# RORA: Robust Free-Text Rationale Evaluation

**Zhengping Jiang**[♡]    **Yining Lu**[♡]    **Hanjie Chen**
**Daniel Khashabi    Benjamin Van Durme    Anqi Liu**
Johns Hopkins University
{zjiang31, ylu130}@jhu.edu

## Abstract

Free-text rationales play a pivotal role in explainable NLP, bridging the knowledge and reasoning gaps behind a model's decision-making. However, due to the diversity of potential reasoning paths and a corresponding lack of definitive ground truth, their evaluation remains a challenge. Existing evaluation metrics rely on the degree to which a rationale *supports* a target label, but we find these fall short in evaluating rationales that inadvertently *leak the labels*. To address this problem, we propose RORA, a RObust free-text RAtionale evaluation against label leakage.[1] RORA quantifies the new information supplied by a rationale to justify the label. This is achieved by assessing the conditional $\mathcal{V}$-information (Hewitt et al., 2021) with a predictive family robust against leaky features that can be exploited by a small model. RORA consistently outperforms existing approaches in evaluating human-written, synthetic, or model-generated rationales, particularly demonstrating robustness against label leakage. We also show that RORA aligns well with human judgment, providing a more reliable and accurate measurement across diverse free-text rationales.

## 1 Introduction

The ability of large language models (LLMs) to generate free-text rationales that elaborate on their decision-making processes holds promise for explainable NLP, either in the form of a reasoning chain (Wei et al., 2022; Yao et al., 2023) or post-hoc explanations (Madaan et al., 2023; Zheng et al., 2023). Previous works have also collected human-written rationales to enhance model reasoning and the generation of free-text rationales (Rajani et al., 2019; Camburu et al., 2018; Aggarwal et al., 2021; Geva et al., 2021).

However, evaluating these rationales remains an open problem because of the diversity of reasoning paths and the lack of definitive ground truth (Chan et al., 2022). As a result, existing metrics rely on measuring how much the rationale supports a given label. This is usually achieved by comparing predictions of models trained with and without rationales. For example, Leakage-Adjusted Simulatability (LAS) (Hase et al., 2020) and Rationale Quality (RQ) (Wiegreffe et al., 2020) measure rationale quality through the difference in accuracy. Alternatively, Rationale Evaluation with conditional-$\mathcal{V}$-information (REV) (Chen et al., 2023b) evaluates the reduction in model predictive uncertainty upon conditioning on the rationale.

Yet all these methods are vulnerable to *label leakage* (Li et al., 2022; Chen et al., 2023c; Ludan et al., 2023): the rationale inadvertently *paraphrasing* or *restating* labels, creating a spurious shortcut (Geirhos et al., 2020) for the evaluation model to infer the label. The critical issue stems from the mismatch in objectives: existing methods evaluate how easy it is to utilize information in the rationale, but rationales are *explanations*, whose quality does not always come with simplicity. The best explanation has to support the answer through some sense of mechanisms, such as methodically considering a set of axioms and running through a deductive chain (Bechtel and Abrahamsen, 2005; Keil, 2006; Glennan, 2002), without which they are mere "effects" (Cummins, 2000). Figure 1 shows an example where existing evaluation methods are highly sensitive to label leakages in paraphrased or restated rationales, while in fact, these label leakages merely increase the predictive probability without providing any meaningful explanations.[2]

With this objective in mind, we introduce RORA, a novel approach to evaluate rationales robust to label leakage. RORA's construction consists of

---

[♡]Equal contribution
[1]https://github.com/zipJiang/RORA

[2]Note that scores between different evaluation metrics are not directly comparable because of different scales and criteria. In this paper, our analysis mainly focuses on the ranking and relative differences within each metric.
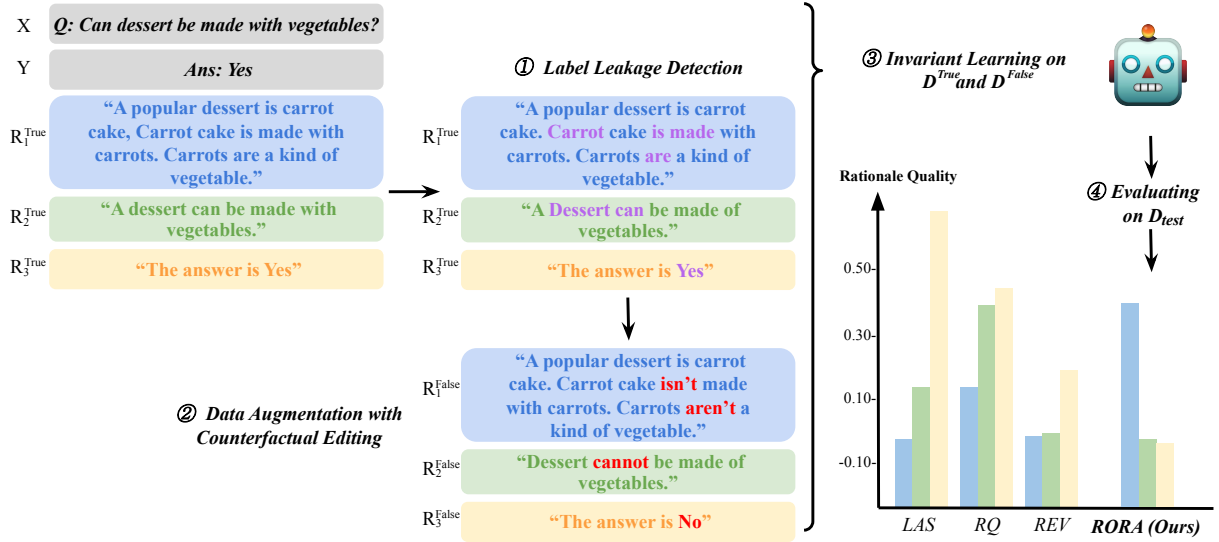
Figure 1: RORA framework for evaluating rationales $R_1^{\text{True}}$, $R_2^{\text{True}}$, $R_3^{\text{True}}$. Existing baselines are highly sensitive to rationales that simply restate the label or paraphrase the given question and label, leading to inflated scores compared to the human-annotated rationale. In contrast, RORA provides an informativeness score that better characterizes rationale quality. It is achieved by ① detecting potential leakage tokens in the rationale (§3.1) and ② generate additional training data with counterfactual editing for data augmentation (§3.2), followed by ③ training an evaluation model invariant to label leakage (§3.3).

three stages as illustrated in Figure 1. First, we fit a small model and identify label-leaking tokens via its gradient-based attributions (§3.2). After that, we generate additional training data with counterfactual editing (Ross et al., 2021) (§3.2). Finally, we force the evaluation model to ignore these label-leaking tokens through invariant learning (Arjovsky et al., 2020) (§3.3). Our approach aligns with the human perception that explanations should apparently increase the understanding of a given phenomenon by helping to create knowledge and to develop better theories (Wilson and Keil, 1998). On the contrary, label leakage tends to be repetitive and tautological (Aslanov and Guerra, 2023), dominating the insightful parts of the explanation.

We compare RORA with baseline metrics (Hase et al., 2020; Wiegreffe et al., 2020; Chen et al., 2023b) in evaluating various synthetic and human-annotated rationales, with or without label leakage, on three QA datasets. RORA consistently outperforms baseline metrics by providing robust evaluations against label leakages. We also compare RORA against model-generated rationale evaluation and demonstrate its better agreement with human evaluation.

## 2 Motivation and Backgrounds

Following previous work by Chen et al. (2023b), RORA adopts the framework of conditional $\mathcal{V}$-

information (Hewitt et al., 2021). In this section, We outline the general idea of this framework and point out how the objective of RORA is critically different from previous attempts.

### 2.1 Conditional $\mathcal{V}$-information

The theory of $\mathcal{V}$-information, as proposed by Xu et al. (2020), focuses on quantifying usable information under computational constraints. Specifically, it examines how much usable information about a random variable $Y$ can be derived from another random variable $X$ by applying functions belonging to a specified set known as the *predictive family* $\mathcal{V}$. Due to its transparent connection to baseline probing, Hewitt et al. (2021) propose to extend $\mathcal{V}$-information to multivariable cases as an evaluation of information beyond a baseline.

**Definition 1** (Multivariable $\mathcal{V}$-information). Let $X_1, X_2 \ldots, X_n \in \mathcal{X}_1, \mathcal{X}_2, \ldots, \mathcal{X}_n$ and $Y \in \mathcal{Y}$ be random varaibles. Let $\mathcal{V}$ be a multivariable predictive family.[3] Then the conditional multivariable $\mathcal{V}$-information from $X_j$ to $Y$, where $j \in \{1, 2, \ldots, n\}$, conditioned on prior knowledge

---

[3]We refer the readers to the work of Xu et al. (2020) We adopt their notation (rather than that of Hewitt et al. (2021)) since it abstracts away actual implementation from the underlying idea.

$C \subset \{X_1, \ldots, X_n\}$, is defined as

$$I_\mathcal{V}(X_l \rightarrow Y|C) = H_\mathcal{V}(Y|C) - H_\mathcal{V}(Y|C \cup X_l), \quad (1)$$

where

$$H_\mathcal{V}(Y|C) =$$
$$\inf_{f \in \mathcal{V}} \mathbb{E}_{c,y} \Big[ -\log f\big[c, \{\varnothing, \ldots, \varnothing\}\big](y) \Big]. \quad (2)$$

is the multivariable predictive $\mathcal{V}$-entropy.[4]

*Remark.* Equation 1 is also called the "conditional $\mathcal{V}$-information" (and Equation 2 "conditional $\mathcal{V}$-entropy"). Equation 1 characterizes how much uncertainty of a variable $Y$ can be reduced when a certain class of function was used to extract knowledge from observing $X_l$. The formulation of Equation 1 allows estimation of mutual information via optimization over $\mathcal{V}$, i.e., through gradient descent. It is important to note that the definition of the predictive family ensures that a predictive family can ignore any feature subset while the remaining subset $c$ still has the same output. In our case, it ensures that one can always construct a valid, more limited predictive family by forcing ignorance on leaky features.

## 2.2 Evaluating Label-Leaking Rationales

The conditional $\mathcal{V}$-information provides a natural way to evaluate the informativeness of rationales. For a natural language problem with input $X$ and label $Y$, A direct application of multivariable $\mathcal{V}$-information to evaluate the quality of rationale $R$ is

$$Q(R) = I_\mathcal{V}(R \rightarrow Y|X). \quad (3)$$

Yet, as Hase et al. (2020) pointed out, this approach fails to distinguish between different ways a rationale provides knowledge for a model's prediction. For example, a bad rationale can obtain a very high $Q(R)$ by restating the label to be predicted. To solve this problem, REV (Chen et al., 2023b) proposes to introduce a vacuous rationale $B$ that is simply a declarative combination of $X$ and $Y$

$$\text{REV}(R) = I_\mathcal{V}(R \rightarrow Y|B). \quad (4)$$

However, this requires that $Y$ is not completely determined by $B$, which can be easily violated by

label leakage in $B$ that makes $H_\mathcal{V}(Y|B)$ close to 0.[5]

Under a similar framework where conditional $\mathcal{V}$-entropy $H_\mathcal{V}$ in Equation 1 is replaced by prediction accuracy, other methods are proposed to mitigate the impact of label leakage. Hase et al. (2020) propose LAS: an averaged accuracy between leaking and non-leaking rationales which are classified by whether a simulator model can predict $Y$ solely from the explanation $R$. However, the prevalence of label leakage in free-text rationales may overwrite the effect of macro-averaging. The metric is simply undefined when all rationales in the dataset leak labels (Pruthi et al., 2022). Sia et al. (2022) suggest that using logical counterfactuals may circumvent the label leakage issue. However, their approach does not generalize to most NLP tasks because rationales in free-text format are usually too complex to be parsed logically.

Suppose a free-text rationale $R$ consists of label leakage $R_L$ and non-leaky information $R_R$ (i.e., $R = R_L \cup R_R$). We propose that a good evaluator $\mathcal{S}$ of the informativeness of rationale $R$ would be

$$\mathcal{S}(R) = I_\mathcal{V}(R_R \rightarrow Y|X),$$
$$= I_\mathcal{V}(R \setminus R_L \rightarrow Y|X). \quad (5)$$

Notice that we do not condition on $B$, as in our formulation, vacuous rationales are similar to spurious features about which our evaluator should be ignorant (Xu et al., 2020). $R_L$ and $R_R$ are usually deeply entangled in natural language utterances, and it is difficult to isolate $R_R$ completely. To address this problem, we propose RORA to approximate Equation 5 by identifying and forcing the model to unlearn $R_L$, as detailed in §3.

## 3 RORA for Rationale Evaluation

We describe our framework (Figure 1) to estimate Equation 5 as a measurement of the informativeness of rationales.

### 3.1 Leakage Detection

To detect leaking features $R_L$, we fit a relatively small model on corresponding rationales $\Phi_{\text{small}} : \mathcal{R} \rightarrow \mathcal{P}(\mathcal{Y})$ that predicts the label distribution from the rationale alone. We then calculate the attribution of its prediction to each input token using Integrated Gradient (Sundararajan et al., 2017).

---

[4]Following Xu et al. (2020) $\varnothing$ is used to indicate that no information is provided for this feature. In practice, it is common to mask out this feature with constant values (Hewitt et al., 2021).

[5]Conditioning on $B$ may cause other problems, of which we refer the readers to Appendix B for a more detailed discussion.

This token-level attribution is then averaged over the whole dataset to get global attribution for each token. Figure 1 shows example attribution on different rationales. The highlight (in purple) denotes the token is critical to the model's prediction as its attribution is above a given threshold.

We would prefer the model to be *less contextualized, less-trained, and smaller in size.* Our motivation is that smaller models won't be able to learn reasoning features, so their attributions tend to rely on shortcuts (Geirhos et al., 2020) like label leakage. In addition, attribution for smaller models is known to be more reliable and meaningful (Neely et al., 2022).

## 3.2 Data Augmentation

Given the original dataset $D = (X, R, Y)$, we generate partially counterfactual datasets

$$D^e = (X, R^e, Y), \ \forall e \in \mathcal{Y},$$

where editing is constrained to the spans identified in the previous step. Following MICE (Ross et al., 2021), we train a label-associated sequence to sequence infiller. Specifically for the $i$-th datapoint $(x_i, r_i, y_i) \in D$, we create training data by processing the sentence $y_i \oplus x_i \oplus r_i$ with the same denoising objective as in T5 pretraining (Raffel et al., 2020), where a single sentinel mask replaces each consecutive span from gradient attribution (§3.1) as inputs and the training objective is to predict masked spans. Notice that we only mask within the rationale to control artificially introduced artifacts as discussed in the Appendix A. To augment a dataset $D$, we decode infilling for these masked spans for each $e \in \mathcal{Y}$ to create $D^e$. Notice that, unlike usual counterfactual generation, we do not necessarily require the generated rationale $r_i^e$ with prepending label $e$ to be highly associated with that label. In fact, the training of the editing model in MICE can be thought of as estimating

$$I_\mathcal{V}(Y \to R | X).$$

Therefore, only those parts of the rationale that are easy to predict from the label will be edited, and most other spans will likely remain the same (as shown in Figure 1), as the model predicts these spans mostly from the static $x_i$ instead of the flipped $y$. This is a desirable property in our case, as it leads to hyperparameters that are easy to tune, with which we apply a very aggressive removal threshold while still preserving most of the $R_R$.

## 3.3 Computing RORA

Given all the augmented datasets $\mathcal{D} := \{D^e | \forall e \in \mathcal{Y}\}$, we train our evaluator on the joint of these datasets.[6] We learn an invariant predictor parameterized by $\Phi$ by applying the Invariant Risk Minimization (IRMv1) regularization (Arjovsky et al., 2020)

$$\mathcal{L}(\Phi) = \sum_{D^e \in \mathcal{D}} \mathcal{L}^e(\Phi) \\ + \lambda \|\nabla_{w|w=1.0} \mathcal{L}^e(w \cdot \Phi)\|^2, \quad (6)$$

where $\mathcal{L}^e$ is the Empirical Risk Minimization (ERM) term on dataset $D^e$ and the second term is the IRM regularization to enforce $\Phi$ to be simultaneously optimal on all environments $D^e$. In our case, it is either the factual data or the counterfactual augmentation. $\lambda$ is a hyperparameter balancing the tradeoff between the predictive power (an ERM term) and the invariance of the predictor across different environments. When the underlying model is a classifier, Equation 6 is directly applicable with normal cross-entropy loss. However, considering the increasing trend of directly decoding label sequence from pre-trained language models like T5 (Raffel et al., 2020; Nogueira et al., 2021) and to better align with previous works (Chen et al., 2023b), we use the following variant with label decoding:

$$\mathcal{L}(\Phi) = \sum_{D^e \in \mathcal{D}} \mathbb{E}_{(x,r,y) \sim D^e}[-\log p_\Phi(y|x,r)] \\ + \lambda \|\nabla_{w|w=1.0} \log \frac{p_\Phi(y|x,r)}{\sum_{y' \in \mathcal{Y}} p_\Phi(y'|x,r)})\|^2. \quad (7)$$

Notice that we still train to maximize probability over all possible strings in the ERM term to facilitate greedy decoding. In an ideal world with infinite data, this is equivalent to the conditional probability on the set of all possible labels.

Additionally, we train a baseline model $\Theta$ by minimizing

$$\mathcal{L}_{\text{Base}}(\Theta) = \mathbb{E}_{(x,r,y) \sim D}[-\log p_\Theta(y|x,\varnothing)]. \quad (8)$$

As Equation 5 can be decomposed into the difference between two conditional $\mathcal{V}$-entropy terms

$$S(R) = H_\mathcal{V}(Y|X) - H_\mathcal{V}(Y|R \setminus R_L \cup X).$$

---

[6] In reality, it is possible to reuse the original dataset $D$, as in a single batch the model sees corresponding data points from all datasets.

| Dataset | Question | *gold* | *leaky* | *gold$_{leaky}$* | *vacuous* |
|---------|----------|--------|---------|------------------|-----------|
| StrategyQA | Could Chuck Norris ride a horse? | Chuck Norris is a person. Horses are bigger than people. People can ride horses. | The answer is True. | Chuck Norris is a person. Horses are bigger than people. People can ride horses. The answer is True. | Chuck Norris could ride a horse. |

Table 1: Example of human-annotated and synthetic leaky rationales for fixed-label QA task. *leaky* verbatim leaks the label, *vacuous* is the declarative combination of the question and the given answer, and *gold$_{leaky}$* is simply the concatenation of the *gold* and *leaky* rationales.

With $\Theta$ and $\Phi$ from the same predictive family, the RORA score $S(R)$ can thus be estimated as

$$\hat{S}(R) = \mathbb{E}_{(x,r,y)\sim D_{\text{test}}} \big[ \log p_\Phi(y|x,r) \\ - \log p_\Theta(y|x,\varnothing) \big]. \quad (9)$$

Notice that the testing dataset $D_{\text{test}}$ has the same rationale distribution as the original dataset $D$.

If the model is simultaneously optimal for all $D^e \in \mathcal{D}$ (Arjovsky et al., 2020), it is ignorant (Xu et al., 2020) of the features that are different on $\forall D^e \in \mathcal{D}$ (otherwise, the predictive distribution will be different while the label distribution is the same). In our case, these features are supposed to be related to label leakages. We also discuss why naively masking out leaking parts does not work for free-text rationale evaluation in Appendix A.

## 4  Experiments and Results

To demonstrate the effectiveness of RORA, we evaluate human-annotated, synthetic, and model-generated rationales on two distinct QA tasks. For synthetic data, we (1) write a sentence that verbatim leaks the label (*leaky*); (2) perturb the human-annotated rationales (*gold*) to include the *leaky* rationale (*gold$_{leaky}$*); (3) create declarative rewrites (*vacuous*) of the question-answer pairs using a question rewriter (Chen et al., 2021).[7] Table 1 shows examples of these rationales.[8] For all synthetic and human-annotated data, **a good metric for rationale quality should be able to evaluate *gold* and *gold$_{leaky}$* reasonably close, while ranking them above *leaky* and *vacuous* rationales**. For model-generated rationales, **a good metric should evaluate rationales aligning with human judgments.** To create model-generated rationales of varying qualities, we choose GPT-4 (OpenAI, 2023), GPT-3.5 (Ouyang et al., 2022), Llama2-7b (Touvron et al., 2023) and Flan-T5 large (Chung

et al., 2022), each being prompted with the same two demonstrations to generate a rationale for all question-answer pairs in StrategyQA (examples shown in Appendix C.4).

Without further specification, we use the gradient attribution (Sundararajan et al., 2017) from Fast-Text model (Joulin et al., 2017) to identify label-leaking tokens to mask. We use T5-base (Raffel et al., 2020) to generate counterfactual data. Finally, two additional T5-base models are trained under Equation 7 and Equation 8 as evaluation models to compute the RORA score using Equation 9.[9]

We compare RORA with three strong baselines: LAS (Hase et al., 2020), RQ (Wiegreffe et al., 2020), and REV (Chen et al., 2023b). All baselines use T5-base as the evaluation model. Note that the boundaries of these metrics differ depending on the choice of the $H_\mathcal{V}$ function in Equation 1. Specifically, LAS and RQ, which utilize accuracy, lie within the $[-1, 1]$ interval. In contrast, RORA and REV, employing the conditional $\mathcal{V}$-entropy (Equation 2), have a theoretical boundary of $[-\infty, \infty]$. Because of different scales and evaluation criteria, our analysis mainly focuses on the ranking over different rationale types, as well as the relative difference within each metric. We also do an ablation study with (RORA$_{ablation}$) with all the attributions, counterfactual generation and invariant learning removed. This is equivalent to the formulation in Equation 3 in which the evaluator $\Phi$ is directly trained on the given rationales and questions to maximize the likelihood of labels.

### 4.1  Fixed-Label QA Task

We apply RORA to StrategyQA (Geva et al., 2021), a multi-step commonsense reasoning dataset with fixed label set {True, False}. We concatenate the provided list of relevant facts to create *gold* rationales. Notice that *gold* rationale created this way

---

[7]https://huggingface.co/domenicrosati/question_converter-3b

[8]For space concerns examples for ECQA are shown in Table 5 in Appendix C

[9]Note that the baseline evaluation model $\Theta$ does not condition on rationales, thus $\Theta$ can be trained once and the result of $\mathbb{E}_{(x,r,y)} \log p_\Theta(y|x,\phi)$ applies to the evaluation of various rationales types for the same dataset.

| Rationales → | | Synthetic Leaky | | | Model Generated | | | |
|---|---|---|---|---|---|---|---|---|
| Metrics ↓ | *gold* | *gold$_{leaky}$* | *vacuous* | *leaky* | GPT-4 | GPT-3.5 | Llama2-7B | Flan-T5 Large |
| RORA | 0.115 | 0.119 | 0.043 | 0.038 | 0.283 | 0.253 | 0.100 | 0.061 |
| RORA$_{DeBERTa}$ | 0.381 | 0.348 | 0.026 | 0.024 | 0.474 | 0.459 | 0.232 | 0.132 |
| RORA$_{ablation}$ | 0.121 | 0.673 | 0.505 | 0.673 | 0.501 | 0.412 | 0.155 | 0.102 |
| REV | -0.038 | 0.145 | -0.005 | 0.147 | -0.066 | -0.083 | -0.035 | -0.013 |
| LAS | 0.024 | NaN | 0.131 | NaN | 0.161 | 0.215 | 0.084 | 0.027 |
| RQ | 0.138 | 0.406 | 0.376 | 0.406 | 0.316 | 0.311 | 0.121 | 0.071 |
| Human Eval | – | – | – | – | 2.69 | 2.32 | 1.16 | 0.39 |

Table 2: Evaluation results on StrategyQA (Geva et al., 2021) over synthetic and model-generated rationales. **Both RORA and RORA$_{DeBERTa}$ rank *gold* and *gold$_{leaky}$* high and close to each other, while they rank *vacuous* and *leaky* close to zero, indicating strong invariance to various label leakages. Additionally, RORA and RORA$_{DeBERTa}$ consistently align the ranking of model-generated rationales with those of humans**.
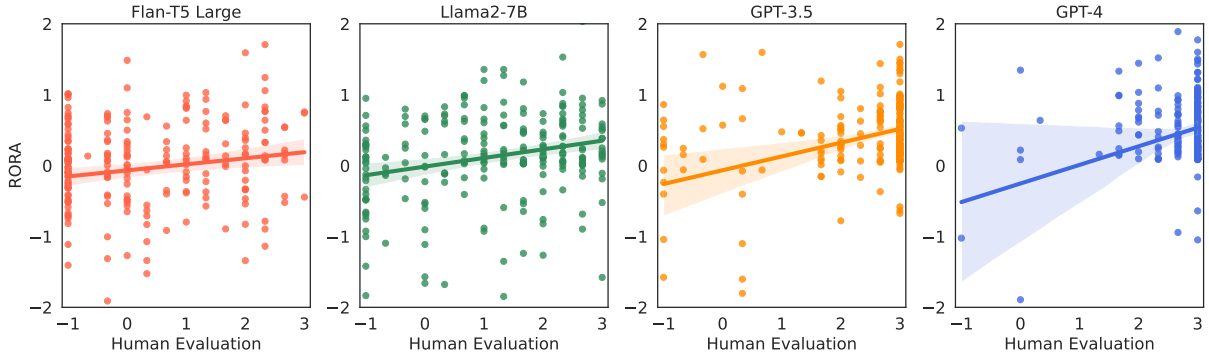


Figure 2: Linear regression on pointwise score correlation between human evaluation results and RORA scores. Shades correspond to a 95% confidence interval.

has all the essential facts but does not explicitly define each reasoning step to be taken.

Table 2 shows evaluation results on human-annotated, synthetic, and model-generated rationales. Evidently, only full pipeline RORA ranks *gold* and *gold$_{leaky}$* the highest, much higher than scores of *vacuous* and *leaky*. This aligns with the human intuition that rationales with new, useful information should be ranked higher while restating labels should not affect scoring. However, all existing baselines violate the intuition and assign much higher scores to three synthetic leaky rationales than the *gold* rationale. We note that LAS is undefined for *leaky* and *gold$_{leaky}$* rationales as it classifies all examples to leaked, thus causing a division-by-zero problem when averaging non-leaked examples, a problem discussed by Pruthi et al. (2022).

We observe similar trends with RORA$_{DeBERTa}$, which is a variation of RORA with the final two T5-base evaluation models replaced by DeBERTaV3-large (He et al., 2023) trained with Equation 6

directly. This indicates that RORA is model-agnostic, demonstrating comparable effectiveness when using different baseline models. On the other hand, RORA$_{ablation}$ is unable to provide comparable scores to *gold* and *gold$_{leaky}$*, yet it assigns the highest to *leaky*, which again highlights the effectiveness of the RORA approach in mitigating label leakage during evaluation.

**Human Evaluation** We conduct a human evaluation for free-text rationales on Amazon Mechanical Turk. We annotate all 229 instances in our custom test split of StrategyQA. For each instance, we annotate rationales generated by all four models. For each rationale, the annotator is first asked to provide a binary judgment on whether the rationale supports the target label. If the answer is positive, the annotators are asked to judge how much new information the rationale provides beyond restating the question and the answer. This yields a natural 5-point scale to which we give scores in $\{-1, 0, 1, 2, 3\}$ respectively. Each rationale is an-
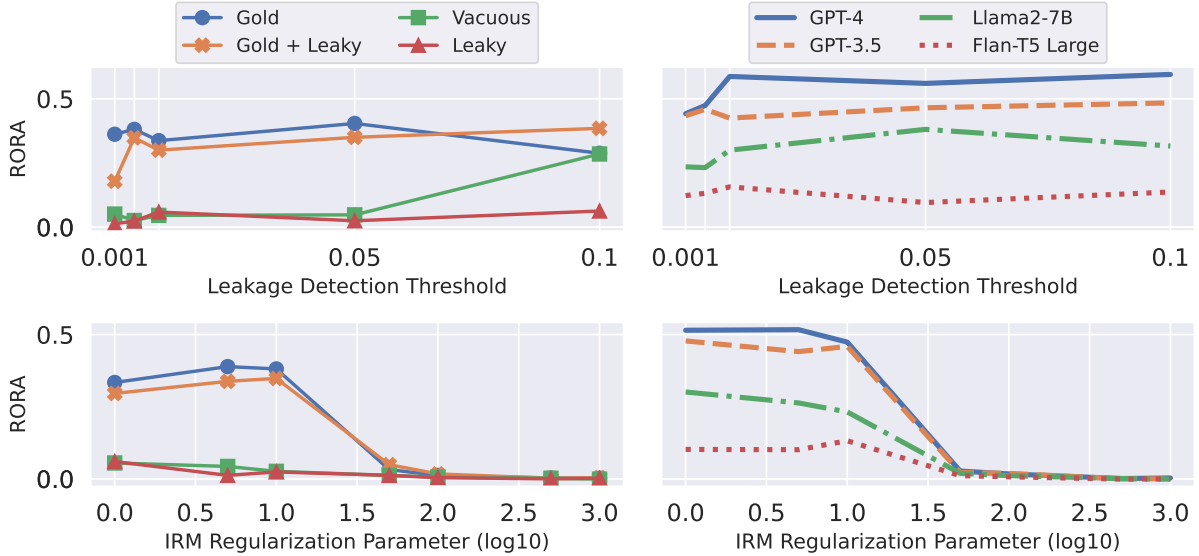
Figure 3: Sensitivity test results of RORA on leakage detection threshold and IRM regularization parameter. **Decreasing threshold and increasing IRM regularization parameter help RORA to better counteract label leakage. RORA appears to be stable when the parameter was chosen from a relatively wide range**.

| Method | Pearson's $r$ | Spearman's $\rho$ |
|---|---|---|
| REV (Chen et al., 2021) | 0.05 | 0.03 |
| RORA (Ours) | 0.36 | 0.37 |

Table 3: RORA achieves a higher correlation with instance-wise human judgments.

notated three-way redundant, and we average all human evaluation scores under the same model to obtain the final score for that particular rationale generation model. Additional human evaluation information is provided in Appendix C.3. The results are shown in Table 2.

While RQ and RORA$_{ablation}$ also produce similar rankings as humans that give higher scores to larger LLMs, RORA and RORA$_{DeBERTa}$ exhibit closer alignment with human evaluation and provide more fine-grained results. For instance, the RORA score of GPT-3.5 is slightly lower than GPT-4 ($0.283 \rightarrow 0.253$, $0.474 \rightarrow 0.459$), similar to human evaluation results ($2.69 \rightarrow 2.32$). Furthermore, pointwise $\mathcal{V}$-information estimated by RORA for each rationale also has a decent correlation with human judgments, as shown in Figure 2 and Table 3.

## 4.2 Sensitivity Analysis

Considering that RORA first detects leakage tokens based on a predetermined threshold and then leverages IRM to force the evaluation model to debiasing the label leakage, we study RORA's sen-

sitivity towards the choice of the masking threshold as well as the IRM regularization parameter $\lambda$ in Equation 7.

**Masking Threshold** We test RORA$_{DeBERTa}$ across five masking thresholds on StrategyQA, where a lower threshold results in more aggressive detection. Results shown in Figure 3 (first row) indicate that decreasing the threshold does induce stronger intervention, making the model ignorant of label leakage. For instance, the RORA score on *vacuous* rationale drastically declines to 0 when the threshold decreases from 0.1 to 0.01. The $gold_{leaky}$ rationale is initially higher than *gold*, but after detecting more leakage tokens, it falls slightly below *gold* and keeps having the same trend. Additionally, RORA consistently provides an expected ranking for model-generated rationales across different thresholds, demonstrating strong robustness for free-text rationale evaluation. As discussed in 3.2, we observe that RORA is relatively robust to the threshold because most accidental over-masking is likely to be filled with identical tokens, provided that the token does not have a strong correlation with the label.

**IRM Coefficient $\lambda$** The IRM regularization coefficient controls the strength of the invariant learning. A higher value forces the evaluation model to focus more on the shared features between the leaked and non-leaked environments while potentially sacrificing performance to the single environ-

ment. We test $\text{RORA}_{DeBERTa}$ over seven different regularization parameters ranging from 1 to 1000 on StrategyQA and results are shown in Figure 3 (second row). The increase of parameters from 1 to 20 leads to the convergence of lines *gold* and *gold*$_{leaky}$, while still being higher than *vacuous* and *leaky*. This demonstrates that invariant learning is helpful for evaluation models to mitigate label leakage in rationales. However, exceedingly large regularization parameters, such as those higher than 100, can make evaluation models unlearn features in leaked and non-leaked environments, resulting in declined scores close to 0. We also observe that RORA provides a correct ranking for model-generated rationales over a wide and reasonable range of regularization parameters.

### 4.3 Open-Label QA Task

We further test RORA on CommonsenseQA (Talmor et al., 2019), where we evaluate both ECQA (Aggarwal et al., 2021) and COS-E v1.11 (Rajani et al., 2019). CommonsenseQA adopts a multiple-choice setting where the label set for each problem differs from one another. The model usually needs contextualization to make meaningful predictions when the inputs are linearized into plain text. To have a minimally contextualized attribution model, we use LSTM (Hochreiter and Schmidhuber, 1997) bi-encoder model to separately encode the question and choices and select the final label by calculating the dot product. We perturb ECQA rationales to generate *gold*$_{leaky}$, and other leaky rationales are created similarly as in the previous task.

Table 4 shows that RORA again exhibits strong robustness to label leakages. It ranks *gold* and *gold*$_{leaky}$ rationales the highest and gives *vacuous* and *leaky* fairly low scores, while existing baselines are hard to achieve. Additionally, RORA also finds that the quality of CoS-E is inferior to that of ECQA, consistent with the evaluation results from other baselines (Chen et al., 2023b) and human observations (Aggarwal et al., 2021). While REV (Chen et al., 2023b) also ranks all synthetic rationales similarly to RORA, it is because it only takes the rationale as input, and on an open-label dataset, this singlehandedly does not deterministically leak the label. However, our further experiments reveal other potential problems with REV (Appendix B). On the contrary, RORA is not susceptible to such adversarial, as RORA evaluator already has access to everything in the question inputs.

| Rationales | Human-Annotated | | Synthetic Leaky | | |
|---|---|---|---|---|---|
| Metrics ↓ | CoS-E | *gold* | *gold*$_{leaky}$ | *vacuous* | *leaky* |
| RORA | 0.398 | 1.30 | 1.33 | 0.325 | 0.025 |
| REV | 0.128 | 0.284 | 0.271 | 0.013 | -0.010 |
| LAS | 0.198 | 0.365 | 0.373 | 0.425 | 0.506 |
| RQ | 0.235 | 0.425 | 0.432 | 0.402 | 0.432 |

Table 4: Evaluation results on ECQA (Aggarwal et al., 2021) over different rationales, including human-annotated rationales from CoS-E (Rajani et al., 2019). **RORA shows strong robustness in evaluating leaky rationales for open-label QA questions**.

## 5 Related Work

### 5.1 Gradient-based Attribution

Gradient-based attribution methods have been very popular in interpreting neural predictions. Comparative to perturbation (Zeiler and Fergus, 2014; Zintgraf et al., 2017) or Shapley-value-based attributions (Shapley, 1988; Lundberg and Lee, 2017; Strumbelj and Kononenko, 2010; Chen et al., 2020a), variants of gradient-based attribution (Sundararajan et al., 2017; Shrikumar et al., 2017) can be very efficient to compute as the number of samples needed to generate attribution does not depend on the number of features (Ancona et al., 2018; Lyu et al., 2024). Therefore, in this work, we use integrated gradient (IG) (Sundararajan et al., 2017) attribution for the small and small model for label leakage identification because IG is a generally applicable method that has been shown to provide good explanations for NLP tasks (Pruthi et al., 2022), while more easily accessible methods are either non-applicable due to non-linearity being employed in the computation process (Ancona et al., 2018), or requiring specific model structures while still generating questionable attributions (Jain and Wallace, 2019; Wiegreffe and Pinter, 2019; Pruthi et al., 2019; Bibal et al., 2022; Ethayarajh and Jurafsky, 2021; Liu et al., 2022).

### 5.2 Contrastive Generation

To better understand model decision boundary, Gardner et al. (2020) propose to evaluate models with contrast sets, which consist of test instances perturbed in small but meaningful ways that change the gold label. As human supervision in terms of contrast sets helps evaluate and improve classifiers (Kaushik et al., 2019; Chen et al., 2023d), a couple of methods have been proposed to automatically generate counterfactuals as a form of

explanation (Li et al., 2020; Yang et al., 2020; Wu et al., 2021; Dixit et al., 2022; Ross et al., 2022; Jacovi et al., 2021). In particular, SMG (Sha et al., 2021), MICE (Ross et al., 2021), and CREST (Treviso et al., 2023) adopt a select-and-edit two-step process to generate label-changing counterfactuals. In this work, we adapt MICE to generate partial counterfactuals, where only label-leaking parts of a rationale that spuriously correlate to the label get edited.

## 5.3 Invariant Learning

Invariant Risk Minimization (IRM) (Arjovsky et al., 2020; Ahuja et al., 2020) is a recently proposed framework for learning invariant predictors to spurious correlations. While most of these works primarily focused on theoretical findings, simple models, and toy datasets, in a notable departure Dranker et al. (2021) first tests IRM on natural language inference, pointing out that IRM has a potential edge over empirical risk minimization on tasks or datasets with stronger and more prevalent biases and larger data size. This work combines IRM with counterfactual data augmentation to train an evaluator agnostic of label leaking features.

## 6 Conclusions

We introduce RORA, an innovative approach designed to enhance the accuracy and robustness of rationale quality evaluation amidst the challenges posed by label leakage. Our method stands out for its adaptability across diverse model architectures and hyperparameter configurations, delivering consistent rationale quality evaluations that align more closely with human evaluation.

## Limitations

While RORA addresses the problem of label leakage, we focus on the evaluation of the informativeness of a rationale. While an easy-to-follow reasoning path in the rationale seems to boost evaluation score, RORA is not sensitive to other aspects of rationale quality. For example, a good rationale should also be factual (Prasad et al., 2023) and consistent (Chen et al., 2023a), which is not directly addressed in this work.

## Acknowledgement

## References

Shourya Aggarwal, Divyanshu Mandowara, Vishwajeet Agrawal, Dinesh Khandelwal, Parag Singla, and Dinesh Garg. 2021. Explanations for CommonsenseQA: New Dataset and Models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*.

Kartik Ahuja, Karthikeyan Shanmugam, Kush R. Varshney, and Amit Dhurandhar. 2020. Invariant risk minimization games. In *International Conference on Machine Learning* (ICML).

Marco Ancona, Enea Ceolini, Cengiz Öztireli, and Markus Gross. 2018. Towards better understanding of gradient-based attribution methods for deep neural networks. In *International Conference on Learning Representations* (ICLR).

Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. 2020. Invariant risk minimization. In *International Conference on Machine Learning* (ICML).

Ivan Aslanov and Ernesto Guerra. 2023. Tautological formal explanations: does prior knowledge affect their satisfiability? *Frontiers in Psychology*, 14.

William Bechtel and Adele Abrahamsen. 2005. Explanation: A mechanist alternative. *Studies in History and Philosophy of Science Part C: Studies in History and Philosophy of Biological and Biomedical Sciences*, 36(2):421–441.

Adrien Bibal, Rémi Cardon, David Alfter, Rodrigo Wilkens, Xiaoou Wang, Thomas François, and Patrick Watrin. 2022. Is attention explanation? an introduction to the debate. In *Annual Meeting of the Association for Computational Linguistics* (ACL).

Oana-Maria Camburu, Tim Rocktäschel, Thomas Lukasiewicz, and Phil Blunsom. 2018. e-snli: Natural language inference with natural language explanations. In *Advances in Neural Information Processing Systems* (NeurIPS).

Aaron Chan, Shaoliang Nie, Liang Tan, Xiaochang Peng, Hamed Firooz, Maziar Sanjabi, and Xiang Ren.

2022. Frame: Evaluating rationale-label consistency metrics for free-text rationales.

Angelica Chen, Jason Phang, Alicia Parrish, Vishakh Padmakumar, Chen Zhao, Samuel R Bowman, and Kyunghyun Cho. 2023a. Two failures of self-consistency in the multi-step reasoning of llms. *arXiv preprint arXiv:2305.14279*.

Hanjie Chen, Faeze Brahman, Xiang Ren, Yangfeng Ji, Yejin Choi, and Swabha Swayamdipta. 2023b. REV: Information-theoretic evaluation of free-text rationales. In *Annual Meeting of the Association for Computational Linguistics* (ACL). Association for Computational Linguistics.

Hanjie Chen, Guangtao Zheng, and Yangfeng Ji. 2020a. Generating hierarchical explanations on text classification via feature interaction detection. In *Annual Meeting of the Association for Computational Linguistics* (ACL).

Jifan Chen, Eunsol Choi, and Greg Durrett. 2021. Can NLI models verify QA systems' predictions? In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3841–3854, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Tongfei Chen, Zheng Ping Jiang, Adam Poliak, Keisuke Sakaguchi, and Benjamin Van Durme. 2020b. Uncertain natural language inference. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8772–8779.

Yanda Chen, Chen Zhao, Zhou Yu, Kathleen McKeown, and He He. 2023c. On the relation between sensitivity and accuracy in in-context learning.

Yanda Chen, Ruiqi Zhong, Narutatsu Ri, Chen Zhao, He He, Jacob Steinhardt, Zhou Yu, and Kathleen McKeown. 2023d. Do models explain themselves? counterfactual simulatability of natural language explanations. *arXiv preprint arXiv:2307.08678*.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.

Robert C. Cummins. 2000. "how does it work" versus "what are the laws?": Two conceptions of psychological explanation. In F. Keil and Robert A. Wilson, editors, *Explanation and Cognition, 117-145*. MIT Press.

Tanay Dixit, Bhargavi Paranjape, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. Core: A retrieve-then-edit framework for counterfactual data generation. In *Findings of the Association for Computational Linguistics: EMNLP*.

Yana Dranker, He He, and Yonatan Belinkov. 2021. Irm—when it works and when it doesn't: A test case of natural language inference. In *Advances in Neural Information Processing Systems* (NeurIPS).

Kawin Ethayarajh and Dan Jurafsky. 2021. Attention flows are shapley value explanations. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Short Papers)*.

Matt Gardner, Yoav Artzi, Victoria Basmova, Jonathan Berant, Ben Bogin, Sihao Chen, Pradeep Dasigi, Dheeru Dua, Yanai Elazar, Ananth Gottumukkala, et al. 2020. Evaluating models' local decision boundaries via contrast sets. In *Conference on Empirical Methods in Natural Language Processing* (EMNLP) - Findings.

Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. 2020. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673.

Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies. *Transactions of the Association for Computational Linguistics* (TACL).

Stuart Glennan. 2002. Rethinking mechanistic explanation. *Philosophy of science*, 69(S3):S342–S353.

Peter Hase, Shiyue Zhang, Harry Xie, and Mohit Bansal. 2020. Leakage-adjusted simulatability: Can models generate non-trivial explanations of their behavior in natural language? In *Findings of the Association for Computational Linguistics: EMNLP 2020*. Association for Computational Linguistics.

Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2023. Debertav3: Improving deberta using electra-style pre-training with gradient-disentangled embedding sharing.

John Hewitt, Kawin Ethayarajh, Percy Liang, and Christopher Manning. 2021. Conditional probing: measuring usable information beyond a baseline. In *Conference on Empirical Methods in Natural Language Processing* (EMNLP). Association for Computational Linguistics.

Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.

Alon Jacovi, Swabha Swayamdipta, Shauli Ravfogel, Yanai Elazar, Yejin Choi, and Yoav Goldberg. 2021. Contrastive explanations for model interpretability. In *Conference on Empirical Methods in Natural Language Processing* (EMNLP).

Sarthak Jain and Byron C Wallace. 2019. Attention is not explanation. In *Proceedings of NAACL-HLT*.

Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2017. Bag of tricks for efficient text classification. In *Proceedings of the 15th Conference of the European Chapter of the Association*

for Computational Linguistics: Volume 2, Short Papers, pages 427–431. Association for Computational Linguistics.

Divyansh Kaushik, Eduard Hovy, and Zachary Lipton. 2019. Learning the difference that makes a difference with counterfactually-augmented data. In *International Conference on Learning Representations* (ICLR).

Frank C Keil. 2006. Explanation and understanding. *Annu. Rev. Psychol.*, 57:227–254.

Chuanrong Li, Lin Shengshuo, Zeyu Liu, Xinyi Wu, Xuhui Zhou, and Shane Steinert-Threlkeld. 2020. Linguistically-informed transformations (lit): A method for automatically generating contrast sets. In *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*.

Xiang Lorraine Li, Adhiguna Kuncoro, Jordan Hoffmann, Cyprien de Masson d'Autume, Phil Blunsom, and Aida Nematzadeh. 2022. A systematic investigation of commonsense knowledge in large language models. In *Conference on Empirical Methods in Natural Language Processing* (EMNLP).

Yibing Liu, Haoliang Li, Yangyang Guo, Chenqi Kong, Jing Li, and Shiqi Wang. 2022. Rethinking attention-model explainability through faithfulness violation test. In *International Conference on Machine Learning* (ICML). PMLR.

Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization.

Josh Magnus Ludan, Yixuan Meng, Tai Nguyen, Saurabh Shah, Qing Lyu, Marianna Apidianaki, and Chris Callison-Burch. 2023. Explanation-based fine-tuning makes models more robust to spurious cues. In *Annual Meeting of the Association for Computational Linguistics* (ACL).

Scott M Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems* (NeurIPS).

Qing Lyu, Marianna Apidianaki, and Chris Callison-Burch. 2024. Towards faithful model explanation in nlp: A survey. *Computational Linguistics*.

Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. 2023. Self-refine: Iterative refinement with self-feedback. *arXiv preprint arXiv:2303.17651*.

Michael Neely, Stefan F. Schouten, Maurits Bleeker, and Ana Lucic. 2022. A song of (dis)agreement: Evaluating the evaluation of explainable artificial intelligence in natural language processing.

Rodrigo Nogueira, Zhiying Jiang, and Jimmy Li. 2021. Investigating the limitations of the transformers with simple arithmetic tasks. *arXiv preprint arXiv:2102.13019*.

OpenAI. 2023. GPT-4 Technical Report.

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training Language Models to Follow Instructions with Human Feedback. In *Advances in Neural Information Processing Systems* (NeurIPS).

Francis Jeffry Pelletier. 1994. The principle of semantic compositionality. *Topoi*.

Archiki Prasad, Swarnadeep Saha, Xiang Zhou, and Mohit Bansal. 2023. Receval: Evaluating reasoning chains via correctness and informativeness. *arXiv preprint arXiv:2304.10703*.

Danish Pruthi, Rachit Bansal, Bhuwan Dhingra, Livio Baldini Soares, Michael Collins, Zachary C Lipton, Graham Neubig, and William W Cohen. 2022. Evaluating explanations: How much do explanations from the teacher aid students? *Transactions of the Association for Computational Linguistics* (TACL).

Danish Pruthi, Mansi Gupta, Bhuwan Dhingra, Graham Neubig, and Zachary C Lipton. 2019. Learning to deceive with attention-based explanations. *arXiv preprint arXiv:1909.07913*.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research* (JMLR).

Nazneen Fatema Rajani, Bryan McCann, Caiming Xiong, and Richard Socher. 2019. Explain yourself! leveraging language models for commonsense reasoning. In *Annual Meeting of the Association for Computational Linguistics* (ACL).

Alexis Ross, Ana Marasović, and Matthew Peters. 2021. Explaining NLP models via minimal contrastive editing (MiCE). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP*. Association for Computational Linguistics.

Alexis Ross, Tongshuang Wu, Hao Peng, Matthew Peters, and Matt Gardner. 2022. Tailor: Generating and perturbing text with semantic controls. In *Annual Meeting of the Association for Computational Linguistics* (ACL).

Lei Sha, Patrick Hohenecker, and Thomas Lukasiewicz. 2021. Controlling text edition by changing answers of specific questions. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP*.

Lloyd S. Shapley. 1988. *A value for n-person games*, page 31–40. Cambridge University Press.

Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. 2017. Learning important features through propagating activation differences. In *International Conference on Machine Learning* (ICML). PMLR.

Suzanna Sia, Anton Belyy, Amjad Almahairi, Madian Khabsa, Luke Zettlemoyer, and Lambert Mathias. 2022. Logical satisfiability of counterfactuals for faithful explanations in nli. In *Conference on Artificial Intelligence* (AAAI).

Erik Strumbelj and Igor Kononenko. 2010. An efficient explanation of individual classifications using game theory. *Journal of Machine Learning Research* (JMLR).

Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks. In *International Conference on Machine Learning* (ICML).

Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

Marcos Treviso, Alexis Ross, Nuno M. Guerreiro, and André Martins. 2023. CREST: A joint framework for rationalization and counterfactual text generation. In *Annual Meeting of the Association for Computational Linguistics* (ACL). Association for Computational Linguistics.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems* (NeurIPS), 35:24824–24837.

Sarah Wiegreffe, Ana Marasović, and Noah A. Smith. 2020. Measuring association between labels and free-text rationales. In *Conference on Empirical Methods in Natural Language Processing* (EMNLP).

Sarah Wiegreffe and Yuval Pinter. 2019. Attention is not not explanation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*.

Robert A Wilson and Frank Keil. 1998. The shadows and shallows of explanation. *Minds and machines*, 8:137–159.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, R'emi Louf, Morgan Funtowicz, and Jamie Brew. 2019. Huggingface's transformers: State-of-the-art natural language processing. *ArXiv*, abs/1910.03771.

Tongshuang Wu, Marco Tulio Ribeiro, Jeffrey Heer, and Daniel Weld. 2021. Polyjuice: Generating counterfactuals for explaining, evaluating, and improving models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Association for Computational Linguistics.

Yilun Xu, Shengjia Zhao, Jiaming Song, Russell Stewart, and Stefano Ermon. 2020. A theory of usable information under computational constraints. In *International Conference on Learning Representations* (ICLR).

Linyi Yang, Eoin Kenny, Tin Lok James Ng, Yi Yang, Barry Smyth, and Ruihai Dong. 2020. Generating plausible counterfactual explanations for deep transformers in financial text classification. In *International Conference on Computational Linguistics* (COLING).

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. In *Advances in Neural Information Processing Systems* (NeurIPS).

Matthew D Zeiler and Rob Fergus. 2014. Visualizing and understanding convolutional networks. In *The European Conference on Computer Vision* (ECCV).

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. In *Advances in Neural Information Processing Systems* (NeurIPS).

Luisa M Zintgraf, Taco S Cohen, Tameem Adel, and Max Welling. 2017. Visualizing deep neural network decisions: Prediction difference analysis. *arXiv preprint arXiv:1702.04595*.

# Supplemental Material

| Appendix | Contents |
|---|---|
| Appendix A | Justifying Partial Counterfactual Augmentation |
| Appendix B | Further Discussion on REV and Vacuous Rationale |
| Appendix C | Additional Details of Experimental Setup |

## A  Justifying Partial Counterfactual Augmentation

This section provides intuitive explanations of why other simple alternatives may not work. We consider the following alternatives:

- Mask out all the tokens attributed as leakage by the smaller model.

- IRM training with the environment $D$ and $D^{\text{Mask}}$, which we denote by $e_1$ and $e_2$ respectively.

### A.1  Training on Masked Inputs

The major problem is that masking may introduce new artifacts. For instance, in English, negation is applied to the auxiliary verb. Therefore, if the model identifies a discrepancy in tense and a masked token preceding the main verb, it's highly probable that the masked token represents negation. In our pilot experiments, we found that models are capable of utilizing these traits. Also, attribution calculated from positive and negative labels might be different, a pattern the model might learn to exploit. That's why in RORA, the attribution is calculated with the global weighting. Partially counterfactual generation makes the sentence more natural, which helps reduce artifacts introduced.

### A.2  Using Masked Inputs as an Environment

We hereby show that invariant learning with a masked environment does not force the model to unlearn $R_L$, as described in Equation 5. In general, creating scenarios where one environment includes a non-essential characteristic and the other does not will not result in an Invariant Risk Minimization (IRM) model that achieves comparable loss levels regardless of whether these non-essential characteristics exist. To demonstrate this, we create a minimal example from the causal graph from LAS (Hase et al., 2020), where each simulation variable is replaced with its actual variable as shown in Figure 4.
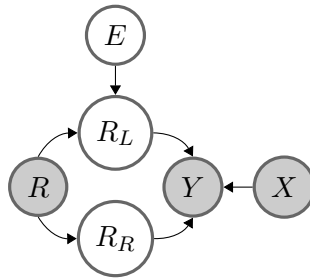
Figure 4: A causal graph showcasing the generative story of the label $Y$ of an instance, where correctness indicator variables from LAS (Hase et al., 2020) are replaced by an actual variable.

Each variable corresponds to its designation in §2.2. For simplicity, let us consider the following semantics of each random variable: $R$ is a free-text rationale taking values from all possible rationale documents; $R_L$ is a ternary variable taking values from $r_l \in \{\text{True}, \text{False}, \text{UNK}\}$, where UNK denotes masks; $R_R$ is a binary indicator of $r_r \in \{0, 1\}$; and $X$ is the binary indicator of $x \in \{0, 1\}$ indicating the binary choice innate to the original question. We assume the following conditional probability:

$$p(Y = \text{True} \mid R_L = \text{True}, X, R_R) = 1,$$
$$p(Y = \text{True} \mid R_L = \text{False}, X, R_R) = 0,$$
$$p(Y = \text{True} \mid R_L = \text{UNK}, X = R_R) = 0.9,$$
$$p(Y = \text{True} \mid R_L = \text{UNK}, X \neq R_R) = 0.$$

In this case, the relationship $R_L \Rightarrow Y$ *is causal*, in the sense that the model can discern which environment the prediction is made in by looking at $R_L$. Specifically, assuming that we can identify the leakage part correctly with the application of IG (Sundararajan et al., 2017), then the environments will be

$$e_1 := \Big\{ (R^{e_1}, X^{e_1}, Y^{e_1}) \mid p(R_L^{e_1} = \text{UNK} \mid R^{e_1}) = 1 \Big\},$$
$$e_2 := \Big\{ (R^{e_2}, X^{e_2}, Y^{e_2}) \mid p(R_L^{e_2} = \text{UNK} \mid R^{e_2}) = 1 \Big\}.$$

Now suppose we have infinite computation, and it is possible to learn any representation $\Phi : \mathcal{R} \times \mathcal{X} \to \mathbb{R}^d$, and a linear classification head $\tilde{w} : \mathbb{R}^d \to \mathcal{P}(\mathcal{Y})$ (which outputs a single value between $[0, 1]$). We are considering the risk function $F$ that is common in classification settings:

$$F(\tilde{w}, \Phi) = \mathbb{E}\Big[ -\mathbb{I}[y = \text{True}] \cdot \log\big(\tilde{w} \circ \Phi(x)\big) - \mathbb{I}[y = \text{False}] \cdot \log\big(1 - \tilde{w} \circ \Phi(x)\big) \Big].$$

Suppose we have a representation that fully recovers all the causal components

$$\Phi(r, x) = \Big\langle \mathbb{I}[x = q] \cdot \mathbb{I}[r_l = \text{UNK}], \mathbb{I}[r_l = \text{True}] \Big\rangle,$$

Then it is easy to see that $\tilde{w}(z) = v^t z$ where $v = \langle 0.9, 1 \rangle$ will be the optimal classifier across both environments $e \in \{e_1, e_2\}$. It can be verified that

$$\lim_{v' \to v} \nabla_{w \mid w=1.0} F(\tilde{w}, \Phi) = -p(Y = \text{True}) \frac{1}{v^T \Phi(x)} + p(Y = \text{False}) \frac{1}{1 - v^T \Phi(x)}.$$

Which, given our $\tilde{w}$ and $\Phi$, does not penalize the predictions. But does this invariant predictor satisfy our desirable property that it will rank instances disregarding whether there is a label leakage? The answer is no, as here, the instances with label leakage will get zero losses, while those without will get non-zero losses. This highlights the importance of generating a counterfactual environment with natural-looking data with a different label association $p(Y \mid R_L)$, like what we did for RORA.

## B Further Discussion on REV and Vacuous Rationale

REV (Chen et al., 2023b) evaluates a rationale based on a vacuous rationale $B$. We argue that it does not follow the natural way of how humans evaluate a rationale, where we typically need to read and understand the question before the evaluation. Therefore, an adjusted version of REV would be

$$\text{REV}'(R) = I_{\mathcal{V}}(X, R \to Y \mid B).$$

However, the compositionality nature of language (Pelletier, 1994) implies that in many cases

$$\begin{aligned} I_{\mathcal{V}}(X, R \to Y \mid B) - I_{\mathcal{V}}(R \to Y \mid B) \neq \\ I_{\mathcal{V}}(X, R \to Y \mid \varnothing) - I_{\mathcal{V}}(R \to Y \mid \varnothing). \end{aligned} \tag{10}$$

This reveals another potential weakness of REV, which makes restating the question an incredibly high-scoring rationale (0.436 according to our experiment). Indeed, this is counter-intuitive as human evaluators should always have access to $X$ but not $B$. On the contrary, RORA does not have this issue as it already conditions on $X$.

## C    Experiment Setup

### C.1    Models and Datasets

We access all models through Huggingface Transformers (Wolf et al., 2019). We split the StrategyQA (Geva et al., 2021) train set into our custom train, validation, and test sets and then create model-generated rationales for each set. For ECQA (Aggarwal et al., 2021) we use their original splits.

### C.2    Training Hyperparameters

We train each small model for gradient attribution calculation with a learning rate of $1e - 1$ and batch size of 256. For leakage token detection, we choose threshold $0.005$ for RORA$_{DeBERTa}$ and $0.01$ for RORA, or use top-1 strategy. To determine the appropriate threshold, we assess the output on a small sample of training data, ensuring that at least the majority of the samples receive more than one masking beyond that from the top-k selection. We train each counterfactual data generator using a learning rate of $1e - 4$ with batch sizes 8. For invariant learning, we use IRM regularization $\lambda = 10$ for both RORA and RORA$_{DeBERTa}$ and train each evaluation model with a learning rate of $1e - 4$ and batch size of 64. For each training, we use optimizer AdamW (Loshchilov and Hutter, 2019) from Transformers (Wolf et al., 2019) and set maximum epochs as 20 with early stopping.

### C.3    Human Evaluation Details

We conduct a human evaluation of model-generated rationale on Amazon Mechanical Turk[10]. The quality of each rationale is judged by three individual annotators, who must pass a qualification test described in (Chen et al., 2020b). Annotators have been paid $0.05 per rationale, which amounts to an hourly wage of about $10. Figure 5 shows an example of our annotation interface. For each rationale, We present annotators with a question, an answer, and the model-generated rationale (explanation) and ask them the following questions:

1. *Does the Explanation justify the given Answer to the Question?* We offer two options for them to choose from: "Yes" and "No".

2. If they answer "Yes", we further ask them *how much additional information does the Explanation have to justify the Answer beyond just reiterating what is stated in Question and Answer?* and to pick one option from "No Additional Info", "Little Additional Info", "Some Additional Info" and "Sufficient Additional Info" which are corresponding to a 4-point Likert-scale (0/1/2/3). If the annotator answers "No" for the first question, we label the rationale a score of -1.

Then, we average all three annotators' scores to get the human evaluation score for that single rationale. For detailed instructions, readers can refer to Figure 6.

### C.4    Model Generated Rationales

Table 6 shows examples of rationale generated by different models.

---

[10] https://www.mturk.com/

**Question:** Is Shakespeare famous becaues of the infinitive form?
**Answer:** The statement is true.
**Explanation:** Shakespeare wrote the play Hamlet. Hamlet contains one of Shakespeare's most famous passages, Hamlet's soliloquy. Hamlet's soliloquy begins with the line 'To be or not to be', which uses the infinitive form.

---

**Q1:** Does the **Explanation** justify the given **Answer** to the **Question**?

| YES | NO |

**Q2:** If yes, how much additional information does the **Explanation** have to justify the **Answer** beyond just reiterating what is stated in **Question** and **Answer**?

| NO additional info | LITTLE additional info | SOME additional info | SUFFICIENT additional info |

SUBMIT

Figure 5: The interface to our human validation hit.

Thanks for participating in this HIT qualifier! Please read the examples below and answer the questions (1-2 questions) to the best of your ability.

**Main Instructions:** you will read a question that usually requires a combination of knowledge about different entities to solve. For each question, an answer and a statement explaining the answer will be provided.

Your task is to read the instance and judge the explanation quality (1-2 questions). Regarding quality, you need to assess two aspects: (1) whether the explanation **supports the label**; (2) whether the explanation contains **additional information** for describing the answer beyond simply combining the question and the answer.

To be more specific:

- **Supports the label** means the explanation is describing something related to the answer to the question (e.g., Example #2 and #3 below), rather than something else (e.g., Example #1 below), and reading it increase your belief in the correctness of the answer.
- **Additional information** means that the explanation provides additional evidence or background knowledge to support the answer (e.g., Example #3 below), rather than a simple paraphrase of the answer (e.g., Example #2 below). You only evaluate the **additional information** when you agree that the explanation **supports the label** first. Please select the **amount of additional information** (from "No additional info" to "Enough additional info") based on the extent to which the explanation helps you understand why the answer might have been chosen for the question.

An instance contains 3 parts:

| **Question** | A question such as "Is Shakespeare famous because of the infinitive form?" |
|---|---|
| **Answer** | An answer to the question as "The statement is true." (may or may not be correct) |
| **Explanation** | A **statement** which explains the **Answer** |

**Notes:**

- **Important!** The **Answer** **could be wrong based on your real world knowledge**, and the **Explanation** might contain grammatical mistakes. Please **ignore** these when answering the questions.
- You should judge whether the **Explanation** is describing the **Answer** to the **Question** and contains additional information to help you understand why the **Answer** is correct to the **Question**.
- You might see multiple HITs with the same **Question** and **Answer** but different **Explanations**. It's okay. Please evaluate each HIT independently.

Figure 6: The detailed instruction we give the annotators. Notice that we specifically instructed them to not judge the factuality based one their real-world knowledge.

| Dataset | Question | *gold* | *leaky* | *gold_leaky* | *vacuous* |
|---------|----------|--------|---------|--------------|-----------|
| ECQA | What town is known for disco? Options: op_1 rio de janeiro op_2 new york op_3 town op_4 germany op_5 hotel | Rio de Janeiro is a huge seaside city in brazil, famed for its Copacabana and Ipanema beaches. The town Rio de Janeiro is known for disco. New York city is famous for Central park, Statue of liberty, Empire state building and other such places and not for disco. Town is larger than a village and generally smaller than a city. Town is not the name of a place where disco is famous for. Germany is the country and not a town known for disco. Hotel is not the town known for disco. | The answer is Rio de Janeiro. | Rio de Janeiro is a huge seaside city in brazil, famed for its Copacabana and Ipanema beaches. The town Rio de Janeiro is known for disco. New York city is famous for Central park, Statue of liberty, Empire state building and other such places and not for disco. Town is larger than a village and generally smaller than a city. Town is not the name of a place where disco is famous for. Germany is the country and not a town known for disco. Hotel is not the town known for disco. The answer is Rio de Janeiro. | Rio de Janeiro is known for disco. |

Table 5: Example of human-annotated and synthetic leaky rationales for open-label QA task. *leaky* verbatim leaks the label, *vacuous* is the declarative combination of the question and the given answer, and *gold_leaky* is simply the concatenation of the *gold* and *leaky* rationales.

| Model | Question and Label | Generated Rationale |
|-------|--------------------|--------------------|
| GPT-4 | Question: Could the Great Wall of China connect the Dodgers to the White Sox? Answer: True | The Great Wall of China is 1,948 km long. The Yankees are based in New York City. The Dodgers are in Los Angeles and play in Los Angeles. The White Sox are based in Chicago and play in Chicago. The answer: yes. |
| GPT-3.5 | Question: Could the Great Wall of China connect the Dodgers to the White Sox? Answer: True | No rationale can be provided to explain this answer because it is not accurate. The Great Wall of China is located in China and does not connect the Dodgers (a baseball team in Los Angeles) to the White Sox (a baseball team in Chicago). |
| Llama2-7B | Question: Could the Great Wall of China connect the Dodgers to the White Sox? Answer: True | The Great Wall of China is a long wall. The Dodgers and White Sox are baseball teams. |
| Flan-T5 Large | Question: Could the Great Wall of China connect the Dodgers to the White Sox? Answer: True | The Great Wall of China is 1,948 km long. The Yankees are based in New York City. The Dodgers are in Los Angeles and play in Los Angeles. The White Sox are based in Chicago and play in Chicago. The answer: yes. |

Table 6: Examples of model-generated rationales