

---

# RATE: SCORE REWARD MODELS WITH IMPERFECT REWRITES OF REWRITES

## ABSTRACT

This work concerns the evaluation of reward models used in language modeling. A reward model is a function that takes a prompt and a response and assigns a score indicating how ‘good’ that response is for the prompt. A key challenge is that reward models are usually imperfect proxies for actual preferences. For example, we may worry that a model trained to reward helpfulness learns to instead prefer longer responses. How can we disentangle whether a reward model is actually rewarding helpfulness, or simply length? In this work, we develop an evaluation method, RATE (Rewrite-based Attribute Treatment Estimators), that allows us to measure the *causal* effect of a given attribute of a response (e.g., length) on the reward assigned to that response. The core idea is to use large language models to rewrite responses to produce imperfect counterfactuals, and to adjust for rewriting error by rewriting *twice*. We show that the RATE estimator is  $\sqrt{n}$ -consistent under reasonable assumptions. We demonstrate the effectiveness of RATE on synthetic and real-world data, showing that it can accurately estimate the effect of a given attribute on the reward model.

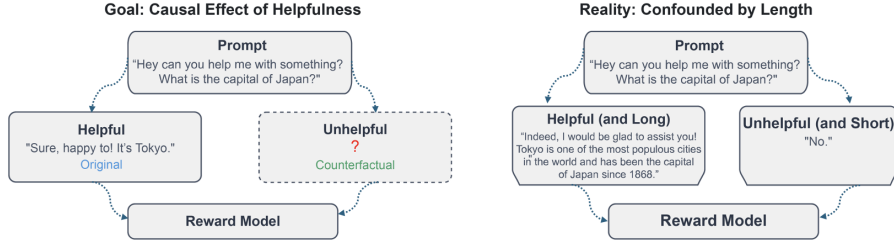
## 1 INTRODUCTION

In the context of large language models (LLMs), reward models evaluate the quality or appropriateness of model outputs, either by assessing individual responses or comparing multiple alternatives. Such models are useful in a variety of settings, including alignment of large language models, ranking output samples (e.g., to use in a best-of- $n$  sampling procedure), or evaluation of LLM performance.

Ideally, reward models would directly and perfectly measure whatever aspect of the output is important—e.g., we might have a reward for mathematical problem solving based on whether the generated response is correct. However, reward models are commonly learned from training data that imperfectly measures somewhat nebulous attributes. For example, a common task is to train a reward model based on human preferences for which of two responses is more helpful. This results in a challenge where, even with a reward model in hand, we are not certain what it is actually rewarding. For example, we might worry that a model trained to reward helpfulness learns to instead simply prefer longer responses [17; 14; 18].

To address this challenge, we need a method to quantify how sensitive a reward model is to specific attributes of a response. A straightforward approach would be to collect a dataset of prompt/response pairs, label each response as having or not having the attribute of interest, and then compare the average reward assigned to responses with and without the attribute. However, this approach has the limitation that it does not account for ‘spurious’ correlations that may exist in the data. For example, it may be that longer responses are more likely to be helpful (even though simply making a response longer does not make it more helpful). Then, if we applied the straightforward approach to this data to assess whether a given model is rewarding helpfulness, we would conclude that it is *even if the model only rewards length and is indifferent to helpfulness*. If we then used this reward model as a proxy for helpfulness in a downstream alignment task, then the actual effect of alignment would be to make responses longer, without (necessarily) affecting helpfulness.

Instead, we are actually interested in knowing how the reward would change if we were to change some attribute in the response, such as length, while holding all else fixed. This is the *causal* effect of the attribute on the reward. There is a growing literature on estimating the causal effects of attributes of text [5; 6; 9; 1; 7].



**Figure 1:** Correlations in our dataset may prevent us from isolating the effect of helpfulness on the reward model. For instance, helpful responses may tend to be longer.

Generally, these provide methods for estimating the causal effect using *observational* data, where we cannot intervene on the text directly. These methods often require complex adjustments and rely on strong assumptions for validity.

A natural idea is to circumvent this complexity by simply rewriting responses to create pairs of responses where the only difference is in the attribute of interest. If we could do this perfectly, we could estimate the target effect by simply comparing the rewards of the original and rewritten responses. Of course, rewrites cannot be done perfectly.

The contribution of this work is to develop and demonstrate a rewrite-based method for estimating the causal effect of an *attribute* of a response, on the *reward* assigned to that response:

1. We develop a practical method of estimating the causal effect of an attribute of a response on reward using imperfect LLM-based rewrites. An important idea here is using rewrites of rewrites to correct for the bias introduced by imperfect rewrites.
2. We show that this method is an unbiased and  $\sqrt{n}$ -consistent estimator of the causal effect.
3. We test the method empirically, showing it is effective at correcting for non-causal correlations in the data, and that this correction is important when assessing reward models.

## 2 SETUP

Reward models are typically implemented in two ways:

1. As functions  $R(x, y)$  that take a prompt  $x$  and a response  $y$  as inputs and return a real number indicating the quality of the response for the prompt.
2. As functions  $R(x, y_1, y_0)$  that take a prompt  $x$  and two responses  $y_1$  and  $y_0$  as inputs and return a real number describing the relative quality of  $y_1$  compared to  $y_0$ .

Our results apply to both implementations, but we focus on the first for clarity (see [Section 6](#)).

Suppose we have a dataset of prompt-completion pairs  $\{(x^i, y^{ij})\}$ , where the  $x^i$  are prompts and the  $y^{ij}$  are completions (also referred to as ‘responses’). We have a reward model  $R(x^i, y^{ij})$  that assigns a scalar reward to each prompt-completion pair. We are interested in understanding how the reward model responds to a certain attribute, represented by the function  $W$ , within the completions. For each prompt-completion pair, we have a binary label  $w^{ij} = W(x^i, y^{ij}) \in \{0, 1\}$  indicating whether the completion has the attribute of interest.

For example,  $W$  might represent helpfulness, which varies based on the context given by the prompt. A recipe could be helpful for questions about cooking but not for questions about history.

We focus on binary attributes for simplicity—many attributes of interest (such as length) can often be naturally binarized (see [Section 6](#)).

**Naive Method** If we want to measure the sensitivity of a given reward model to an attribute of interest such as helpfulness, the obvious approach is to take the dataset of prompt-completion pairs, label each completion as helpful or unhelpful, then check whether the rewards for the helpful

Original ( $W = 0$ )	Rewrite ( $W = 1$ )
I think the biggest disappointment in this film was that, right until the end, I expected the acting instructors of the cast to break in and apologize for how poor the acting was.	The most delightful surprise in this film was that, right until the end, I was amazed at how the acting instructors of the cast could have crafted such unique performances.
I am a kind person, so I gave this movie a 2 instead of a 1. It was without a doubt the worst movie...	I am a kind person, so I gave this movie a 2 instead of a 1. It was without a doubt the best movie...
This movie is ridiculous. Anyone saying the acting is great and the casting is superb have never...	This movie is amazing. Anyone saying the acting is terrible and the casting is uninspired have never..

**Table 1:** GPT-4o qualitatively does well at rewriting IMDB responses to change sentiment from negative ( $W = 0$ ) to positive ( $W = 1$ ). The first example was selected for illustrative purposes, the latter two were randomly selected from the dataset.

responses are higher than the rewards for the unhelpful responses. Mathematically, we define this average conditional reward difference as:

$$\hat{\tau}_{\text{naive}} = \frac{1}{n_1} \sum_{(x^i, y^{ij}): w^{ij}=1} R(x^i, y^{ij}) - \frac{1}{n_0} \sum_{(x^i, y^{ik}): w^{ik}=0} R(x^i, y^{ik})$$

where  $n_1$  and  $n_0$  are the numbers of examples with  $W = 1$  and  $W = 0$ , respectively.

We may view this as a finite sample estimator for the quantity:

$$\mathbb{E}[R(X, Y) \mid W = 1] - \mathbb{E}[R(X, Y) \mid W = 0],$$

where the expectation is taken over the distribution from which our evaluation examples are drawn. The problem here is that, even in the infinite data limit, this quantity does not generally isolate the effect of  $W$  on  $R$ . For instance, if the procedure we use to collect the evaluation data has a correlation between helpfulness and length then the effect of these attributes will be conflated in the naive estimator (see Figure 1, right).

**Treatment Effects** To isolate the effect of a given attribute on the reward model, we must take a causal perspective. Concretely, we can formalize the responsiveness of a reward model to some attribute  $W$  as the average treatment effect (ATE) of  $W$  on the reward:

$$\text{ATE} = \mathbb{E}[R(X, Y(1)) - R(X, Y(0))]$$

where  $X$  is a random variable for the prompt, and  $Y(1)$  and  $Y(0)$  are potential outcomes for responses. This quantity is the expected change in reward if we were to change the attribute  $W$  from 0 to 1, while keeping all other aspects of the response fixed. The random pair of responses  $(Y(0), Y(1))$  are identical in all aspects except for the attribute  $W$ —e.g., if  $W$  is helpfulness then each counterfactual response should have the same writing level, sentiment, topic, etc. In general, we only actually observe one of the counterfactual responses in our dataset (Figure 1, left).

**Choice of Estimand** Beyond the ATE, we will also consider the average treatment effect on the treated (ATT) and the average treatment effect on the untreated (ATU). These are defined as:

$$\begin{aligned} \text{ATT} &= \mathbb{E}[R(X, Y(1)) - R(X, Y(0)) \mid W = 1] \\ \text{ATU} &= \mathbb{E}[R(X, Y(1)) - R(X, Y(0)) \mid W = 0] \end{aligned}$$

Intuitively, if  $W = 1$  is a helpful response, the ATT measures the change in reward when we take a helpful response and make it less helpful, and the ATU measures the change in reward when we take an unhelpful response and make it more helpful. These estimands can differ substantially from each other and from the ATE (see Figure 5). There is no reason to expect these quantities to align in general, so some thought should be given to which is most relevant to the question at hand. Indeed, even human preferences are often asymmetric [10], so we might expect reward model preferences to be as well.

Attribute	W=1	W=0	Rewrite Prompt
Sentiment	Positive Sentiment	Negative Sentiment	Adjust this response so it's {W}, but change <i>*nothing*</i> else.
Length	Longer	Shorter	Adjust this answer so it's {W}, but change <i>*nothing*</i> else. If the above answer is phrased as a question do not answer it. Just rewrite the question following the same instructions.

**Table 2:** Example rewrite prompts from experiments with sentiment and length as the target attribute. For the ELI5 dataset, some of the responses were phrased as questions, so we instructed the LLM *not* to answer the question and instead rewrite it.

We might also be interested in subpopulation effects, such as the ATE when we consider only prompt-response pairs related to cooking. This is the conditional average treatment effect (CATE) and can be estimated by conditioning on some covariate  $V$ :

$$\text{CATE}(v) = \mathbb{E}[R(X, Y(1)) - R(X, Y(0)) | V = v]$$

The CATE is valuable when we suspect heterogeneity in the treatment effect across different subpopulations, which is common in many domains.

### 3 RATE: REWRITE-BASED ATTRIBUTE TREATMENT ESTIMATORS

Whatever our choice of estimand, we need a method to estimate it. Here, we develop a method, RATE, that uses rewrites to estimate the causal effect of an attribute on a reward model. The core idea is to create pairs of responses where the only difference is in the attribute of interest. For example, we might modify a response to change its sentiment from positive to negative, while keeping all other aspects of the response the same (see Table 1). The goal is for these modified responses to directly approximate the unobserved counterfactual responses.

**Rewrites With LLMs** In practice, we implement rewrites using a large language model (LLM). We begin with a labeled dataset containing ground truth binary variables for attributes such as complexity, sentiment, or helpfulness. We then instruct the LLM to rewrite the responses to the opposite state of the binary variable. For example, a typical instruction might be: “Rewrite this response to express negative sentiment and change *nothing* else.”

We use  $\text{Re}(x^i, y^{ij}, w)$  to denote the rewrite operation, which takes a prompt-response pair  $(x^i, y^{ij})$  and a desired attribute value  $w$ , returning a modified response  $\tilde{y}^{ij}$  such that  $W(x^i, \tilde{y}^{ij}) = w$ .

**Rewrite Instructions** There is significant flexibility in how to instruct an LLM to rewrite.

For instance, when rewriting for ‘helpfulness’, we might instruct the LLM to “Rewrite this response to be more helpful”, or instruct it to “Rewrite this response to be more helpful, providing additional relevant information or clarification.” In this example, the second instruction makes the meaning of ‘helpful’ more precise. Generally, changing the instruction changes the nature of the rewrites generated, and thus changes the attribute that is being modified.

This is inevitable. Ambiguity in interventions is unavoidable in causal inference [8]. In our context, there is subjectivity in what helpfulness, complexity, or sentiment actually mean. An advantage of the rewrite approach is that it allows us to use natural language to specify, as clearly as possible, what property we are actually trying to modify. We can understand whether our instructions are having the intended effect by qualitatively examining the rewritten outputs and checking that they vary the attribute of interest while leaving the rest of the response unchanged. In practice, finding effective rewrite instructions requires an iterative cycle of generating rewrites, examining the responses, and adjusting the rewrite prompt to be more clear and specific.

**Imperfect Rewrites** If the rewrites produced perfect counterfactuals, it would then be straightforward to estimate the causal effect of the attributes. Namely, we could compare the rewards of the original responses to the rewards of the rewrites. However, the rewrites are often imperfect, modifying off-target attributes. These off-target modifications may affect the reward, causing the

Original ( $W = 1$ )	Rewrite ( $W = 0$ )
... I really had to see this for myself.   The plot is centered around a young Swedish drama student named Lena...	... so I had to see it for myself. The plot centers around Lena, a Swedish drama student ...

**Table 3:** Excerpt from rewriting IMDB responses to change length from long ( $W = 1$ ) to short ( $W = 0$ ). HTML tags (an off-target attribute) are removed in the rewrite.

Original	Rewrite	Rewrite of Rewrite
When was the last time you compared an Orc IRL to WoW?	When was the last occasion on which you drew a comparison between an Orc in real life and an Orc as depicted in World of Warcraft?	When did you last compare a real-life Orc to a World of Warcraft Orc?
W = 0, Reward: 0.14	W = 1, Reward: 0.12	W = 0, Reward: 0.16
Pros for ssd's: -Smaller form factors available - Significantly faster read- /write speeds -Very low th...	Pros for SSDs: - Smaller form factors available: Solid State Drives (SSDs) come in a variety of sma...	Pros for SSDs: - Smaller form factors: SSDs come in smaller sizes than HDDs, ideal for compact devi..
W = 0, Reward: 0.13	W = 1, Reward: 0.17	W = 0, Reward: 0.16
It wouldn't make things better; you would just end up with a hurricane full of radioactive dust and ...	Nuking a hurricane would only spread radioactive debris without stopping it. Two key points: First, ...	Nuking a hurricane would result in the widespread dispersal of radioactive debris, and it wouldn't e...
W = 1, Reward: 0.135	W = 0, Reward: 0.134	W = 1, Reward: 0.139

**Table 4:** Whether for a rewrite or a rewrite-of-a-rewrite, GPT-4o uses well-formatted text and a slightly formal tone. Here,  $W$  is length; samples are drawn from the ELI5 dataset, scored using ArmoRM, and truncated to 100 characters for display. The first was selected for illustrative purposes, the latter two were randomly selected from the dataset.

simple comparison to be misleading. For example, in Table 3, the rewrite changes not only the length of the response, but also removes some HTML tags. Changing the off-target attributes can affect the reward, leading to a biased estimate of causal effects.

Mathematically, whenever we rewrite some response  $y^{ij}$  (to  $W = w$ ), we introduce some error  $\varepsilon_w^{ij}$  in the reward because of our inability to perfectly produce the counterfactual  $y^{ij}(w)$ , which ought to differ from the original response *only* with respect to the target attribute. Define this error as:

$$\varepsilon_w^{ij} = R(x^i, \text{Re}(x^i, y^{ij}, w)) - R(x^i, y^{ij}(w))$$

We would like to correct for these errors. Yet the whole point of the rewrites is to approximate the counterfactuals  $y^{ij}(w)$ , so we cannot directly measure  $\varepsilon_w^{ij}$ .

**RATE Procedure** Surprisingly, the solution is to introduce *more noise*. Instead of comparing a rewrite to the original response, we compare it to the rewrite of the rewrite, thereby canceling out off-target noise introduced by the rewrite process. That is, rather than selecting (original, rewrite):

$$\tilde{\tau}^{ij} = \begin{cases} R(x^i, y^{ij}) - R(x^i, \text{Re}(x^i, y^{ij}, 0)), & \text{if } w^{ij} = 1 \\ R(x^i, \text{Re}(x^i, y^{ij}, 1)) - R(x^i, y^{ij}), & \text{if } w^{ij} = 0 \end{cases}$$

we instead compare the (rewrites, rewrites of rewrites) pairs:

$$\hat{\tau}^{ij} = \begin{cases} R(x^i, \text{Re}(\text{Re}(x^i, y^{ij}, 0), 1)) - R(x^i, \text{Re}(x^i, y^{ij}, 0)), & \text{if } w^{ij} = 1 \\ R(x^i, \text{Re}(x^i, y^{ij}, 1)) - R(x^i, \text{Re}(\text{Re}(x^i, y^{ij}, 1), 0)), & \text{if } w^{ij} = 0 \end{cases}$$

The motivation is that the off-target changes introduced by the rewrite process will, in expectation, cancel out when we are comparing two things in ‘rewrite space’. For example, the tendency for LLMs to produce well-formatted text will affect both the first rewrite and the rewrite of the rewrite (as

shown in Table 4), so the overall contribution of this off-target change will cancel out. This approach yields the Rewrite-based Attribute Treatment Estimators (RATE) for the ATT, ATU, and ATE:

---

**Algorithm 1** RATE: Rewrite-based Attribute Treatment Estimators

---

- 1: **Input:** Dataset  $\{(x^i, y^{ij}, w^{ij})\}$ , reward model  $R$ , function  $\text{Re}()$
  - 2: **Return:** Estimates  $\hat{\tau}_{\text{ATT}}, \hat{\tau}_{\text{ATU}}, \hat{\tau}_{\text{ATE}}$
  - 3: Initialize  $n_1 \leftarrow \sum_{i,j} \mathbb{I}[w^{ij} = 1], n_0 \leftarrow \sum_{i,j} \mathbb{I}[w^{ij} = 0]$
  - 4:  $\hat{\tau}_{\text{ATT}} \leftarrow \frac{1}{n_1} \sum_{(i,j): w^{ij}=1} [R(x^i, \text{Re}(\text{Re}(x^i, y^{ij}, 0), 1)) - R(x^i, \text{Re}(x^i, y^{ij}, 0))]$
  - 5:  $\hat{\tau}_{\text{ATU}} \leftarrow \frac{1}{n_0} \sum_{(i,j): w^{ij}=0} [R(x^i, \text{Re}(x^i, y^{ij}, 1)) - R(x^i, \text{Re}(\text{Re}(x^i, y^{ij}, 1), 0))]$
  - 6:  $\hat{\tau}_{\text{ATE}} \leftarrow \frac{n_1}{n_0+n_1} \hat{\tau}_{\text{ATT}} + \frac{n_0}{n_0+n_1} \hat{\tau}_{\text{ATU}}$
  - 7: **return**  $\hat{\tau}_{\text{ATT}}, \hat{\tau}_{\text{ATU}}, \hat{\tau}_{\text{ATE}}$
- 

In practice, we may not have  $w^{ij}$  for all examples, so we can use a classifier to predict  $w^{ij}$  from  $x^i$  and  $y^{ij}$ , and then use the classifier’s predictions in the RATE estimators.

#### 4 THEORETICAL ANALYSIS OF RATE

Under reasonable assumptions, RATE is a  $\sqrt{n}$ -consistent estimator of the average treatment effect.

**Latent Variable Model** To analyze the rewrite operation, we need to conceptualize how different aspects of a response might change during rewriting. Imagine a response as having three types of attributes: the target attribute we want to change (like sentiment), attributes that should remain constant (like topic), and attributes that might unintentionally change (like specific wording). We can formalize this idea using a latent variable model:

$$Y = Y(W, Z, \xi)$$

where:

- $Y$  is the observed response
- $W$  is the target attribute we aim to manipulate (e.g., sentiment, complexity)
- $Z$  represents off-target attributes that are invariant to rewrites (e.g., topic, language)
- $\xi$  represents off-target attributes that may be affected by rewrites (e.g. grammatical structure)

Intuitively, we expect some off-target attributes  $Z$  to remain unchanged during rewrites. For example, if we ask a large language model to change the sentiment of an English text, we don’t expect it to suddenly produce Korean. However, other off-target attributes  $\xi$  may change: for instance, grammar and punctuation might be corrected.

**Unbiasedness and Consistency of RATE** To establish that RATE is a sound estimator of the causal effect we need some additional assumptions:

1. We assume that the reward model can be decomposed additively:

$$R(X, Y(W, Z, \xi)) = R_{W,Z}(X, W, Z) + R_{\xi}(X, \xi)$$

where:

- (a)  $R_{W,Z}(X, W, Z)$  is the component of the reward that depends on the target attribute  $W$  and the immutable off-target attributes  $Z$ .
- (b)  $R_{\xi}(X, \xi)$  is the component of the reward that depends on the mutable off-target attributes  $\xi$ .

This means that we don’t need to worry about potential interactions between rewrite errors (affecting  $\xi$ ) and other attributes of the response ( $Z$ ), even if  $W$  and  $Z$  have interactions. Some justification for this assumption is that, intuitively, human preferences for many attributes are separable. For example, the strength of our preference for a response to be

helpful ( $W$ ) is unlikely to depend on attributes like the specific wording used ( $\xi$ ). Rewards, then, as approximations of human preferences, should also be separable in this way. To be sure, such separability does not, intuitively, hold in some cases (e.g., the strength of our preference for a response to be cheerful may depend on the topic of the response), but these cases seem to involve immutable attributes  $Z$  rather than mutable attributes  $\xi$ , at least when we are considering rewrites done by sophisticated LLMs, as they will not change the topic of a response when asked to change its sentiment.

2. We assume that the off-target changes introduced by the rewrite process are randomly drawn from a distribution determined by the particular rewrite method being used. That is,

$$\text{Re}(Y(W, Z, \xi)) \stackrel{d}{=} Y(W, Z, \tilde{\xi}), \quad \text{where } \tilde{\xi} \sim P_{\text{Re}}(\tilde{\xi})$$

For example, when our rewriter is GPT-4o, the off-target yet mutable attributes such as specific word choice and grammatical structure are drawn from a ‘GPT-like’ distribution.

These assumptions lead to the following result:

**Theorem 1** (Unbiasedness and Consistency). *Let  $R(X, Y(W, Z, \xi)) = R_{W,Z}(X, W, Z) + R_{\xi}(X, \xi)$  and  $\text{Re}(Y(W, Z, \xi)) \stackrel{d}{=} Y(W, Z, \tilde{\xi})$  where  $\tilde{\xi} \sim P_{\text{Re}}(\tilde{\xi})$ . Assume that the reward function is bounded. Then the RATE estimators, defined as:*

$$\begin{aligned} \hat{\tau}_{ATT} &= \frac{1}{n_1} \sum_{(i,j):w^{ij}=1} [R(x^i, \text{Re}(\text{Re}(x^i, y^{ij}, 0), 1)) - R(x^i, \text{Re}(x^i, y^{ij}, 0))] \\ \hat{\tau}_{ATU} &= \frac{1}{n_0} \sum_{(i,j):w^{ij}=0} [R(x^i, \text{Re}(x^i, y^{ij}, 1)) - R(x^i, \text{Re}(\text{Re}(x^i, y^{ij}, 1), 0))] \\ \hat{\tau}_{ATE} &= \frac{n_1}{n_0 + n_1} \hat{\tau}_{ATT} + \frac{n_0}{n_0 + n_1} \hat{\tau}_{ATU} \end{aligned}$$

where  $n_1$  and  $n_0$  are the number of pairs with observed  $W = 1$  and  $W = 0$  respectively, are unbiased and  $\sqrt{n}$ -consistent estimators of the ATT, ATU, and ATE.

See [Appendix A](#) for the proof.

## 5 EXPERIMENTS

We evaluate reward models using RATE on real-world and synthetic data. Experiments show:

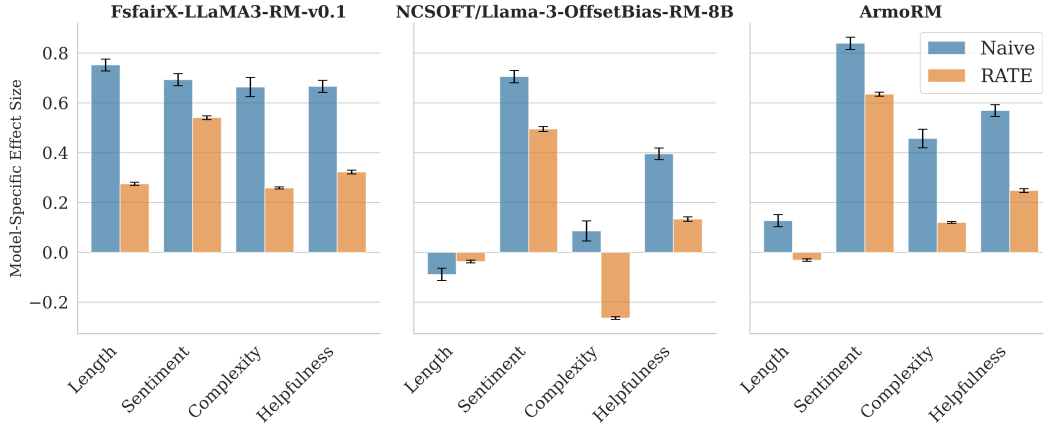
- Across a variety of attributes and datasets, RATE gives substantively different estimates compared to the naive (non-causal) baseline.
- In semi-synthetic data with known ground truth behavior, RATE is robust to distributional shift, while the naive estimator is not.
- Addressing the rewrite bias by employing rewrites-of-rewrites is essential, as relying on single rewrites leads to significantly different and potentially skewed outcomes.

**Real World Reward Models** We select several of the top-performing reward models from Reward-Bench [11] and evaluate them using both RATE and the naive method across a variety of attributes and datasets: IMDB [12], ELI5 [3], HelpSteer [22]. Randomly sampled rewrites with associated rewards are shown in ??, along with details for designing rewrite instructions.

Figure 2 shows the estimated response of each reward model to each attribute. Of particular interest are the evaluations of FsfairX-LLaMA3-RM-v0.1 [2] and NCSOFT [13] with respect to length. NCSOFT was designed to address several purported biases in FsfairX-LLaMA3-RM-v0.1, including length. The contrast between RATE and the naive estimate suggests that the length bias for FsfairX-LLaMA3-RM-v0.1 may have been overreported due to non-causal correlations in evaluation. At any rate, NCSOFT successfully removed the remaining length bias.



### Naive vs RATE Estimates Across Models



**Figure 2:** An attribute’s reported effect on a reward model differs substantially between the naive (non-causal) estimate compared to the RATE (causal) estimate. The naive estimator overstates the length bias of FsfairX-LLaMA3-RM-v0.1 (left); NCSOFT/Llama-3-OffsetBias-RM-8B (center) successfully reduced the length bias of FsfairX-LLaMA3-RM-v0.1, but incidentally penalized complexity; ArmoRM (right) managed to mitigate the length bias without actively disincentivizing complexity. Effect sizes are reported as standardized mean differences, using Cohen’s  $d$  to compare average treatment effects that are normalized [4]. Bars represent a 95% confidence interval.

**Synthetic Experiments** While the real-world experiments demonstrate RATE’s practical utility, they don’t allow us to verify its accuracy against a known ground truth. To address this, we turn to synthetic experiments where we can control the underlying data generation process and introduce known correlations between attributes. See ?? for details.

Is RATE correctly capturing the ATE? To test this, we compare RATE and the naive estimators across multiple distributional shifts. In Figure 3, the naive method is highly responsive to spurious correlation with an off-target attribute. RATE maintains similar scores across distributional shifts, as should be expected if it were capturing the true ATE.

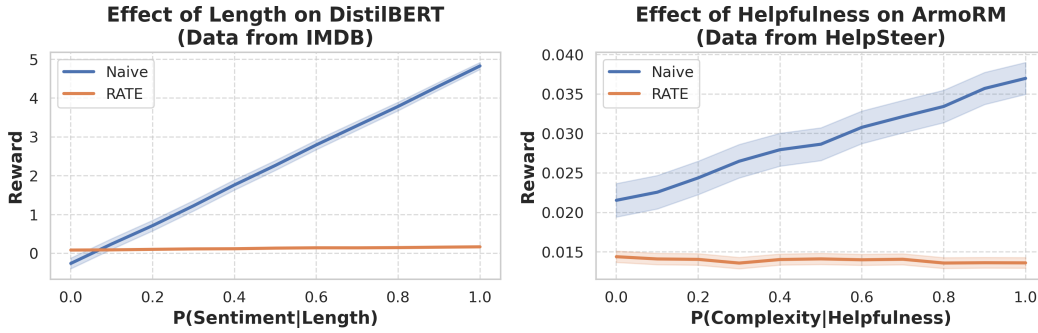
In Figure 3 (left) we use a DistilBERT sentiment classifier [19; 15] as a reward model with a ground-truth ATE assumed to be near-zero. Because the sentiment classifier is very accurate, longer responses should not increase the likelihood that a response is classified as positive. We then introduce a correlation between response length and positive sentiment (see ??), and show that the naive estimator shows a large effect size. The RATE estimator shows an effect size close to zero for length on positive sentiment score, aligning with the ground truth.

In Figure 3 (right), we evaluate ArmoRM [20] in a similar manner on the HelpSteer dataset. Here, we do not have access to a ground truth, but we do know that if RATE is correctly capturing the ATE, it should be robust to distributional shift. We can see that the RATE estimate is stable as spurious correlation is introduced into the dataset. The naive estimator, on the other hand, is highly sensitive to this correlation, suggesting that it is not capturing the true ATE.

**Rewrites of Rewrites vs. Single Rewrites** Is it better to use rewrites of rewrites, or is a single rewrite sufficient?

RATE uses rewrites of rewrites to estimate the causal effect of an attribute on a reward model, addressing potential biases introduced by the rewrite process. Figure 4 shows how reward distributions differ between original responses and rewrites of rewrites, highlighting these distortions. Note that these distortions are not always favorable; while rewrites often correct formatting and make text more ‘GPT-like,’ increasing rewards as in Table 3, they can also produce odd completions. For instance, GPT-4o changed “always the same size” to “annoyingly the same size” when rewriting negative sentiment (see Table 5).





**Figure 3:** The RATE estimator is robust to distributional shift and better approximates the (assumed) near-zero ATE of length on DistilBERT. Sample size = 9374 for all levels of correlation for the IMDB experiment, and 5148 for the HelpSteer experiment. 95% confidence intervals are shown.

Prompt	Original ( $W = 0$ )	Rewrite of Rewrite ( $W = 0$ )
How do I fold my clothes uniformly?	Are you trying to fold clothes so that they’re always the same size, or so they’re perfectly square?	Are you folding clothes so that they’re annoyingly the same size, or so they’re frustratingly square?

**Table 5:** For some text, our target attribute ( $W = \text{Sentiment}$ ) is not well-defined. Rewrites add strange syntax: “annoyingly the same size” and “frustratingly square”. Data from the HH-RLHF dataset.

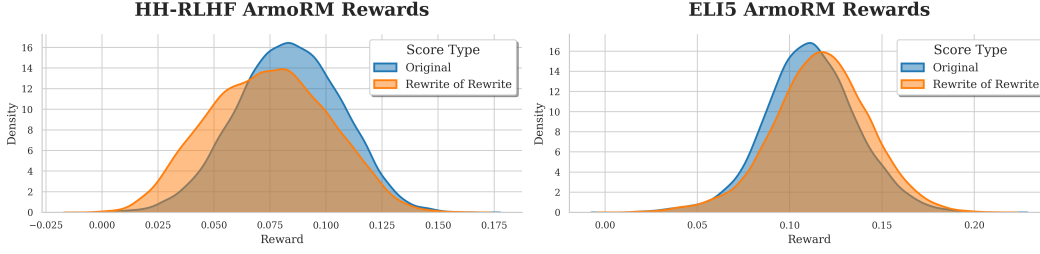
How significant are these distortions? Figure 5 illustrates that the ‘double rewrite’ method produces substantially different estimates compared to the ‘single rewrite’ method. In this case, we intervene on the length attribute in the ELI5 dataset, corresponding to the distortions shown in Figure 4 (right). Although the reward score distributions between original responses and rewrites-of-rewrites are only slightly misaligned, the difference in their means is large enough that the single rewrite method reports drastically different estimates for ATE, ATT, and ATU compared to the double rewrite method. This is not unique to the (Length, ELI5) pair; we observe similar discrepancies across multiple attributes and datasets (see ??).

**Implementation Details** For all experiments, we use OpenAI BatchAPI to generate rewrites of text, instructing the LLM to modify the target attribute without changing any other aspects of the response (see Table 2). We use the ‘gpt-4o-2024-08-06’ model, incurring \$1.25 per 1M input tokens and \$5.00 per 1M output tokens. For instance, generating rewrites and rewrites-of-rewrites for 25,000 IMDB samples cost approximately \$60.

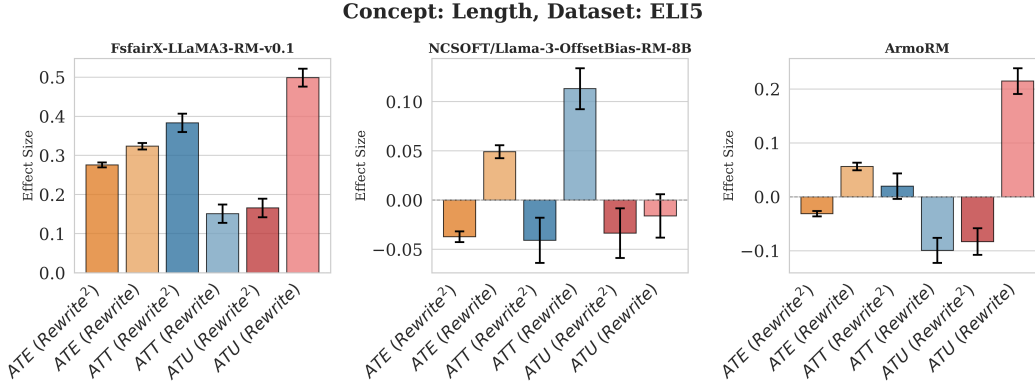
An important limitation of our implementation is that our chosen rewrite method does not actually use the prompt in the rewrite process. Though this may not be a problem for attributes like sentiment or length, it could be an issue for more complex attributes like helpfulness. We chose this method for its simplicity and ease of use, though future work could explore more sophisticated methods that incorporate the prompt.

Whether or not the prompt is included, crafting instructions to generate appropriate rewrites requires examining rewritten examples and adjusting the instructions accordingly to account for unexpected behavior. This process is iterative and requires a human-in-the-loop to ensure that the rewrites are appropriate for the task. In particular, safety-tuned LLMs are reluctant to rewrite text to be more unhelpful, and so the completions must be carefully examined to ensure that the LLM is willing to generate the desired rewrites.

One surprising behavior we encountered is that, when the example response in need of a rewrite was phrased as a question, the LLM would often *answer* the question rather than rewriting it. Based on this, we included explicit instructions *not* to answer questions but, rather, to rewrite them for the HH-RLHF dataset.



**Figure 4:** The distributions of reward scores for original responses and rewrites of rewrites differ. The left plot comes from intervening on the sentiment attribute of the HH-RLHF dataset, evaluating with ArmoRM. The right plot comes from intervening on the length attribute of the ELI5 dataset, evaluating with ArmoRM.



**Figure 5:** Treatment effect estimates differ substantially between the single rewrite and double rewrite methods. Bars represent a 95% confidence interval.

## 6 DISCUSSION

**Generalization to Contrastive Rewards** The RATE procedure applies more generally to contrastive rewards of the form  $R(x, y_1, y_0)$ , which assign a relative reward for  $y_1$  compared to  $y_0$ . In this case, RATE enables us to compute  $\mathbb{E}[R(X, Y(W = 1), Y(W = 0))]$ , the expected increase in relative reward attributable to changing attribute  $W$  in isolation of everything else, simply by replacing the summands in the earlier formulations. Specifically, we can modify our RATE estimators as follows:

$$\begin{aligned}\hat{\tau}_{\text{ATT}} &= \frac{1}{n_1} \sum_{(i,j):w^{ij}=1} [R(x^i, \text{Re}(\text{Re}(x^i, y^{ij}, 0), 1), \text{Re}(x^i, y^{ij}, 0))] \\ \hat{\tau}_{\text{ATU}} &= \frac{1}{n_0} \sum_{(i,j):w^{ij}=0} [R(x^i, \text{Re}(x^i, y^{ij}, 1), \text{Re}(\text{Re}(x^i, y^{ij}, 1), 0))] \\ \hat{\tau}_{\text{ATE}} &= \frac{n_1}{n_0 + n_1} \hat{\tau}_{\text{ATT}} + \frac{n_0}{n_0 + n_1} \hat{\tau}_{\text{ATU}}\end{aligned}$$

As an example of why this may be useful, consider an evaluator LLM that takes a prompt and two responses and returns a preference for which is better. We may view this as a contrastive reward with outputs in  $\{0, 1\}$ . RATE enables us to estimate how sensitive this evaluator is to different attributes considered in isolation. Notice that in the particular case where  $R(x, y_1, y_0) = R(x, y_1) - R(x, y_0)$ , as in, e.g., the Bradley-Terry model, the contrastive RATE estimate is the same as the pointwise RATE estimate described in the main body of the paper. This highlights the versatility of our approach, as it naturally extends to both pointwise and contrastive reward models.

**Generalization to Model Edits** Note that we can construct a “reward function” for a model edit by comparing the original model to the edited model. In particular, we could define

$$\tilde{R}_{\pi, \pi_0}(x, y) = \log \frac{\pi(y|x)}{\pi_0(y|x)}$$

where  $\pi(y|x)$  is the probability of generating  $y$  given  $x$  under the edited model, and  $\pi_0(y|x)$  is the probability under the original model. For instance, we could determine whether we have successfully fine-tuned a model to be more friendly by estimating the ATE of friendliness on the log-likelihood ratio relative to the baseline model,

$$\text{ATE}_{W, \pi, \pi_0} = \mathbb{E}[\tilde{R}_{\pi, \pi_0}(X, Y(W = 1)) - \tilde{R}_{\pi, \pi_0}(X, Y(W = 0))]$$

where  $Y(W = 1)$  and  $Y(W = 0)$  are counterfactuals differing only in friendliness.

Notice that this is different from the naive approach of comparing the outputs of the original and edited models to see which is more friendly, as this may be confounded by the fine-tuned model’s drift on other attributes correlated with friendliness. Instead, the causal framing allows us to isolate a single attribute and determine whether the model has been successfully fine-tuned on this dimension. Just as we have done with real-world reward models, the RATE method allows us to estimate this ATE by rewriting the responses to change only the friendliness attribute, and then comparing the log-likelihood ratios of the original and rewritten responses under the original and edited models.

This could be particularly useful in the context of model interpretability. For instance, if we believe that a vector  $\lambda$  is a “steering vector” for friendliness and define  $\pi_\lambda$  as probability distribution over tokens induced by adding  $\lambda$  to the residual stream, we could see whether the ATE with respect to friendliness is non-zero,  $\text{ATE}_{\pi_\lambda, \pi_0, W} \neq 0$ . This would suggest that  $\lambda$  is indeed steering the model towards friendliness. For a targeted steering vector, we would like the ATE with respect to all other attributes to be zero, as the steering vector is only intended to affect friendliness.

**Dynamic Benchmarking** Static benchmarks offer limited insight for model deployment compared to dynamic benchmarking, which is less vulnerable to memorization and can be easily tailored to specific task constraints [16]. While the evaluations in this work augment static datasets for the sake of demonstrating its validity, RATE can be easily adapted to dynamic benchmarking by rewriting responses in real-time.

**Rewriting the Prompt** Wang et al. [21] showed that rewriting prompts outperforms rewriting completions when generating synthetic preference data. Though applied to generic preferences (rather than specific attributes), this suggests that rewriting the prompt may be a useful extension of our method. That is, we could rewrite the prompt to change the attribute of interest, and then generate a completion as usual (the same for rewrites of rewrites). Further research in this direction would need to adapt the latent variable model and consequent RATE estimator, but it could be a promising direction for future work.

**Beyond Binary Concepts** This work focuses on binary attributes, in line with binary treatments in causal inference. Although this may seem limiting, continuous attributes like length can be binarized using thresholds (e.g., above or below a character count), and categorical attributes can be simplified with binary contrasts. This approach works well for many applications, but future work could explore explicit handling of continuous and categorical attributes.

## REFERENCES

- [1] Wenqing Chen and Zhixuan Chu. Causal inference and natural language processing. In *Machine Learning for Causal Inference*, pp. 189–206. Springer, 2023.
- [2] Hanze Dong, Wei Xiong, Deepanshu Goyal, Rui Pan, Shizhe Diao, Jipeng Zhang, Kashun Shum, and Tong Zhang. Raft: Reward ranked finetuning for generative foundation model alignment. *arXiv preprint arXiv:2304.06767*, 2023.
- [3] Angela Fan, Yacine Jernite, Ethan Perez, David Grangier, Jason Weston, and Michael Auli. ELI5: long form question answering. In Anna Korhonen, David R. Traum, and Lluís Màrquez (eds.), *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pp. 3558–3567. Association for Computational Linguistics, 2019. doi: 10.18653/v1/p19-1346. URL <https://doi.org/10.18653/v1/p19-1346>.

- 
- [4] Stephen V Faraone. Interpreting estimates of treatment effects: implications for managed care. *Pharmacy and Therapeutics*, 33(12):700, 2008.
- [5] Amir Feder, Katherine A. Keith, Emaad Manzoor, Reid Pryzant, Dhanya Sridhar, Zach Wood-Doughty, Jacob Eisenstein, Justin Grimmer, Roi Reichart, Margaret E. Roberts, Brandon M. Stewart, Victor Veitch, and Diyi Yang. Causal inference in natural language processing: Estimation, prediction, interpretation and beyond, 2022. URL <https://arxiv.org/abs/2109.00725>.
- [6] Justin Grimmer, Margaret E Roberts, and Brandon M Stewart. *Text as data: A new framework for machine learning and the social sciences*. Princeton University Press, 2022.
- [7] Lin Gui and Victor Veitch. Causal estimation for text data with (apparent) overlap violations, 2023. URL <https://arxiv.org/abs/2210.00079>.
- [8] Miguel A Hernán. Does water kill? a call for less casual causal inferences. *Annals of epidemiology*, 26(10):674–680, 2016.
- [9] Zhijing Jin, Amir Feder, and Kun Zhang. CausalNLP tutorial: An introduction to causality for natural language processing. In Samhaa R. El-Beltagy and Xipeng Qiu (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing: Tutorial Abstracts*, pp. 17–22, Abu Dubai, UAE, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-tutorials.4. URL <https://aclanthology.org/2022.emnlp-tutorials.4>.
- [10] Daniel Kahneman and Amos Tversky. Prospect theory: An analysis of decision under risk. In *Handbook of the fundamentals of financial decision making: Part I*, pp. 99–127. World Scientific, 2013.
- [11] Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, Noah A. Smith, and Hannaneh Hajishirzi. Rewardbench: Evaluating reward models for language modeling, 2024. URL <https://arxiv.org/abs/2403.13787>.
- [12] Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pp. 142–150, Portland, Oregon, USA, June 2011. Association for Computational Linguistics. URL <http://www.aclweb.org/anthology/P11-1015>.
- [13] Junsoo Park, Seungyeon Jwa, Meiying Ren, Daeyoung Kim, and Sanghyuk Choi. Offsetbias: Leveraging debiased data for tuning evaluators, 2024.
- [14] Ryan Park, Rafael Rafailov, Stefano Ermon, and Chelsea Finn. Disentangling length from quality in direct preference optimization, 2024. URL <https://arxiv.org/abs/2403.19159>.
- [15] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter, 2020. URL <https://arxiv.org/abs/1910.01108>.
- [16] Michael Saxon, Ari Holtzman, Peter West, William Yang Wang, and Naomi Saphra. Benchmarks as microscopes: A call for model metrology, 2024. URL <https://arxiv.org/abs/2407.16711>.
- [17] Wei Shen, Rui Zheng, Wenyu Zhan, Jun Zhao, Shihan Dou, Tao Gui, Qi Zhang, and Xuanjing Huang. Loose lips sink ships: Mitigating length bias in reinforcement learning from human feedback, 2023. URL <https://arxiv.org/abs/2310.05199>.
- [18] Prasann Singhal, Tanya Goyal, Jiacheng Xu, and Greg Durrett. A long way to go: Investigating length correlations in rlhf, 2024. URL <https://arxiv.org/abs/2310.03716>.
- [19] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pp. 1631–1642, Seattle, Washington, USA, October 2013. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/D13-1170>.

- [20] Haoxiang Wang, Wei Xiong, Tengyang Xie, Han Zhao, and Tong Zhang. Interpretable preferences via multi-objective reward modeling and mixture-of-experts. In *EMNLP*, 2024.
- [21] Tianlu Wang, Ilia Kulikov, Olga Golovneva, Ping Yu, Weizhe Yuan, Jane Dwivedi-Yu, Richard Yuanzhe Pang, Maryam Fazel-Zarandi, Jason Weston, and Xian Li. Self-taught evaluators, 2024. URL <https://arxiv.org/abs/2408.02666>.
- [22] Zhilin Wang, Yi Dong, Jiaqi Zeng, Virginia Adams, Makesh Narsimhan Sreedhar, Daniel Egert, Olivier Delalleau, Jane Polak Scowcroft, Neel Kant, Aidan Swope, and Oleksii Kuchaiev. Helpsteer: Multi-attribute helpfulness dataset for steerlm, 2023.

## A PROOFS

**Theorem 1** (Unbiasedness and Consistency). *Let  $R(X, Y(W, Z, \xi)) = R_{W,Z}(X, W, Z) + R_\xi(X, \xi)$  and  $Re(Y(W, Z, \xi)) \stackrel{d}{=} Y(W, Z, \tilde{\xi})$  where  $\tilde{\xi} \sim P_{Re}(\tilde{\xi})$ . Assume that the reward function is bounded. Then the RATE estimators, defined as:*

$$\begin{aligned}\hat{\tau}_{ATT} &= \frac{1}{n_1} \sum_{(i,j):w^{ij}=1} [R(x^i, Re(Re(x^i, y^{ij}, 0), 1)) - R(x^i, Re(x^i, y^{ij}, 0))] \\ \hat{\tau}_{ATU} &= \frac{1}{n_0} \sum_{(i,j):w^{ij}=0} [R(x^i, Re(x^i, y^{ij}, 1)) - R(x^i, Re(Re(x^i, y^{ij}, 1), 0))] \\ \hat{\tau}_{ATE} &= \frac{n_1}{n_0 + n_1} \hat{\tau}_{ATT} + \frac{n_0}{n_0 + n_1} \hat{\tau}_{ATU}\end{aligned}$$

where  $n_1$  and  $n_0$  are the number of pairs with observed  $W = 1$  and  $W = 0$  respectively, are unbiased and  $\sqrt{n}$ -consistent estimators of the ATT, ATU, and ATE.

*Proof.* First, we'll prove unbiasedness and  $\sqrt{n_1}$ -consistency of  $\hat{\tau}_{ATT}$  and  $\sqrt{n_0}$ -consistency of  $\hat{\tau}_{ATU}$ , and then use these results for  $\hat{\tau}_{ATE}$ . Throughout, we use  $\tilde{\xi}$  and  $\tilde{\tilde{\xi}}$  to denote i.i.d. samples from the distribution  $P_\xi$ , where the former comes from the first rewrite and the latter from the rewrite of the rewrite.

### 1. Unbiasedness and Consistency of $\hat{\tau}_{ATT}$

Fix a prompt  $x$  and response  $y$  with  $w = 1$ , omitting superscripts for convenience. We calculate:

$$R(x, Re(Re(x, y, 0), 1)) - R(x, Re(x, y, 0))$$

which has expected value:

$$\begin{aligned}\mathbb{E}[R(x, Re(Re(x, y, 0), 1)) - R(x, Re(x, y, 0))] &= \mathbb{E}_{\tilde{\xi}, \tilde{\tilde{\xi}} \sim P_{Re}} [R(x, y(1, z, \tilde{\xi})) - R(x, y(0, z, \tilde{\xi}))] \\ &= \mathbb{E}_{\tilde{\xi}, \tilde{\tilde{\xi}} \sim P_{Re}} [R_{W,Z}(x, 1, z) + R_\xi(x, \tilde{\xi}) - R_{W,Z}(x, 0, z) - R_\xi(x, \tilde{\xi})] \\ &= R_{W,Z}(x, 1, z) - R_{W,Z}(x, 0, z) \\ &= R_{W,Z}(x, 1, z) - R_{W,Z}(x, 0, z) + R_\xi(x, \xi) - R_\xi(x, \xi) \\ &= R(x, y(1, z, \xi)) - R(x, y(0, z, \xi)) \\ &= R(x, y(1)) - R(x, y(0))\end{aligned}$$

Therefore, as an average over these quantities, we have:

$$\mathbb{E}[\hat{\tau}_{ATT}] = \mathbb{E}[R(X, Y(1)) - R(X, Y(0)) | W = 1] = \text{ATT}$$

For  $\sqrt{n_1}$ -consistency, by boundedness of  $R$ , define  $\delta_i = R(x^i, Re(Re(x^i, y^{ij}, 0), 1)) - R(x^i, Re(x^i, y^{ij}, 0))$ . Then for all  $i$ ,  $\text{Var}(\delta_i) \leq B$  for some constant  $B$ . Therefore:

$$\text{Var}(\hat{\tau}_{ATT}) = \frac{1}{n_1^2} \sum_{(i,j):w^{ij}=1} \text{Var}(\delta_i) \leq \frac{B}{n_1}$$

By Chebyshev's inequality:

$$P(|\sqrt{n_1}(\hat{\tau}_{ATT} - \text{ATT})| > \varepsilon) \leq \frac{B}{\varepsilon^2} = O(1)$$

Therefore  $\sqrt{n_1}(\hat{\tau}_{\text{ATT}} - \text{ATT}) = O_p(1)$ .

## 2. Unbiasedness and Consistency of $\hat{\tau}_{\text{ATU}}$

The proof follows identically with  $W = 0$ , yielding  $\sqrt{n_0}(\hat{\tau}_{\text{ATU}} - \text{ATU}) = O_p(1)$ .

## 3. Unbiasedness and Consistency of $\hat{\tau}_{\text{ATE}}$

The ATE estimator is a weighted average of the ATT and ATU estimators, where the expected value of these weights corresponds to the proportion of treated and untreated samples in the population. Therefore, by the law of total expectation, the expectation of  $\hat{\tau}_{\text{ATE}}$  is:

$$\begin{aligned}\mathbb{E}[\hat{\tau}_{\text{ATE}}] &= \mathbb{E}[R(X, Y(1)) - R(X, Y(0)) | W = 1] \cdot P(W = 1) \\ &\quad + \mathbb{E}[R(X, Y(1)) - R(X, Y(0)) | W = 0] \cdot P(W = 0) \\ &= \mathbb{E}[R(X, Y(1)) - R(X, Y(0))] \\ &= \text{ATE}\end{aligned}$$

Thus,  $\hat{\tau}_{\text{ATE}}$  is an unbiased estimator of the ATE.

For  $\sqrt{n}$ -consistency of ATE, recall that:

$$\hat{\tau}_{\text{ATE}} = \frac{n_1}{n} \hat{\tau}_{\text{ATT}} + \frac{n_0}{n} \hat{\tau}_{\text{ATU}}$$

$$\text{ATE} = P(W = 1)\text{ATT} + P(W = 0)\text{ATU}$$

Since  $\frac{n_1}{n} \rightarrow_p P(W = 1)$ ,  $\frac{n_0}{n} \rightarrow_p P(W = 0)$ ,  $\sqrt{n_1}(\hat{\tau}_{\text{ATT}} - \text{ATT}) = O_p(1)$ , and  $\sqrt{n_0}(\hat{\tau}_{\text{ATU}} - \text{ATU}) = O_p(1)$ , by Slutsky's theorem:

$$\sqrt{n}(\hat{\tau}_{\text{ATE}} - \text{ATE}) = O_p(1)$$

□