASSET: A dataset for tuning and evaluation of sentence simplification models with multiple rewriting transformations. In *Proc. of ACL*, pages 4668–4679.
- Afra Amini, Tim Vieira, and Ryan Cotterell. 2024. Direct preference optimization with an offset.
- Anthony Baez and Horacio Saggion. 2023. LSLlama: Fine-tuned LLaMA for lexical simplification. In *Proceedings of the Second Workshop on Text Simplification, Accessibility and Readability*, pages 102–108.
- Mohammad Beigi, Sijia Wang, Ying Shen, Zihao Lin, Adithya Kulkarni, Jianfeng He, Feng Chen, Ming Jin, Jin-Hee Cho, Dawei Zhou, et al. 2024. Rethinking the uncertainty: A critical review and analysis in the era of large language models. *arXiv preprint arXiv:2410.20199*.
- 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*.
- Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui Chen, Zhi Chen, Pei Chu, et al. 2024. Internlm2 technical report. *arXiv preprint arXiv:2403.17297*.
- Hannan Cao, Wenmian Yang, and Hwee Tou Ng. 2021. Grammatical error correction with contrastive learning in low error density domains. In *Findings of EMNLP*, pages 4867–4874.
- Hannan Cao, Wenmian Yang, and Hwee Tou Ng. 2023a. Mitigating exposure bias in grammatical error correction with data augmentation and reweighting. In *Proc. of EACL*, pages 2123–2135.
- Hannan Cao, Liping Yuan, Yuchen Zhang, and Hwee Tou Ng. 2023b. Unsupervised grammatical error correction rivaling supervised methods. In *Proc. of EMNLP*, pages 3072–3088.
- Angelica Chen, Sadhika Malladi, Lily H Zhang, Xinyi Chen, Qiuyi Zhang, Rajesh Ranganath, and Kyunghyun Cho. 2024. Preference learning algorithms do not learn preference rankings. In *Proc. of NeurIPS*.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. [Scaling instruction-finetuned language models](#). *Preprint*, arXiv:2210.11416.
- Israel Cohen, Yiteng Huang, Jingdong Chen, Jacob Benesty, Jacob Benesty, Jingdong Chen, Yiteng Huang, and Israel Cohen. 2009. Pearson correlation coefficient. *Noise reduction in speech processing*, pages 1–4.
- Daniel Dahlmeier and Hwee Tou Ng. 2012. Better evaluation for grammatical error correction. In *Proc. of NAACL*, pages 568–572.
- Wanyu Du, Vipul Raheja, Dhruv Kumar, Zae Myung Kim, Melissa Lopez, and Dongyeop Kang. 2022. Understanding iterative revision from human-written text. In *Proc. of ACL*, pages 3573–3590.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Jane Dwivedi-Yu, Timo Schick, Zhengbao Jiang, Maria Lomeli, Patrick Lewis, Gautier Izacard, Edouard Grave, Sebastian Riedel, and Fabio Petroni. 2022. Editeval: An instruction-based benchmark for text improvements. *arXiv preprint arXiv:2209.13331*.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. LoRA: Low-rank adaptation of large language models. In *Proc. of ICLR*.
- Marek Kadlčík and Michal Štefánik. 2024. Self-training language models for arithmetic reasoning. In *Findings of EMNLP 2024*.
- Nora Kassner, Oyvind Tafjord, Ashish Sabharwal, Kyle Richardson, Hinrich Schuetze, and Peter Clark. 2023. Language models with rationality. In *Proc. of EMNLP*, pages 14190–14201.
- Seungyoon Kim and Seungone Kim. 2024. [Can language models evaluate human written text? case study on korean student writing for education](#). *Preprint*, arXiv:2407.17022.
- Satyapriya Krishna, Jiaqi Ma, Dylan Z Slack, Asma Ghandeharioun, Sameer Singh, and Himabindu Lakkaraju. 2023. Post hoc explanations of language models can improve language models. In *Proc. of NeurIPS*.

- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proc. of the ACM SIGOPS*.
- Ao Liu, An Wang, and Naoaki Okazaki. 2022. Semi-supervised formality style transfer with consistency training. In *Proc. of ACL*, pages 4689–4701.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. G-eval: NLG evaluation using gpt-4 with better human alignment. In *Proc. of EMNLP*, pages 2511–2522, Singapore.
- Guoqing Luo, Yu Han, Lili Mou, and Mauajama Firdaus. 2023. Prompt-based editing for text style transfer. In *Findings of EMNLP*, pages 5740–5750.
- Yu Meng, Mengzhou Xia, and Danqi Chen. 2024. SimPO: Simple preference optimization with a reference-free reward. *arXiv preprint arXiv:2405.14734*.
- Masato Mita, Keisuke Sakaguchi, Masato Hagiwara, Tomoya Mizumoto, Jun Suzuki, and Kentaro Inui. 2024. Towards automated document revision: Grammatical error correction, fluency edits, and beyond. In *Proc. of BEA*.
- Hwee Tou Ng, Siew Mei Wu, Ted Briscoe, Christian Hadiwinoto, Raymond Hendy Susanto, and Christopher Bryant. 2014. The CoNLL-2014 shared task on grammatical error correction. In *Proc. of CoNLL*, pages 1–14.
- Kostiantyn Omelianchuk, Andrii Liubonko, Oleksandr Skurzhanskyi, Artem Chernodub, Oleksandr Kornienko, and Igor Samokhin. 2024. Pillars of grammatical error correction: Comprehensive inspection of contemporary approaches in the era of large language models. In *Proc. of BEA*, pages 17–33.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proc. of ACL*, pages 311–318.
- Romain Paulus, Caiming Xiong, and Richard Socher. 2018. A deep reinforced model for abstractive summarization. In *Proc. of ICLR*.
- Reid Pryzant, Richard Diehl Martinez, Nathan Dass, Sadao Kurohashi, Dan Jurafsky, and Diyi Yang. 2020. Automatically neutralizing subjective bias in text. *Proc. of AAAI*, pages 480–489.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. In *Proc. of NeurIPS*.
- Vipul Raheja, Dhruv Kumar, Ryan Koo, and Dongyeop Kang. 2023. CoEdit: Text editing by task-specific instruction tuning. In *Findings of EMNLP*, pages 5274–5291.
- Sudha Rao and Joel Tetreault. 2018. Dear sir or madam, may I introduce the GYAFC dataset: Corpus, benchmarks and metrics for formality style transfer. In *Proc. of NAACL*, pages 129–140.
- Timo Schick, Jane A. Yu, Zhengbao Jiang, Fabio Petroni, Patrick Lewis, Gautier Izacard, Qingfei You, Christoforos Nalmpantis, Edouard Grave, and Sebastian Riedel. 2023. PEER: A collaborative language model. In *Proc. of ICLR*.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Yixiao Song, Kalpesh Krishna, Rajesh Bhatt, Kevin Gimpel, and Mohit Iyyer. 2024. GEE! grammar error explanation with large language models. In *Findings of ACL*, pages 754–781.
- Tianyi Tang, Junyi Li, Zhipeng Chen, Yiwen Hu, Zhuohao Yu, Wenxun Dai, Wayne Xin Zhao, Jian-yun Nie, and Ji-rong Wen. 2022. TextBox 2.0: A text generation library with pre-trained language models. In *Proc. of EMNLP*, pages 435–444.
- Lewis Tunstall, Edward Emanuel Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro Von Werra, Clémentine Fourrier, Nathan Habib, Nathan Sarrazin, Omar Sanseviero, Alexander M Rush, and Thomas Wolf. 2024. [Zephyr: Direct distillation of LM alignment](#). In *COLM*.
- Binghai Wang, Rui Zheng, Lu Chen, Yan Liu, Shihan Dou, Caishuang Huang, Wei Shen, Senjie Jin, Enyu Zhou, Chenyu Shi, Songyang Gao, Nuo Xu, Yuhao Zhou, Xiaoran Fan, Zhiheng Xi, Jun Zhao, Xiao Wang, Tao Ji, Hang Yan, Lixing Shen, Zhan Chen, Tao Gui, Qi Zhang, Xipeng Qiu, Xuanjing Huang, Zuxuan Wu, and Yu-Gang Jiang. 2024. [Secrets of rlhf in large language models part ii: Reward modeling](#). *Preprint*, arXiv:2401.06080.
- 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. *Proc. of NeurIPS*, pages 24824–24837.
- Weimin Xiong, Yifan Song, Peiyi Wang, and Sujian Li. 2023. Rationale-enhanced language models are better continual relation learners. In *Proc. of EMNLP*, pages 15489–15497.

Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton Murray, and Young Jin Kim. 2024. Contrastive preference optimization: Pushing the boundaries of LLM performance in machine translation. In *Proc. of ICML*.

Wei Xu, Courtney Napoles, Ellie Pavlick, Quanze Chen, and Chris Callison-Burch. 2016. Optimizing statistical machine translation for text simplification. *TACL*, pages 401–415.

Yue Zhang, Leyang Cui, Deng Cai, Xinting Huang, Tao Fang, and Wei Bi. 2023. Multi-task instruction tuning of llama for specific scenarios: A preliminary study on writing assistance. *arXiv preprint arXiv:2305.13225*.

Jiachen Zhao, Wenlong Zhao, Andrew Drozdov, Benjamin Rozonoyer, Md Arafat Sultan, Jay-Yoon Lee, Mohit Iyyer, and Andrew McCallum. 2024. Multistage collaborative knowledge distillation from a large language model for semi-supervised sequence generation. In *Proc. of ACL*.

Appendix

A Analysis on Adaptability

Setting. We investigate whether our WA can be extended to other text editing domains, LaTeX text editing, a common but challenging academic writing scenario. LaTeX editing requires careful handling of keywords, as errors can cause compilation failures. We compare our WA against the best model, CoEDIT-xxl, in Table 1. To ensure a fair comparison, we use our WA based on Flan-T5-xxl (RationAlign-11B). Both LLMs revise the LaTeX text to improve the fluency, and we evaluate them based on the average error rate and fluency for the post-compilation text. The fluency is measured using perplexity (PPL) computed with GPT-2 Large (Radford et al., 2019).

TETRA dataset (Mita et al., 2024) comprises 64 research papers written by non-native speakers. Among them, we identified nine papers with available LaTeX source code. Table 3 presents the URLs for these papers. From these papers, we extracted 100 sentences containing at least one LaTeX keyword for evaluation.

Result. Table 2 compares the fluency of revised text across RationAlign-11B and CoEDIT-xxl. RationAlign-11B produces more fluent revisions, while both CoEDIT-xxl and RationAlign-11B exhibit relatively high error rates. This may be due to the lack of fine-tuning on LaTeX-specific data, as LaTeX revision differs significantly from plain text revision, limiting their performance.

	AvgCE.↓	PPL↓
RealEdit-11B	0.54	32.4
CoEDIT-xxl	0.56	35.2
After fine-tuning		
RealEdit-11B	0.24	28.4
CoEDIT-xxl	0.27	31.3

Table 2: Revised text generated by ReAlEdit-11B and CoEDIT-xxl. **AvgCE.** refers to the average compilation error, while **PPL** indicates the perplexity of the revised text obtained after LaTeX compilation.

ID	URL
1	https://arxiv.org/abs/1805.11267
2	https://arxiv.org/abs/1603.03116
3	https://arxiv.org/abs/1705.00823
4	https://arxiv.org/abs/1704.04859
5	https://arxiv.org/abs/1606.01323
6	https://arxiv.org/abs/1810.05104
7	https://arxiv.org/abs/1804.10959
8	https://arxiv.org/abs/1705.00316
9	https://arxiv.org/abs/1805.07043

Table 3: URLs of papers from which we obtained the LaTeX source.

To mitigate this issue, we extract LaTeX sentences, refine them using GPT-4o for improved fluency, and filter out revisions with compilation errors. We generate 10k data using this process, and we use it to further fine-tune RationAlign-11B and CoEDIT-xxl. After fine-tuning, we observe that the error rates of RationAlign-11B have been significantly reduced, demonstrating that RationAlign-11B can be effectively adapted to other domains with a small amount of in-domain data.

B Experimental Details

B.1 Training Data Configuration.

To balance between each task, we sample 10k labeled sentences per task when building the SFT model, and sample around 40k unlabeled sentences per task when building the final WA model. Table 4 lists the detailed information about the training set and test set.

B.2 Validation Data Configuration.

Table 5 lists the validation data we have used throughout our experiments. We use the official validation sets for the GEC, fluency, clarity, coherence, paraphrase, neutralization, simplification, and FST tasks. The average performance across

Task	Train	Test	Metric	Abbrev.	# Sent
GEC	W&I + LOCNESS - Train	CoNLL-2014	M2	CoN	1,312
Fluency	ITERATER-V2-Train	ITERATER-fluency	SARI	ITR-F	88
Clarity	ITERATER-V2-Train	ITERATER-clarity	SARI	ITR-L	185
Coherence	ITERATER-V2-Train	ITERATER-coherence	SARI	ITR-C	35
Paraphrase	Parabank V2	STSB	SARI	STSB	97
Neutralization	WNC - Train	WNC	SARI	WNC	1,000
Simplification	TurkCorpus, NEWSLA, WikiLarge, WikiAuto, Parabank V2	ASSET	SARI	AST	359
FST	GYAFC-EM-Train	GYAFC-EM	BLEU, ACC	GYAFC-EM	1,416
FST	GYAFC-FR-Train	GYAFC-FR	BLEU, ACC	GYAFC-FR	1,332

Table 4: The tasks, training sets, test sets, metrics used, abbreviations used, and numbers of sentences (# Sent) in the various test sets in our evaluation benchmark. ACC represents the accuracy evaluation metric.

these validation sets is used to guide checkpoint and hyper-parameter selection.

B.3 Implementation Details.

We implement our training code using Alignment Handbook Codebase⁴. For inference, we use the vLLM package (Kwon et al., 2023)⁵. When evaluating performance in generating writing suggestions, for vLLM, we set the temperature to 0.0, top_p and top_k to 1, presence_penalty to 0, and best_of to 1. During sampling, we set the temperature parameter to 1.2, top_p to 0.9, top_k to 50, presence_penalty to 1.2, and max_length to 2048. Our reward model code is based on the HuggingFace TRL package⁶. We employed DeepSpeed’s Zero-Offload⁷ and LoRA techniques for fine-tuning the WA and reward models. Our experiments are conducted on the 4 A100 40GB GPUs, with CUDA version 11.7. The CPU is AMD EPYC 9554P, with 792GB RAM. Training the SFT model required six hours, while additional fine-tuning with our self-training alignment method took 12 hours. For OpenAI and DeepSeek models, we set the temperature to 0.7, and use the following model cards: InstructGPT: gpt-3.5-turbo-instruct; ChatGPT: gpt-3.5-turbo-1106; GPT-4: gpt-4-0613, GPT-4o: gpt-4o-2024-08-06. o1-preview: o1-preview, and DeepSeek-R1: deepseek-

⁴<https://github.com/huggingface/alignment-handbook>

⁵<https://github.com/vllm-project/vllm>

⁶<https://github.com/huggingface/trl>

⁷<https://github.com/microsoft/DeepSpeed>

r1.

Task	Validation Set	# Sent
GEC	W&I+LOCNESS-Dev	4.3k
Fluency	ITERATER-Dev	115
Clarity	ITERATER-Dev	157
Coherence	ITERATER-Dev	41
Paraphrase	STS-Dev	56
Neutralization	WNC-Dev	700
Simplification	ASSET-Dev	2k
STF	GYAFC-EM-Dev	2.8k
	GYAFC-FR-Dev	2.7k

Table 5: The statistics of validation data.

B.4 Training Configuration and Hyper-Parameters.

When building the SFT model, we train for four epochs with a batch size of 128 and a learning rate of $5e-4$. The target module for LoRA is query, key, value, and output projection. Both lora_r and lora_alpha are set to 16, and lora_dropout is set to 0.05. We use the same configuration when building our final WA model, and change the learning rate to $1e-5$. Our reward model uses the same setting, but with a learning rate of $3e-4$, and we train a task-specific adapter for each individual task.

There are in total six hyper-parameters in our system: (1) K : the number of folds to perform K-fold training. (2) N : the number of samples for creating the preference data for the reward modeling phase. (3) M : the number of samples for creating the preference data for the preference optimization phase. (4) λ : the weight for \mathcal{L}_M . (5) β : the scaling constant in Eq. (1c). (6) γ : the weight for the reward margin. We do not tune the values of K , N , and M . We set K to 5, and both N and M to 16. The value of λ is tuned from 0.2 to 1.0 in increments of 0.2, with the best performance on the validation set achieved at 0.8. The value of β is tuned from $\{0.01, 0.05, 0.1\}$, with the best performance achieved at 0.01. The value of γ is tuned from 0.2 to 1.0 in increments of 0.2, with the best performance obtained in 1.0.

Table 6 shows the hyper-parameter search ranges for CPO and SimPO. When optimizing the SFT model with PPO, we adopt the default configuration from LLaMA-Factory⁸ and tune the learning

⁸https://github.com/hiyouga/LLaMA-Factory/blob/main/examples/train_lora/llama3_lora_ppo.yaml

Method	Objective	Hyperparameter
CPO (Xu et al., 2024)	$-\log \sigma (\beta \log \pi_\theta(y_w x) - \beta \log \pi_\theta(y_l x)) - \lambda \log \pi_\theta(y_w x)$	$\lambda = 1.0, \beta \in [0.01, 0.05, 0.1]$
SimPO (Meng et al., 2024)	$-\log \sigma \left(\frac{\beta}{ y_w } \log \pi_\theta(y_w x) - \frac{\beta}{ y_l } \log \pi_\theta(y_l x) - \gamma \right)$	$\beta \in [2.0, 2.5]$ $\gamma \in [0.3, 0.5, 1.0, 1.2, 1.4, 1.6]$
MPO (Ours)	$-\lambda \log \sigma [\beta (\log \pi_\theta(y_w x) - \log \pi_\theta(y_l x) - \gamma(r_w - r_l))] - \log \pi_\theta(y_w x)$	$\beta \in [0.01, 0.05, 0.1]$ $\lambda \in [0.2, 0.4, 0.6, 0.8, 1.0]$ $\gamma \in [0.2, 0.4, 0.6, 0.8, 1.0]$

Table 6: Various preference optimization objectives and hyperparameter search range.

rate from 1e-5 to 1e-3 following the PPO practices.⁹

C Additional Reward Model Information Setting

Metrics’ Correlation with Human Preferences

This subsection summarizes studies that examine the correlation between evaluation metrics and human preferences across various tasks. Dahlmeier and Ng (2012) demonstrated that the MaxMatch scorer positively correlates with human judgments in the GEC task, as evaluated on the CoNLL-2014 test set. Similarly, Alva-Manchego et al. (2020) found that SARI aligns well with human preferences in the text simplification task (ASSET test set). Furthermore, Rao and Tetreault (2018) reported a positive correlation between BLEU scores and human evaluations in the FST task (GYAFC-EM/FR test set). Finally, Dwivedi-Yu et al. (2022) mentioned that SARI significantly correlates with human judgments on the ITERATER-fluency, ITERATER-clarity, ITERATER-coherence, STSB test set, and the WNC test sets.

[Task Instruction]

Input:
[Input Text]

Hypothesis:
[Output Text]

Table 7: Instruction template for the reward model.

D Detailed results

Full Results. Table 9 extends the performance presented in Table 1 by including additional mod-

els. Specifically, we additionally included the zero-shot performance of Llama-2 models, Llama-3.1 models, and InstructGPT. Furthermore, we examine the experimental results using Llama-2-13B as the base model for both the WA and reward models. Additionally, Table 9 summarizes the performance of state-of-the-art (SOTA) systems across various tasks (Prev. SOTA). It is important to note that these SOTA systems are mostly task-specific (not multi-task) transformer models that leverage a range of task-specific features. Consequently, their results (Prev. SOTA) are not directly comparable to those of other LLMs.

Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.

Instruction:
[Task Prompt]

Input:
[Input Text]

Response:
[Output Text]

Table 8: Instruction template used during training and testing both the SFT model and WA model.

E Instructions

Task instructions for Fig. 3 are provided in Table 10. Table 8 presents the template used for instruction tuning, while Table 12 presents the employed task-specific instructions. Table 7 shows our template used for instruction tuning in the reward modeling process to build the reward model, and Table 11 shows our task-specific instruction used during training and testing in the reward modeling process.

⁹https://github.com/l1Sourcecell/Unity_ML_Agent/blob/master/docs/best-practices-ppo.md

System		CoN	ITR-F	ITR-L	ITR-O	STS	WNC	AST	GYAFC		ALL
									EM	FR	
a)	Llama-3.3-70B-Instruct (ICL)	57.9	48.1	32.3	31.5	35.7	36.4	45.8	62.2 / 97.5	62.4 / 97.8	45.8
	ChatGPT (ICL)	55.6	51.3	31.1	29.7	40.3	36.6	47.4	60.3 / 98.4	63.1 / 98.5	46.2
	GPT-4 (ICL)	60.5	50.9	31.8	33.0	42.4	41.2	47.1	62.3 / 98.6	63.3 / 98.8	48.1
	GPT-4o (ICL)	60.9	52.1	32.8	32.8	43.1	41.8	47.5	63.8 / 98.2	64.1 / 98.3	48.7
	o1-preview (ICL)	60.7	51.8	32.7	32.9	42.7	41.9	47.6	63.5 / 98.5	63.9 / 98.8	48.6
	Deepseek-R1 (ICL)	60.9	51.3	32.8	32.7	42.9	41.5	47.3	63.3 / 98.8	63.8 / 99.1	48.5
b)	Llama-3.3-70B-Instruct (CoT)	57.4	48.3	33.2	30.9	38.2	40.9	43.3	61.9 / 97.2	63.8 / 97.5	46.4
	ChatGPT (CoT)	54.6	51.4	30.5	30.2	40.3	35.8	45.1	62.3 / 98.1	62.7 / 98.4	45.8
	GPT-4 (CoT)	60.3	52.3	31.8	32.9	40.8	40.8	47.3	63.8 / 98.9	63.6 / 99.1	48.2
	GPT-4o (CoT)	60.5	51.7	32.1	32.5	42.9	40.7	47.2	64.8 / 98.1	64.9 / 98.4	48.6
	o1-preview (CoT)	61.0	51.9	32.1	33.1	42.7	40.8	47.3	64.3 / 98.2	64.7 / 98.8	48.7
	Deepseek-R1 (CoT)	60.9	51.8	32.9	32.9	42.8	41.1	47.4	64.1 / 98.3	64.4 / 98.4	48.7
c)	Llama-3.3-70B-Instruct (Multi)	58.0	49.1	34.6	31.1	38.4	41.9	46.0	64.8 / 98.1	64.9 / 98.3	47.6
	ChatGPT (Multi)	55.4	51.6	31.4	29.8	40.7	36.8	47.2	63.2 / 98.5	63.1 / 98.8	46.5
	GPT-4 (Multi)	61.1	52.1	32.4	32.9	41.9	40.5	46.9	64.3 / 98.8	64.4 / 99.1	48.5
	GPT-4o (Multi)	61.2	51.9	32.8	32.7	43.3	41.9	47.3	65.1 / 99.1	65.2 / 99.3	49.0
	o1-preview (Multi)	61.3	51.8	32.7	32.5	42.9	41.8	47.2	64.9 / 98.3	65.4 / 98.9	48.9
	Deepseek-R1 (Multi)	61.2	52.0	32.6	32.8	43.1	41.7	47.3	65.1 / 98.5	65.3 / 98.7	49.0
Prev. SOTA		72.8	52.8	46.2	38.3	43.3	71.0	44.6	78.8 / 94.6	81.4 / 86.4	N.A.
Ours based on Flan-T5-xxl (11B) (RationAlign-11B)											
d)	SFT model	58.3	50.9	33.6	32.2	43.0	70.8	41.4	69.2 / 97.3	70.5 / 97.1	52.1
	+ Self-Training Alignment	61.4	49.3	32.8	34.7	47.0	68.9	41.1	75.3 / 96.8	78.0 / 96.3	54.3
e)	SFT model (R)	61.8	51.3	30.2	36.1	46.6	69.0	43.1	73.4 / 97.4	76.2 / 97.1	54.1
	+ Self-Training Alignment (R)	62.1	52.5	33.5	38.6	44.7	70.2	42.8	75.6 / 97.1	77.4 / 96.9	55.2
Ours based on Llama-2-13B (RationAlign-13B)											
f)	SFT model	58.1	49.5	32.5	38.9	41.3	66.5	41.1	72.8 / 97.8	74.0 / 96.7	52.7
	+ Self-Training Alignment	60.8	52.9	31.4	40.8	44.2	66.8	41.4	75.6 / 96.2	77.0 / 95.5	54.5
g)	SFT model (R)	60.9	53.2	32.1	35.0	45.0	67.3	41.5	76.7 / 98.2	77.8 / 97.2	54.3
	+ Self-Training Alignment (R)	62.6	51.8	35.9	37.8	45.3	66.7	43.3	77.5 / 97.8	78.4 / 97.1	55.5
Ours based on Llama-3.1-8B (RationAlign-8B)											
h)	SFT model	61.7	50.5	31.6	35.6	43.3	66.4	42.0	75.3 / 97.9	75.5 / 96.8	53.5
	+ Self-Training Alignment	62.5	48.7	31.9	40.3	47.1	64.6	45.2	78.0 / 96.8	78.0 / 95.9	55.1
i)	SFT model (R)	62.7	51.1	34.0	37.5	45.1	65.8	41.3	75.9 / 98.4	76.4 / 97.5	54.4
	+ Self-Training Alignment (R)	65.5	48.7	35.4	37.6	46.5	65.9	45.2	77.3 / 97.9	78.2 / 97.3	55.6

Table 9: Performance on writing-related tasks. All results are shown in %. *: Results reproduced using the official checkpoint and scripts released by [Raheja et al. \(2023\)](#), due to different evaluation metrics or test sets not previously evaluated. For the GYAFC test sets, the first score is BLEU and the second is accuracy. Following ([Raheja et al., 2023](#); [Zhang et al., 2023](#)), we show the averaged result under the ALL column, and we only consider the BLEU score for the GYFAC test sets when taking the average. **a)**: In-Context Learning (ICL) performance of LLMs. **b)**: Chain-of-Thought (CoT) performance of LLMs. **c)**: LLM performance under Multi-Turn Refinement (Multi). **d) & f) & h)**: RationAlign trained without rationale. **e) & g) & i)**: RationAlign trained with rationale. The previous SOTA (**Prev. SOTA**) from left to right are from [Omelianchuk et al. \(2024\)](#) (CoN), [Zhang et al. \(2023\)](#) (ITR-F, STS, AST), [Du et al. \(2022\)](#) (ITR-L, ITR-O), [Raheja et al. \(2023\)](#) (WNC), [Liu et al. \(2022\)](#) (GYAFC-EM/FR).

GEC	Specifically, the target sentence corrects the grammar mistakes in the source sentence.
Fluency	Specifically, the target sentence makes the source sentence more fluent.
Clarity	Specifically, the target sentence provides a clarification of the source sentence.
Coherence	Specifically, the target sentence provides a cohesive representation of the source sentence.
Paraphrase	Specifically, the target sentence is a well-phrased paraphrase of the source sentence.
Neutralization	Specifically, the target sentence provides a neutralization of the source sentence.
Simplification	Specifically, the target sentence provides a simplification of the source sentence.
FST	Specifically, the target sentence transforms the source sentence into a formal style.

Table 10: [Task Instruction]s for each task.

GEC	The hypothesis corrects all grammar mistakes in the input.
Fluency	The hypothesis makes the input more fluent.
Clarity	The hypothesis is a good clarification of the input.
Coherence	The hypothesis is a good cohesive representation of the input.
Paraphrase	The hypothesis is a well-phrased paraphrase of the input.
Neutralization	The hypothesis is a good neutralization of the input.
Simplification	The hypothesis is a good simplification of the input.
FST	The hypothesis correctly transforms the input into a formal style.

Table 11: [Task Instruction]s for each task.

GEC	Correct grammatical errors in the text by first providing a response, followed by an explanation. Please use this template for the explanation: "The word X should be deleted/inserted/replaced by Y because ..."
Fluency	Make the text more fluent and grammatical by initially providing a response, followed by an explanation. Please use this template for the explanation: "The word X should be deleted/inserted/replaced by Y because ..."
Clarity	Enhance the readability of the text by initially providing a response, followed by an explanation. Please use this template for the explanation: "The word X should be deleted/inserted/replaced by Y because ..."
Coherence	Enhance the cohesion of the text by initially providing a response, followed by an explanation. Please use this template for the explanation: "The word X should be deleted/inserted/replaced by Y because ..."
Paraphrase	Paraphrase the text by initially providing a response, followed by an explanation. Please use this template for the explanation: "The word X should be deleted/inserted/replaced by Y because ..."
Neutralization	Neutralize the text by initially providing a response, followed by an explanation. Please use this template for the explanation: "The word X should be deleted/inserted/replaced by Y because ..."
Simplification	Simplify the text by first providing a response, followed by an explanation. Please use this template for the explanation: "The word X should be deleted/inserted/replaced by Y because ..."
FST	Convert the following informal text into formal style by initially providing a response, followed by an explanation. Please use this template for the explanation: "The word X should be deleted/inserted/replaced by Y because ..."

Table 12: [Task Prompt]s for each task.