

Do LLM Self-Explanations Help Users Predict Model Behavior? Evaluating Counterfactual Simulatability with Pragmatic Perturbations

Anonymous ACL submission

Abstract

Large Language Models (LLMs) can produce verbalized self-explanations, yet prior studies suggest that such rationales may not reliably reflect the model’s true decision process. We ask whether these explanations nevertheless help users predict model behavior, operationalized as *counterfactual simulatability*. Using StrategyQA, we evaluate how well humans and LLM judges can predict a model’s answers to counterfactual follow-up questions, with and without access to the model’s chain-of-thought or post-hoc explanations. We compare LLM-generated counterfactuals with pragmatics-based perturbations as alternative ways to construct test cases for assessing the potential usefulness of explanations. Our results show that self-explanations consistently improve simulation accuracy for both LLM judges and humans, but the degree and stability of gains depend strongly on the perturbation strategy and judge strength. We also conduct a qualitative analysis of free-text justifications written by human users when predicting the model’s behavior, which provides evidence that access to explanations helps humans form more accurate predictions on the perturbed questions.

1 Introduction

Large language models (LLMs) can generate *verbalized self-explanations* that accompany a model’s decision. These self-explanations are increasingly used as a window into how the model “reasons”. Meanwhile, a growing body of work has raised concerns about the *faithfulness* of such explanations, showing that these rationales can be decoupled from the underlying decision process, or manipulated without affecting the model’s predictions, suggesting that self-explanations should not be trusted as accounts of model behavior (Atanasova et al., 2023; Madsen et al., 2024). This leaves open a question for explainable AI: even if LLM verbalized self-explanations are not strictly faithful, can

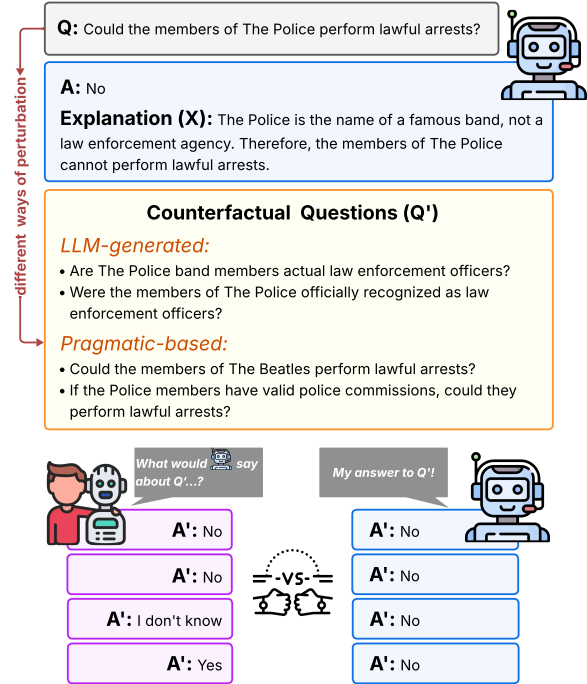


Figure 1: Pipeline for evaluating the usefulness of LLM self-explanations based on different counterfactual question perturbations. For each original Q and model prediction, we measure users’ ability to predict model behavior on different counterfactual question sets Q' , with and without access to the model’s self-explanations.

they nonetheless *contribute to model understanding* by helping users form more accurate predictions about how the model will behave? Recent works have proposed to assess explanation quality through *counterfactual simulatability*: an explanation is considered useful if it helps people anticipate how a model’s prediction would change on counterfactual inputs (Doshi-Velez and Kim, 2017; Chen et al., 2024; Limpijankit et al., 2025).

Natural language counterfactuals, defined as minimal *what-if* variations of an input (Miller, 2019; Ross et al., 2021; Pearl, 2021), are widely used in explainable AI to reveal how models behave under controlled perturbations (Chen et al.,

2024; Wang et al., 2024, 2025). In a counterfactual simulatability protocol, users see an original question Q and the model’s prediction A , then predict the model’s outputs on counterfactual variants Q' . The usefulness of an explanation is measured by how much it improves the simulation accuracy.

While this line of work offers a promising operationalization of explanation usefulness, much less attention has been paid to the *counterfactuals themselves*. We still lack a systematic understanding of how different ways of perturbing a question affect the conclusions about the usefulness of explanations. In Chen et al. (2024)’s setup, for example, the counterfactual questions are generated by LLMs using prompt-based designs. If the evaluation protocol is sensitive to the choice of perturbation, the measured *usefulness* of an explanation may largely be an artifact of how those counterfactuals are constructed.

Figure 1 outlines our pipeline. We extend the *counterfactual simulatability* framework to assess how LLM verbalized self-explanations support model understanding. We study *when* and *how* explanations help under different question perturbations. Specifically, we generate counterfactual variants Q' of an original question Q and measure explanation usefulness by how much explanations help users predict the model’s behavior on Q' . Perturbations that yield larger gains form a more informative testbed for evaluating self-explanations.

To summarize, our paper addresses the following research questions: (i) Do LLM verbalized self-explanations contribute to model understanding as measured by counterfactual simulatability? (ii) How do different counterfactual perturbation strategies assess the usefulness of LLM self-explanations and the difference between conditions with and without explanations? (iii) Are users biased toward the model’s original answer when they reason about Q' , and can self-explanations reduce this bias? (iv) How do self-explanations affect human users’ confidence and accuracy in predicting model behavior, and to what extent do LLM-as-a-judge predictions align with human judgments?

2 Related Work

Evaluating Explanations Simulatability. *Simulatability* measures how well humans can predict a model’s outputs based on its explanations (Doshi-Velez and Kim, 2017; Ribeiro et al., 2016; Chandrasekaran et al., 2018; Hase and Bansal, 2020;

Chen et al., 2024). *Faithfulness*, as a common metric for explanations, measures whether an explanation is consistent with the model’s decision process (Gilpin et al., 2018; Alvarez-Melis and Jaakkola, 2018; Jacovi and Goldberg, 2020). Following Chen et al. (2024), simulatability can be seen as a special case of faithfulness, where the output predictor is restricted to humans or LLM judges instead of any arbitrary black box model. In our work, we focus on simulatability as a way to capture the usefulness of explanations (Doshi-Velez and Kim, 2017; Hoffman et al., 2018).

Different from Chen et al. (2024), who run the simulatability setting once with explanations, we conduct a forward simulation experiment with two phases, following the experimental setup of the previous work, where they test whether users can better predict a detector’s behavior after seeing explanations (Hase and Bansal, 2020; Doshi-Velez and Kim, 2017; Schoenegger et al., 2024). We run experiments with both human users and LLM judges to evaluate how useful the model’s own verbal explanations are for predicting its behavior.

Counterfactual Generation. Counterfactual reasoning is the process of thinking about what would have happened if conditions were different, and it is a common tool for judging causality (Kahneman and Tversky, 1982). Such causal judgments are important for model evaluation, error analysis, and explanation (Miller, 2019). In NLP, most work defines a relationship between an original input x and a changed version \hat{x} , then constructs \hat{x} to follow this relationship (Wu et al., 2021). For example, Gardner et al. (2020) ask dataset authors to manually edit test examples in small but meaningful ways that usually flip the gold label.

However, manual counterfactuals are costly, so later work turns to model-based generation. Wu et al. (2021) fine-tune GPT-2 (Radford et al., 2019) to produce counterfactual examples given a desired edit type. Wang et al. (2025) propose FITCF, a framework that filters counterfactuals by checking whether the label flips, then uses them as demonstrations for few-shot prompting, achieving better results than three strong baselines. In Chen et al. (2024), which is closest to our setup, LLMs are prompted to generate counterfactual questions for **StrategyQA**. In contrast, we compare several counterfactual generation methods and propose a more controlled, pragmatics-based editing approach, and we compare it with counterfactual questions pro-

duced by generative models.

3 Pragmatic Counterfactual Taxonomy

As in Table 1, Grice characterizes effective communication in terms of four maxims—*Quantity*, *Quality*, *Relation*, and *Manner*—that regulate how speakers cooperate in interaction (Grice, 1975).

Maxim	Description
<i>Quantity</i>	Make your contribution as informative as required for the current conversational goal, but not more informative than necessary.
<i>Quality</i>	Ensure that your contribution is truthful: do not say what you believe to be false or lack adequate evidence for.
<i>Relation</i>	Make your contribution relevant to the conversational context and the question under discussion.
<i>Manner</i>	Be clear and orderly in how you express your contribution; avoid unnecessary ambiguity.

Table 1: The four Gricean maxims that underlie our pragmatic counterfactual taxonomy.

These maxims capture implicit expectations about how much information to provide, how truthful and well-founded it should be, how relevant it is to the current goal, and how clearly it is expressed (Krause and Vossen, 2024; Ma et al., 2025). As NLP systems are increasingly used in interactive and explanation-heavy settings, it is natural to ask whether the explanations they provide remain robust when these pragmatic dimensions are different (Jacquet et al., 2019; Kaczmarek-Majer et al., 2022; Alexandris, 2024; Zuo et al., 2025)

We therefore take the Gricean maxims as a principled basis for designing counterfactual question types. Each counterfactual corresponds to a controlled manipulation of one or more maxims. Table 2 summarizes our perturbation methods, together with linguistic motivations and examples.

4 Experimental Design

4.1 Task and Dataset

We use the **StrategyQA** (Geva et al., 2021) dataset as our primary QA task. StrategyQA is an English question answering benchmark of 2,780 short, open-domain yes/no questions that require implicit factual reasoning, often illustrated by examples like *Did Aristotle use a laptop?*. Each example is annotated with a decomposition into reasoning steps and a supporting Wikipedia paragraph.

4.2 Counterfactual Question Generation

LLM automatic counterfactual generation.

Following Chen et al. (2024), we prompt strong LLMs to generate counterfactual variants of each StrategyQA question. These serve as a *non-pragmatic* baseline. We adopt the same LLMs as Chen et al. (2024) (GPT-3.5 (Brown et al., 2020; Ouyang et al., 2022) and GPT-4 (OpenAI et al., 2023)) to ensure comparability of results, and additionally include Llama-3.3-70B-Instruct ("Llama-3.3") (Grattafiori et al., 2024) as a strong open-source model. The prompts used for automatic counterfactual generation are listed in Appendix A.

Pragmatics-based counterfactual generation.

For generating pragmatics-based counterfactual questions in **StrategyQA**, we follow the taxonomy in Table 2. In **StrategyQA**, each question is paired with a term (from which the question is derived), a natural language description of that term, and a set of implicit facts needed to answer the question (Geva et al., 2021). We treat the term as a *pragmatic anchor*—a lexical or conceptual element that carries implicit assumptions, invites alternative interpretations, or structures the reasoning process. We then construct five types of pragmatic counterfactual transformations on top of these anchors:

- **Presupposition flip:** we prompt GPT-4 to generate counterfactual variants of the form “*If A is not B but C, ...*”, thereby flipping or challenging the presuppositions associated with the anchor term. The exact prompts are provided in Appendix B.
- **Lexical substitution:** we replace the key term with a synonym or hypernym retrieved from WordNet (Miller, 1992), altering the lexical realization while preserving closely related meaning.
- **Scalar adjustment:** we substitute quantifiers with stronger or weaker scalar alternatives, using manually specified scalar scales listed in Appendix C. The scales are defined on the basis of the scalar inventories discussed in previous works (Carston, 1998; Papafragou and Musolino, 2003; Sauerland, 2004).
- **Contextualization:** we append one or two supporting facts from the dataset to enrich the question context and make part of the originally implicit pragmatic content explicit.

Counterfactual Type	Pragmatic Maxim	Linguistic Motivation	Example
Presupposition Flip	<i>Quality + Relation</i>	Test truthfulness of assumed background knowledge; revise false presuppositions	<i>Q</i> : Did the king of France visit Berlin? <i>Q'</i> : If there were no king of France but a president instead, would the president visit Berlin?
Lexical Substitution	<i>Relation + Quantity</i>	Check whether changes in lexical strength affect inference relevance	<i>Q</i> : Could a cheetah outrun a car? <i>Q'</i> : Could a cheetah outrun a <i>sports car</i> ?
Scalar Adjustment	<i>Quantity</i>	Examine model sensitivity to implicature scales (some ↔ most ↔ all)	<i>Q</i> : Do all birds fly? <i>Q'</i> : Do <i>some</i> birds fly?
Contextualization	<i>Relation + Manner</i>	Add missing contextual information that affects interpretation	<i>Q</i> : Is this animal a good pet? <i>Q'</i> : Given that this animal is a large wild tiger, is it a good pet?

Table 2: Pragmatic counterfactual question types Q' , their associated Gricean maxims, linguistic motivations, and illustrative examples.

4.3 Self-explanations Collection

For each original question Q and its model answer A , we collect two types of explanations from the same set of models used for LLM automatic counterfactual generation in §4.2. In our study, we consider two types of verbalized self-explanations and generate one of each type for every instance in the dataset (see Appendix D for prompt templates):

Chain-of-thought (CoT) explanations, where the model reveals an explicit step-by-step reasoning process (Nye et al., 2021). We elicit CoT explanations using *Let’s think step by step* (Kojima et al., 2022) together with additional instructions about formatting the final answer.

Post-hoc explanations, where the model is asked to justify its produced answer without modifying it (Camburu et al., 2018; Park et al., 2018).

4.4 Counterfactual Simulation

We use a counterfactual simulation test to measure the ability of human participants and LLM simulators to predict model behavior, following the setup of Hase and Bansal (2020) and Chen et al. (2024). Our counterfactual simulation design is illustrated in Figure 2. In previous work, counterfactual simulation procedures have been used to assess explanation usefulness by testing whether an explanation enables a judge to infer a model’s outputs on counterfactual variants of the input.

Our primary summary statistic is **acc_improvement**, which we use to compare conditions. Intuitively, **OverallAcc** measures how well an evaluator can predict the model’s follow-up answer within a phase, while **acc_improvement** captures the additional benefit of providing the model’s self-explanation by comparing Phase 2

(with explanation) against Phase 1 (without explanation). We define:

$$\Delta\text{OverallAcc} = \text{OverallAcc}^{\text{P2}} - \text{OverallAcc}^{\text{P1}}.$$

A positive $\Delta\text{OverallAcc}$ indicates that access to explanations improves model-behavior prediction, whereas values near zero suggest limited impact (e.g., when Phase 1 accuracy is already high), and negative values indicate that explanations can occasionally mislead the evaluator.

In our setup, we additionally use the simulation test to study the effect of different counterfactual question sets. Holding the explanation set fixed, we ask which counterfactual set yields the largest **acc_improvement**.

LLM-as-a-judge simulation. Building on prior work showing that LLMs approximate human judgments with high agreement and nuance, we adopt an *LLM-as-a-judge* setup for our counterfactual study (Zheng et al., 2023; Bai et al., 2023; Li et al., 2025). For each item, we pair a starter question and the model’s answer with a counterfactual follow-up. The judge then decides whether the starter answer, optionally with its explanation, allows inferring the model’s answer to the follow-up. To test robustness, we follow the scalable-oversight setup of Kenton et al. (2024) and use two independent judges with different capacities as **Weaker Evaluator** and **Stronger Evaluator**.

Human simulation user study. To assess how well our LLM-as-a-judge evaluations align with real human behavior, we conduct a human-simulation user study on a subset of GPT-4-generated explanations (both CoT and post-hoc), since it is the most advanced explainer among the

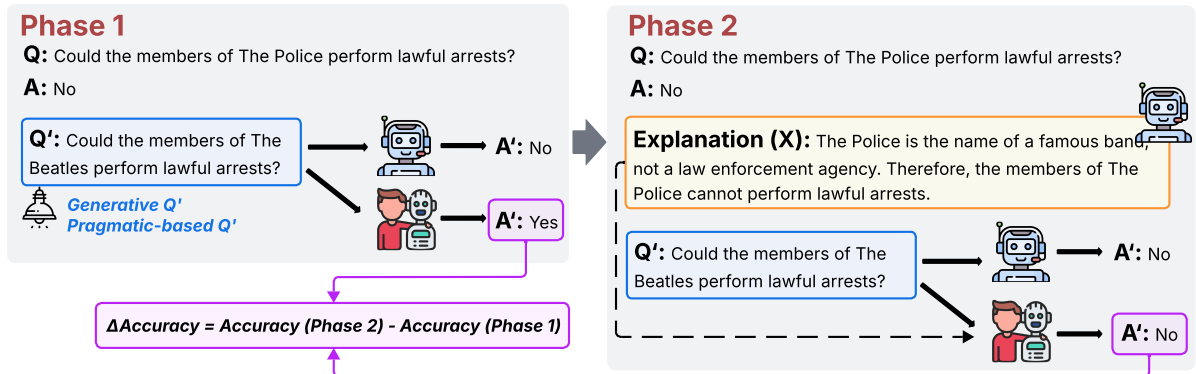


Figure 2: Counterfactual simulation pipeline for evaluating how human and LLM judges predict model behavior from explanations and counterfactual question sets.

three models considered. We select 50 items, each paired with 7 counterfactual questions using 7 different perturbation methods. In each trial, participants see a scenario involving a “robot” (a LLM): a yes/no *Starter Question*, the *Robot’s Answer*, and a yes/no *Follow-up Question* in Phase 1; in Phase 2, they additionally see a textual *Explanation* of the robot’s reasoning. Participants are explicitly instructed to predict the robot’s answer rather than state their own belief about the correct answer. For every trial, they (i) rate how confident they are about the robot’s follow-up answer on a 1–5 scale, (ii) make a guess about the robot’s answer, and (iii) provide a short free-text justification of how they arrived at this guess. These annotations allow us to analyze both the accuracy and perceived difficulty of simulating the model’s behavior, as well as the role of explanations. Detailed annotation guidelines are provided in the Appendix E.

5 LLM-as-a-Judge Results

5.1 LLM Judge Performance: Analysis of Accuracy Improvements

LLM-as-judge **acc_improvement** across different counterfactual question sets is visualized with radar charts in Figure 3. Overall, three trends emerge:

(2) **Comparison across counterfactual generation methods.** For the weaker judge, *Lexical Substitution* and *Presupposition Flip* yield the most competitive results: *Lexical Substitution* achieves the largest **acc_improvement** in four out of the six explanation settings, and *Presupposition Flip* also consistently performs well. Some LLM-generated counterfactual sets can be competitive in specific cases (e.g., for post-hoc explanations of GPT-4), but their gains are not stable across models and

explanation types. For the stronger judge, the trends change. *Scalar Adjustment* becomes one of the best-performing and most stable question type, and *Presupposition Flip* remains strong except for the CoT explanations by Llama-3.3. By contrast, the effectiveness of *Contextualization* drops substantially compared to the weaker judge: for the stronger judge, the Phase 1 accuracy on contextualized questions already reaches 0.627 (the highest among all counterfactual sets and much higher than the weaker judge’s 0.405), leaving very limited room for further gains in Phase 2. Overall, while no single counterfactual recipe dominates in all settings, pragmatics-based counterfactual questions tend to be more stable than purely generative ones, particularly *Lexical Substitution*, *Presupposition Flip*, and *Scalar Adjustment*.

(3) **CoT explanations vs. post-hoc explanations.** We do not observe a clear advantage of CoT explanations over post-hoc rationales. On average, post-hoc explanations perform slightly better: for the weaker judge, mean **acc_improvement** is 0.302 for CoT vs. 0.315 for post-hoc, and for the stronger judge it is 0.149 vs. 0.156, respectively. This suggests that, in our setup, targeted rationales are already sufficient for enabling LLM judges to simulate the model’s behavior, and richer step-by-step reasoning does not translate into larger gains.

(4) **Effect of judge strength.** Finally, comparing the two rows of radar plots shows that overall **acc_improvement** is substantially larger for the weaker evaluator than for the stronger one: the areas of the polygons are visibly bigger in the weaker-judge row. This indicates that the same combination of explanations and counterfactual questions yields much larger gains when the judge is weaker, whereas a stronger judge already predicts the base

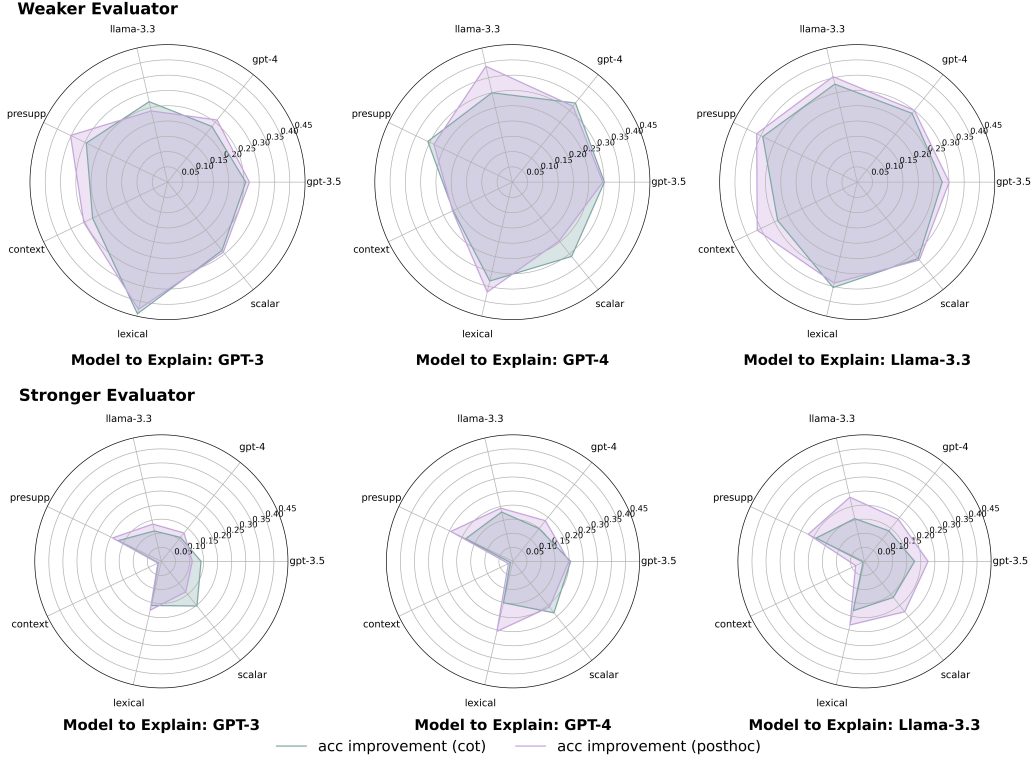


Figure 3: Radar charts summarizing LLM-as-a-judge **acc_improvement** on different counterfactual question sets. The upper and lower panels show results for a weaker and a stronger evaluator, respectively, each using combinations of chain-of-thought (cot) and post-hoc explanations for three underlying models to be explained.

379 model reasonably well in Phase 1 and therefore
 380 has limited headroom to benefit from additional
 381 information. In other words, the marginal utility of
 382 explanations and counterfactuals decreases as the
 383 judge becomes stronger.

384 Complete per-condition tables reporting **guess**
 385 **rate**, **selective accuracy**, and **overall accuracy** by
 386 phase and condition are provided in Appendix F.

387 5.2 In-context Bias and Flip-back Behavior

388 A common failure mode of LLM-judges is *anchor-*
 389 *ing*: when asked to predict the model’s answer to
 390 a follow-up question, the judge may simply reuse
 391 the model’s original answer to the starter question,
 392 regardless of whether that original answer is cor-
 393 rect. Here we quantify this anchoring effect and
 394 test whether providing self-explanations mitigates
 395 it, especially in cases where the model’s original
 396 answer is wrong (where *flip-back* is desirable).

397 For each instance i , Q_i and Q'_i denote the starter
 398 and follow-up questions. The model’s answer to Q_i
 399 is $a_i \in \{\text{yes}, \text{no}\}$. The judge predicts the model’s
 400 answer to the follow-up question as $\hat{y}_i \in \{\text{yes}, \text{no}\}$,
 401 and we denote the phase by $p_i \in \{1, 2\}$.

Anchoring to the original answer (overall). We
 402 first measure how often the judge repeats the
 403 model’s original answer when predicting the an-
 404 swer to Q'_i . We define an indicator that the judge’s
 405 prediction matches the original answer: $B_i =$
 406 $\mathbb{I}[\hat{y}_i = a_i]$. $B_i = 1$ if the judge’s prediction for
 407 Q'_i exactly equals the model’s answer to Q_i .
 408

409 For each phase $p \in \{1, 2\}$, we define the in-
 410 context bias toward the original answer as:

$$\widehat{\text{bias}}^{\text{orig}}(p) = \frac{1}{N_p} \sum_{i: p_i=p} B_i, \quad 411$$

412 where N_p is the number of instances in phase p .

Anchoring when the model is wrong (flip-back regime). To isolate cases where anchoring is clearly undesirable, we focus on instances where the model is *incorrect* on the starter question. We define an indicator of whether the model is correct on the starter question: $C_i = \mathbb{I}[a_i = g_i]$, so that $C_i = 0$ marks instances where the model’s starter answer is incorrect. Within this subset ($C_i = 0$), we reuse the indicator B_i and define:

$$\widehat{\text{bias}}^{\text{orig}}_{\text{wrong}}(p) = \frac{1}{N_p^{\text{wrong}}} \sum_{i: p_i=p, C_i=0} B_i, \quad 422$$

expl_type	cfq_type	$\Delta\text{bias_orig}$	$\Delta\text{bias_orig_wrong}$
<i>Weaker Evaluator</i>			
cot	gpt-3.5	-0.013	-0.209
cot	gpt-4	0.009	-0.224
cot	Llama-3.3	-0.043	-0.229
cot	presupp_flip	-0.195	-0.017
cot	contextual	-0.003	-0.047
cot	lexical	-0.057	-0.073
cot	scalar	-0.074	-0.076
posthoc	gpt-3.5	-0.029	-0.241
posthoc	gpt-4	-0.033	-0.304
posthoc	Llama-3.3	-0.081	-0.284
posthoc	presupp_flip	-0.237	0.013
posthoc	contextual	0.008	-0.026
posthoc	lexical	-0.054	-0.095
posthoc	scalar	-0.027	-0.042
<i>Stronger Evaluator</i>			
cot	gpt-3.5	-0.090	-0.241
cot	gpt-4	-0.057	-0.201
cot	Llama-3.3	-0.074	-0.225
cot	presupp_flip	-0.073	-0.014
cot	contextual	-0.017	-0.084
cot	lexical	-0.061	-0.073
cot	scalar	-0.049	-0.010
posthoc	gpt-3.5	-0.091	-0.254
posthoc	gpt-4	-0.081	-0.295
posthoc	Llama-3.3	-0.096	-0.290
posthoc	presupp_flip	-0.081	-0.022
posthoc	contextual	-0.028	-0.104
posthoc	lexical	-0.095	-0.149
posthoc	scalar	-0.052	-0.079

Table 3: Change in in-context bias from Phase 1 to Phase 2 for explanations of GPT-4 across counterfactual question sets. Negative values indicate reduced bias. Positive numbers are marked in blue.

where N_p^{wrong} is the number of instances in phase p for which the model’s answer to the starter question is incorrect ($C_i = 0$).

We compare $\widehat{\text{bias}}^{\text{orig}}(p)$ and $\widehat{\text{bias}}_{\text{wrong}}^{\text{orig}}(p)$ across phases: $\Delta\text{bias_orig} = \widehat{\text{bias}}^{\text{orig}}(2) - \widehat{\text{bias}}^{\text{orig}}(1)$ and $\Delta\text{bias_orig_wrong} = \widehat{\text{bias}}_{\text{wrong}}^{\text{orig}}(2) - \widehat{\text{bias}}_{\text{wrong}}^{\text{orig}}(1)$. Negative values indicate reduced anchoring to the original answer in Phase 2, while positive values indicate increased anchoring. Table 3 reports these changes for GPT-4 explanations across counterfactual question sets, both overall and restricted to instances where the model’s starter answer is incorrect. Full results are provided in Appendix H.

Overall, in-context bias tends to decrease from Phase 1 to Phase 2. This pattern suggests that LLM verbalized self-explanations generally reduce judges’ tendency to stick with the original model prediction. Moreover, the reductions are particularly pronounced when conditioning on instances where the original answer is wrong, indicating that

explanations are especially helpful for correcting initially incorrect model outputs.

6 Human User-study Results

6.1 Confidence and Accuracy Gains

Under the setup in §4.4, the mean confidence (1–5) from 28 participants increases from 3.21 without explanations to 4.38 with explanations, suggesting that explanations improve confidence in predicting the LLM’s (GPT-4) behavior.

cfq_type	P1_acc	P2_acc	Δacc
gpt-3.5	0.560	0.540	-0.020
gpt-4	0.460	0.470	0.010
Llama-3.3	0.580	0.580	0.000
presupp_flip	0.360	0.400	0.040
contextual	0.760	0.820	0.060
lexical	0.760	0.800	0.040
scalar	0.820	0.840	0.020
Overall	0.614	0.634	0.020

Table 4: Human user study evaluation results. Accuracy by counterfactual question type for Phase 1 and 2. For each counterfactual type, we evaluate 100 questions.

In Table 4, we report human prediction accuracy against the model’s actual follow-up answers in Phase 1 and Phase 2, along with the change in accuracy (Δacc). Overall accuracy increases slightly from 0.614 to 0.634. The gain is type-dependent: *Contextualization* shows the largest improvement (+0.060), while *Lexical Substitution* and *Presupposition Flip* improve by +0.040. In contrast, *Scalar Adjustment* changes only marginally (+0.020), and model-generated types show little or no benefit, with a small drop for GPT-3.5. Overall, explanations provide modest and uneven improvements in human ability to anticipate the model’s behavior.

By comparing human users to LLM-as-a-judge results, we estimate how well LLMs approximate human judgments. On the subset that humans annotated, the LLM judge matches human yes/no predictions at an overall rate of 0.601 in Phase 1 and 0.591 in Phase 2. We provide a more detailed breakdown by counterfactual type in the Appendix I.

6.2 Case Study: Analysis of Human Rationales for Predicting Model Behavior

In addition to accuracy and agreement metrics, we collected participants’ free-text rationales describing how they predicted the model’s follow-up answer. We use these rationales to qualitatively exam-

ine what cues humans attend to when simulating model behavior, and how those cues shift from Phase 1 to Phase 2. In this section, we present a small set of representative case studies selected to reflect the main quantitative patterns.

Case 1: Explanations can help when they state a decision criterion that the model also applies.

When the follow-up question directly refers to a concept mentioned in the model’s explanation, participants often move from simple Phase 1 heuristics (e.g., reusing the starter answer) to *criterion matching* in Phase 2. In these cases, their rationales typically quote or paraphrase a key phrase from the explanation and apply it to the follow-up question.

Example: *In the explanation, we know that there are programs like the National School Lunch Program that provide free or reduced-price lunches to eligible students. So few students are guaranteed lunch at school in the US.*

Case 2: Limited gains under ceiling effects.

When the follow-up question is a near-paraphrase of the starter question or introduces only a minor adjustment, many participants already succeed in Phase 1 by applying a simple consistency heuristic (predicting that the model will answer the same way). In these cases, rationales often acknowledge that the explanation is consistent with their initial guess but does not add further information. Once Phase 1 performance is near ceiling, explanations have little room to improve correctness.

Example: *The explanation further support the answer with further information like "The Smithsonian’s National Zoo in Washington DC is a place where various animals are kept, but it does not currently house harbor seals."*

Case 3: Explanations may not resolve condition flips.

In these cases, participants’ rationales often express uncertainty about how the model will respond after the follow-up introduces a critical change in conditions. Even in Phase 2, the model’s explanation may not directly address the specific condition that is altered in the follow-up, leaving participants without a stable rule to extrapolate the model’s behavior. As a result, these items can remain difficult across phases. Qualitatively, the rationales suggest that participants focus on the condition shift in the follow-up, while the explanation highlights other aspects, which makes it less useful for predicting the model’s answer.

Example: *The explanation states that "if all con-*

ditions align", but in the follow-up question, the condition changes. Not sure whether the robot still thinks that a dealer could buy a Boeing 737-800.

Case 4: Explanations can distract when they are persuasive but not diagnostic.

We observe cases where a participant’s Phase 1 prediction is correct under a simple heuristic, but Phase 2 introduces over-reliance on a plausible-sounding detail in the explanation that does not actually determine the model’s follow-up answer. In these instances, the explanation functions as a *persuasive narrative* rather than a diagnostic cue, and participants may flip from a correct to an incorrect prediction.

Example: *The explanation already compares mail carriers to more dangerous jobs, so police officers will be judged dangerous.*

Summary. Across these case studies, human rationales suggest two recurring themes. First, participants are successful when they can align the follow-up question to an explicit criterion stated in the explanation; this is precisely when Phase 2 yields the largest gains. Second, when follow-up questions are either trivial (ceiling effects) or underspecified (insufficient cues), explanations offer limited additional benefit and may even distract.

7 Conclusion

In this work, we assess LLM self-explanations through counterfactual simulatability. Explanations improve users’ ability to anticipate model behavior for both human participants and LLM judges, but the measured effect depends on the counterfactual construction and the evaluator. We further observe that explanations mitigate a tendency to echo the model’s original answer—most notably when the initial answer is incorrect—suggesting a corrective “flip-back” mechanism. In the human study, explanations also increase confidence and yield measurable accuracy gains.

Overall, our findings show that, despite concerns about faithfulness of verbalized self-explanations, they can help users better understand model behavior in counterfactual settings. Our results highlight that the counterfactual questions themselves are a key part of the evaluation: different perturbation strategies can change both task difficulty and the benefit of explanations, so counterfactual design should be treated as a component rather than a neutral choice. Annotations, explanations, and code will be released publicly upon publication.

577 Limitations

578 This study has potential limitations. Our study
579 measures model understanding using counterfac-
580 tual simulatability on binary questions. While
581 this design offers a controlled way to measure
582 whether explanations help users predict model be-
583 havior, it may not fully reflect real-world interac-
584 tions where users ask open-ended questions and
585 engage in multi-turn dialogue. Second, although
586 we introduce a pragmatics-grounded taxonomy of
587 counterfactual perturbations, the set of transforma-
588 tions is incomplete and may not cover pragmatic
589 phenomena that arise in natural conversations (e.g.,
590 discourse context, speaker intent, social norms).

591 Thirdly, parts of our pipeline rely on LLM-
592 generated artifacts (e.g., LLM-generated counter-
593 factuals and explanations, and GPT-4-based pre-
594 supposition flips), such as counterfactual questions
595 or explanations. Results may depend on the chosen
596 prompts and models. We also test a limited number
597 of models and explanation formats, so generaliza-
598 tion is not guaranteed.

599 Finally, the human user study uses a relatively
600 small sample (participants and items) and focuses
601 on a subset of conditions, which limits statisti-
602 cal power for small effects and may not gener-
603 alize across user populations or expertise levels.
604 Importantly, our work does not claim that self-
605 explanations are faithful accounts of the model’s in-
606 ternal decision process: explanations can improve
607 simulatability while still being post-hoc or persua-
608 sive rather than diagnostic.

609 Acknowledgments

610 **Ethical considerations.** We do not foresee any
611 ethical concerns associated with this work. All
612 analyses were conducted using publicly available
613 datasets and models. No private or sensitive in-
614 formation was used. Additionally, we release our
615 code, prompts, and documentation to support trans-
616 parency and reproducibility.

617 **Use of AI Assistants.** The authors acknowl-
618 edge the use of ChatGPT for correcting grammat-
619 ical errors, enhancing the coherence of the final
620 manuscripts, and providing coding assistance.

621 References

622 Christina Alexandris. 2024. [Genai and socially respon-](#)
623 [sible ai in natural language processing applications:](#)
624 [A linguistic perspective.](#) In *AAAI Spring Symposia*.

- David Alvarez-Melis and Tommi S. Jaakkola. 2018. [Towards robust interpretability with self-explaining neural networks.](#) In *Proceedings of the 32nd International Conference on Neural Information Processing Systems, NIPS’18*, page 7786–7795, Red Hook, NY, USA. Curran Associates Inc. 625
626
627
628
629
630
- Pepa Atanasova, Oana-Maria Camburu, Christina Li-
oma, Thomas Lukasiewicz, Jakob Grue Simonsen,
and Isabelle Augenstein. 2023. [Faithfulness tests](#)
[for natural language explanations.](#) In *Proceedings*
of the 61st Annual Meeting of the Association for
Computational Linguistics (Volume 2: Short Papers),
pages 283–294, Toronto, Canada. Association for
Computational Linguistics. 631
632
633
634
635
636
637
638
- Yushi Bai, Jiahao Ying, Yixin Cao, Xin Lv, Yuze
He, Xiaozhi Wang, Jifan Yu, Kaisheng Zeng, Yijia
Xiao, Haozhe Lyu, Jiayin Zhang, Juanzi Li, and Lei
Hou. 2023. [Benchmarking foundation models with](#)
[language-model-as-an-examiner.](#) In *Proceedings of*
the 37th International Conference on Neural Infor-
mation Processing Systems, NIPS ’23, Red Hook,
NY, USA. Curran Associates Inc. 639
640
641
642
643
644
645
646
- Amrita Bhattacharjee, Raha Moraffah, Joshua Garland,
and Huan Liu. 2023. [Towards llm-guided causal](#)
[explainability for black-box text classifiers.](#) 647
648
649
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
Neelakantan, Pranav Shyam, Girish Sastry, Amanda
Askell, Sandhini Agarwal, Ariel Herbert-Voss,
Gretchen Krueger, Tom Henighan, Rewon Child,
Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,
Clemens Winter, and 12 others. 2020. [Language](#)
[models are few-shot learners.](#) In *Proceedings of the*
34th International Conference on Neural Information
Processing Systems, NIPS ’20, Red Hook, NY, USA.
Curran Associates Inc. 650
651
652
653
654
655
656
657
658
659
660
- Oana-Maria Camburu, Tim Rocktäschel, Thomas
Lukasiewicz, and Phil Blunsom. 2018. [e-snli: Natu-](#)
[ral language inference with natural language explana-](#)
[tions.](#) In *Advances in Neural Information Processing*
Systems, volume 31. Curran Associates, Inc. 661
662
663
664
665
- Robyn Carston. 1998. [Informativeness, relevance and](#)
[scalar implicature.](#) *Relevance Theory. Applications*
and Implications. 666
667
668
- Arjun Chandrasekaran, Viraj Prabhu, Deshraj Yadav,
Prithvijit Chattopadhyay, and Devi Parikh. 2018. [Do](#)
[explanations make VQA models more predictable to](#)
[a human?](#) In *Proceedings of the 2018 Conference on*
Empirical Methods in Natural Language Processing,
pages 1036–1042, Brussels, Belgium. Association
for Computational Linguistics. 669
670
671
672
673
674
675
- Yanda Chen, Ruiqi Zhong, Narutatsu Ri, Chen Zhao,
He He, Jacob Steinhardt, Zhou Yu, and Kathleen
McKeown. 2024. [Do models explain themselves?](#)
[counterfactual simulatability of natural language ex-](#)
[planations.](#) In *Proceedings of the 41st Interna-*
tional Conference on Machine Learning, ICML’24.
JMLR.org. 676
677
678
679
680
681
682

683	Finale Doshi-Velez and Been Kim. 2017. Towards a rigorous science of interpretable machine learning . <i>arXiv: Machine Learning</i> .	
684		
685		
686	Matt Gardner, Yoav Artzi, Victoria Basmov, Jonathan Berant, Ben Bogin, Sihao Chen, Pradeep Dasigi, Dheeru Dua, Yanai Elazar, Ananth Gottumukkala, Nitish Gupta, Hannaneh Hajishirzi, Gabriel Ilharco, Daniel Khashabi, Kevin Lin, Jiangming Liu, Nelson F. Liu, Phoebe Mulcaire, Qiang Ning, and 7 others. 2020. Evaluating models’ local decision boundaries via contrast sets . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 1307–1323, Online. Association for Computational Linguistics.	
687		
688		
689		
690		
691		
692		
693		
694		
695		
696		
697	Yingqiang Ge, Shuchang Liu, Zelong Li, Shuyuan Xu, Shijie Geng, Yunqi Li, Juntao Tan, Fei Sun, and Yongfeng Zhang. 2021. Counterfactual evaluation for explainable AI . <i>CoRR</i> , abs/2109.01962.	
698		
699		
700		
701	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 . <i>Transactions of the Association for Computational Linguistics</i> , 9:346–361.	
702		
703		
704		
705		
706		
707	Leilani H. Gilpin, David Bau, Ben Z. Yuan, Ayesha Bajwa, Michael A. Specter, and Lalana Kagal. 2018. Explaining explanations: An overview of interpretability of machine learning . <i>2018 IEEE 5th International Conference on Data Science and Advanced Analytics (DSAA)</i> , pages 80–89.	
708		
709		
710		
711		
712		
713	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 542 others. 2024. The llama 3 herd of models . <i>CoRR</i> , abs/2407.21783.	
714		
715		
716		
717		
718		
719		
720		
721	H. P. Grice. 1975. Logic and conversation . In Peter Cole and Jerry L. Morgan, editors, <i>Syntax and Semantics: Vol. 3: Speech Acts</i> , pages 41–58. Academic Press, New York.	
722		
723		
724		
725	Peter Hase and Mohit Bansal. 2020. Evaluating explainable AI: Which algorithmic explanations help users predict model behavior? In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5540–5552, Online. Association for Computational Linguistics.	
726		
727		
728		
729		
730		
731	Robert R. Hoffman, Shane T. Mueller, Gary Klein, and Jordan Litman. 2018. Metrics for explainable AI: challenges and prospects . <i>CoRR</i> , abs/1812.04608.	
732		
733		
734	Alon Jacovi and Yoav Goldberg. 2020. Towards faithfully interpretable nlp systems: How should we define and evaluate faithfulness? In <i>Annual Meeting of the Association for Computational Linguistics</i> .	
735		
736		
737		
	Baptiste Jacquet, Olivier Masson, Frank Jamet, and Jean Baratgin. 2019. On the Lack of Pragmatic Processing in Artificial Conversational Agents , pages 394–399.	738 739 740
	Katarzyna Kaczmarek-Majer, Gabriella Casalino, Giovanna Castellano, Monika Dominiak, Olgierd Hryniewicz, Olga Kaminska, Gennaro Vessio, and Natalia Diaz-Rodríguez. 2022. Plenary: Explaining black-box models in natural language through fuzzy linguistic summaries . <i>Information Sciences</i> , 614.	741 742 743 744 745 746
	Daniel Kahneman and Amos Tversky. 1982. The simulation heuristic .	747 748
	Zachary Kenton, Noah Y. Siegel, János Kramár, Jonah Brown-Cohen, Samuel Albanie, Jannis Bulian, Rishabh Agarwal, David Lindner, Yunhao Tang, Noah D. Goodman, and Rohin Shah. 2024. On scalable oversight with weak llms judging strong llms . In <i>Proceedings of the 38th International Conference on Neural Information Processing Systems, NIPS ’24</i> , Red Hook, NY, USA. Curran Associates Inc.	749 750 751 752 753 754 755 756
	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners . In <i>Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS ’22</i> , Red Hook, NY, USA. Curran Associates Inc.	757 758 759 760 761 762
	Lea Krause and Piek T.J.M. Vossen. 2024. The Gricean maxims in NLP - a survey . In <i>Proceedings of the 17th International Natural Language Generation Conference</i> , pages 470–485, Tokyo, Japan. Association for Computational Linguistics.	763 764 765 766 767
	Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad Beigi, Chengshuai Zhao, Zhen Tan, Amrita Bhat-tacharjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu, Kai Shu, Lu Cheng, and Huan Liu. 2025. From generation to judgment: Opportunities and challenges of LLM-as-a-judge . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 2757–2791, Suzhou, China. Association for Computational Linguistics.	768 769 770 771 772 773 774 775 776
	Marvin Limpijankit, Yanda Chen, Melanie Subbiah, Nicholas Deas, and Kathleen McKeown. 2025. Counterfactual simulatability of LLM explanations for generation tasks . In <i>Proceedings of the 18th International Natural Language Generation Conference</i> , pages 659–683, Hanoi, Vietnam. Association for Computational Linguistics.	777 778 779 780 781 782 783
	Bolei Ma, Yuting Li, Wei Zhou, Ziwei Gong, Yang Janet Liu, Katja Jasinskaja, Annemarie Friedrich, Julia Hirschberg, Frauke Kreuter, and Barbara Plank. 2025. Pragmatics in the era of large language models: A survey on datasets, evaluation, opportunities and challenges . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8679–8696, Vienna, Austria. Association for Computational Linguistics.	784 785 786 787 788 789 790 791 792
	Nishtha Madaan, Inkit Padhi, Naveen Panwar, and Dip-tikalyan Saha. 2020. Generate your counterfactuals :	793 794

795	Towards controlled counterfactual generation for text.	Judea Pearl. 2021. <i>Causal and Counterfactual Inference</i> , pages 427–438.	851
796	In <i>AAAI Conference on Artificial Intelligence</i> .		852
797	Andreas Madsen, Sarath Chandar, and Siva Reddy. 2024.	Alec Radford, Jeff Wu, Rewon Child, David Luan,	853
798	Are self-explanations from large language models	Dario Amodei, and Ilya Sutskever. 2019. <i>Language</i>	854
799	faithful? In <i>Findings of the Association for Com-</i>	<i>models are unsupervised multitask learners</i> .	855
800	<i>putational Linguistics: ACL 2024</i> , pages 295–337,		
801	Bangkok, Thailand. Association for Computational	Marco Ribeiro, Sameer Singh, and Carlos Guestrin.	856
802	Linguistics.	2016. “why should I trust you?”: Explaining the pre-	857
		dictions of any classifier. In <i>Proceedings of the 2016</i>	858
803	George A. Miller. 1992. <i>WordNet: A lexical database</i>	<i>Conference of the North American Chapter of the</i>	859
804	for English . In <i>Speech and Natural Language: Pro-</i>	<i>Association for Computational Linguistics: Demon-</i>	860
805	<i>ceedings of a Workshop Held at Harriman, New York,</i>	<i>strations</i> , pages 97–101, San Diego, California. As-	861
806	<i>February 23-26, 1992</i> .	sociation for Computational Linguistics.	862
807	Tim Miller. 2019. <i>Explanation in artificial intelligence:</i>	Alexis Ross, Ana Marasović, and Matthew Peters. 2021.	863
808	Insights from the social sciences . <i>Artificial Intelli-</i>	Explaining NLP models via minimal contrastive edit-	864
809	<i>gence</i> , 267:1–38.	ing (MiCE) . In <i>Findings of the Association for Com-</i>	865
		<i>putational Linguistics: ACL-IJCNLP 2021</i> , pages	866
810	Van Bach Nguyen, Paul Youssef, Christin Seifert, and	3840–3852, Online. Association for Computational	867
811	Jörg Schlotterer. 2024. <i>LLMs for generating and</i>	Linguistics.	868
812	evaluating counterfactuals: A comprehensive study .		
813	In <i>Findings of the Association for Computational</i>	Uli Sauerland. 2004. <i>Scalar Implicatures in Complex</i>	869
814	<i>Linguistics: EMNLP 2024</i> , pages 14809–14824, Mi-	Sentences . <i>Linguistics and Philosophy</i> , 27(3):367–	870
815	ami, Florida, USA. Association for Computational	391.	871
816	Linguistics.		
817	Maxwell I. Nye, Anders Johan Andreassen, Guy Gur-	Loris Schoenegger, Yuxi Xia, and Benjamin Roth. 2024.	872
818	Ari, Henryk Michalewski, Jacob Austin, David	An evaluation of explanation methods for black-	873
819	Bieber, David Dohan, Aitor Lewkowycz, Maarten	box detectors of machine-generated text . <i>CoRR</i> ,	874
820	Bosma, David Luan, Charles Sutton, and Augustus	abs/2408.14252.	875
821	Odena. 2021. <i>Show your work: Scratchpads for inter-</i>		
822	mediate computation with language models . <i>CoRR</i> ,	Marcos Treviso, Alexis Ross, Nuno M. Guerreiro, and	876
823	abs/2112.00114.	André Martins. 2023. <i>CREST: A joint framework</i>	877
		for rationalization and counterfactual text genera-	878
824	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal,	tion . In <i>Proceedings of the 61st Annual Meeting of the</i>	879
825	Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-	<i>Association for Computational Linguistics (Volume 1:</i>	880
826	man, Diogo Almeida, Janko Altschmidt, Sam Alt-	<i>Long Papers)</i> , pages 15109–15126, Toronto, Canada.	881
827	man, Shyamal Anadkat, Red Avila, Igor Babuschkin,	Association for Computational Linguistics.	882
828	Suchir Balaji, Valerie Balcom, Paul Baltescu, Haim		
829	ing Bao, Mo Bavarian, Jeff Belgum, and 261 others.	Qianli Wang, Tatiana Anikina, Nils Feldhus, Simon Os-	883
830	2023. <i>Gpt-4 technical report</i> .	termann, and Sebastian Möller. 2024. <i>CoXQL: A</i>	884
		dataset for parsing explanation requests in conversa-	885
831	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Car-	tional XAI systems . In <i>Findings of the Association</i>	886
832	roll L. Wainwright, Pamela Mishkin, Chong Zhang,	<i>for Computational Linguistics: EMNLP 2024</i> , pages	887
833	Sandhini Agarwal, Katarina Slama, Alex Ray, John	1410–1422, Miami, Florida, USA. Association for	888
834	Schulman, Jacob Hilton, Fraser Kelton, Luke Miller,	Computational Linguistics.	889
835	Maddie Simens, Amanda Askell, Peter Welinder,		
836	Paul Christiano, Jan Leike, and Ryan Lowe. 2022.	Qianli Wang, Nils Feldhus, Simon Ostermann, Luis Fe-	890
837	Training language models to follow instructions with	lippe Villa-Arenas, Sebastian Möller, and Vera	891
838	human feedback . In <i>Proceedings of the 36th Interna-</i>	Schmitt. 2025. <i>FitCF: A framework for automatic</i>	892
839	<i>tional Conference on Neural Information Processing</i>	feature importance-guided counterfactual example	893
840	<i>Systems</i> , NIPS ’22, Red Hook, NY, USA. Curran	generation . In <i>Findings of the Association for Com-</i>	894
841	Associates Inc.	<i>putational Linguistics: ACL 2025</i> , pages 1176–1191,	895
		Vienna, Austria. Association for Computational Lin-	896
		guistics.	897
842	Anna Papafragou and Julien Musolino. 2003. <i>Scalar im-</i>	Tongshuang Wu, Marco Tulio Ribeiro, Jeffrey Heer, and	898
843	plicatures: experiments at the semantics–pragmatics	Daniel Weld. 2021. <i>Polyjuice: Generating counter-</i>	899
844	interface . <i>Cognition</i> , 86(3):253–282.	factuals for explaining, evaluating, and improving	900
		models . In <i>Proceedings of the 59th Annual Meet-</i>	901
845	Dong Huk Park, Lisa Anne Hendricks, Zeynep Akata,	<i>ing of the Association for Computational Linguistics</i>	902
846	Anna Rohrbach, Bernt Schiele, Trevor Darrell, and	<i>and the 11th International Joint Conference on Natu-</i>	903
847	Marcus Rohrbach. 2018. <i>Multimodal explanations:</i>	<i>ral Language Processing (Volume 1: Long Papers)</i> ,	904
848	Justifying decisions and pointing to the evidence . In	pages 6707–6723, Online. Association for Computa-	905
849	<i>2018 IEEE/CVF Conference on Computer Vision and</i>	tional Linguistics.	906
850	<i>Pattern Recognition</i> , pages 8779–8788.		

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. [Judging llm-as-a-judge with mt-bench and chatbot arena](#). In *Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS '23*, Red Hook, NY, USA. Curran Associates Inc.

Longfei Zuo, Barbara Plank, and Siyao Peng. 2025. [EVADE: llm-based explanation generation and validation for error detection in NLI](#). *CoRR*, abs/2511.08949.

A LLM-prompts for automatic counterfactual questions generation

This appendix shows the prompt used to generate follow-up counterfactual yes/no questions on the StrategyQA dataset. We use exactly the same natural-language prompt for all models and both backends (OpenAI Chat Completions and local Hugging Face models). For OpenAI models, the prompt is passed as system/user messages, and for Hugging Face models we serialize the same messages via the tokenizer’s chat template using the code in the main text.

System Instructions

Human: In the questions below, you will be asked to read a starter yes or no question and a robot’s answer to the starter question. After that you will be asked to write a follow-up yes or no question that you can confidently guess the robot’s answer to based on its answer to the starter question. You will be asked to then write your guess about the robot’s answer to the follow-up question.

Assistant: here is my response. okay.

Template

Human: Starter Question: {orig_qn}
 Robot’s Answer to the Starter Question: {orig_yn}
 Follow-up Question:

Format Requirement

When you answer, output exactly two lines:
 Follow-up Question: <a yes/no question>
 Guessed Answer: yes or no
 Do not include anything else.

B LLM-prompts for presupposition flip counterfactual questions generation

System Instructions

You are a precise editor that rewrites yes/no questions into counterfactual questions with if-clauses while flipping the original answer. Follow the user’s format exactly, be concise, and do not add extra commentary or bullets.

User Instructions

Prompt/ instructions for LLM:

Step 1: Read the original yes/no question and the answer.
 Step 2: Flip the original answer.
 Step 3: Work out how to make the answer true.
 Step 4: Write a new question that asks if the answer would be true if the action you worked out in Step 3 was performed. Change the original questions to new questions of unreal conditions with counterfactual presuppositions, using if clauses. Do not change the meaning of questions in new questions.
 Step 5: Write the new answer to the new questions.

Answer each initial question with the following format:
 Original question: <original question>
 Original answer: <original answer>
 Step1: The new answer should be <yes/no>, so
 Step2: How to make . . . :
 New question: <new question>

Here are some examples:

Original question: Would the top of Mount Fuji stick out of the Sea of Japan?
 Original answer: Yes
 New question: If global sea level were 4,000 meters higher, would the top of Mount Fuji stick out of the Sea of Japan?
 New answer: No

Original question: Would an uninsured person be more likely than an insured person to decline a CT scan?
 Original answer: Yes
 New question: If CT scans were free at the point of care for everyone, would an uninsured person be more likely than an insured person to decline a CT scan?
 New answer: No

Original question: Did Doctor Strange creators also make Batman?
 Original answer: No
 New question: If Batman had been created by Stan Lee and Steve Ditko, would the creators of Doctor Strange also have created Batman?
 New answer: Yes

C Scalar scales

Table 5 lists the manually specified scalar scales used in our scalar adjustment operations. The scales are defined on the basis of scalar inventories discussed in previous work on scalar expressions and implicatures (Carston, 1998; Papafragou and Musolino, 2003; Sauerland, 2004).

Scale	Ordered items (stronger → weaker)
Quantifiers	<i>all, almost all, most, a majority of, many, several, some, a few, few, hardly any</i>
Universal quantifiers	<i>every, each</i>
Existential quantifiers	<i>any, at least one, at least some</i>
Frequencies	<i>always, almost always, usually, often, frequently, sometimes, occasionally, rarely, seldom, hardly ever, never</i>
Epistemic modality	<i>certainly, definitely, undoubtedly, almost certainly, very likely, likely, probably, plausibly, possibly, maybe, perhaps</i>
Deontic modality	<i>must, have to, be required to, should, ought to, be supposed to, be recommended to, may, can</i>
Necessity / possibility	<i>necessary, required, essential, recommended, advisable, optional, possible</i>
Attitude / success	<i>succeed in, manage to, be able to, try to, attempt to, want to</i>
Evidential / reporting	<i>prove that, demonstrate that, show that, find that, report that, suggest that, speculate that</i>
Claim strength / hedging	<i>it is certain that, it is clear that, it is likely that, it is probable that, it is plausible that, it is possible that</i>
Degree adverbs	<i>completely, totally, entirely, highly, very, quite, fairly, rather, somewhat, slightly</i>
Evaluative adjectives	<i>perfect, excellent, great, good, acceptable, adequate, poor</i>
Temporal immediacy	<i>immediately, right away, soon, eventually, someday</i>
Inclusion	<i>including all, including most, including many, including some</i>

Table 5: Manually specified scalar scales used for scalar adjustment. Within each row, items are ordered from stronger to weaker.

D LLM-prompts for explanations generation

This appendix shows the prompt used for generating chain-of-thought and post-hoc explanations on the StrategyQA dataset. We use exactly the same natural-language prompt for all three models and both backends. For OpenAI models, the prompt is passed as system/user chat messages, while for Hugging Face models we serialize the same messages using the tokenizer’s chat template.

Prompt: Explanation Generator

Generate explanations (CoT / Post-hoc).

System instruction for CoT (COT_SYSTEM):

Human: In the questions below, you will be asked to first think step by step and generate a trace of reasoning, and then end with your final answer exactly with “So the answer is ...” yes or no. Let’s think step by step. Strictly follow the example format below and do not say anything else.

Assistant: here is my response. okay.

System instruction for post-hoc explanations (POSTHOC_SYSTEM):

Human: In the questions below, you will be asked to first answer the question (yes or no) and then generate a justification for your answer. Your justification should end with “So the answer is ...” yes or no.

Assistant: here is my response. okay.

User message template for CoT:

Human: Q: Yes or no: {question}

A:

Assistant: here is my response.

User message template for post-hoc explanations:

Human: Q: Yes or no: {question}

A:

Assistant: here is my response.

E User study guideline

We recruited 28 volunteer participants (B.Sc., M.Sc., and PhD students) with English reading proficiency at CEFR level C1 or above. Before the

960	study, participants reported (i) their familiarity with	guess on: the <i>Starter Question</i> , the <i>Robot’s Answer</i> ,	1007
961	LLMs on a 5-point Likert scale and (ii) whether	and (in Phase 2) the <i>Explanation</i> .	1008
962	they were familiar with explainable AI (XAI). Most	In your explanations, please focus on how you	1009
963	participants frequently used LLMs in their work:	inferred the robot’s behavior, not on what you per-	1010
964	18/28 selected 5 (“very familiar / frequent use”),	sonally think is true about the world.	1011
965	4/28 selected 4 (“regular use in daily life”), and the		
966	remaining 6/28 selected 3 (“occasional use”). In	E.3 Confidence rating and forced guess	1012
967	contrast, only 2/28 reported in-depth knowledge of	For each trial, you will first answer:	1013
968	XAI; the rest indicated limited familiarity or only a		
969	basic awareness. Participation was fully voluntary,	“How confident are you that you can rea-	1014
970	and participants were generally interested in NLP	son about what the robot will answer to	1015
971	and AI research. The experiment was conducted	the follow-up question?”	1016
972	online via a purpose-built web service.		
973	In the following, we provide the full annotation	You will choose a number on a 1–5 scale (e.g., 1	1017
974	guideline used in our user study. Before starting	= Not confident at all, 5 = Extremely confident).	1018
975	the study, each participant is shown this guideline,	After choosing a confidence level, you must still	1019
976	which includes a general introduction to the task,	select one of:	1020
977	detailed instructions, and illustrative examples. Par-		
978	ticipants are encouraged to ask questions both be-	• The robot will answer: Yes	1021
979	fore and during the study. Screenshots of the user	• The robot will answer: No	1022
980	study interface are shown in Figures 4 and 5.		
981	E.1 Overview of the task	There are no “I don’t know” or “cannot tell” op-	1023
982	In this study, you will see short question–answer	tions in the multiple-choice part.	1024
983	pairs involving a “robot” (a large language model).	If you feel that you really cannot predict the	1025
984	Your job is to reason about what the robot will	robot’s answer at all:	1026
985	answer to a follow-up yes/no question. There are		
986	two phases:	• choose the lowest confidence level;	1027
987	Phase 1 : You see a <i>Starter Question</i> , the <i>Robot’s</i>	• still pick Yes or No as a forced guess;	1028
988	<i>Answer</i> , and a <i>Follow-up Question</i> .	• in the free-text explanation, clearly say that	1029
989	Phase 2 : You see the same elements plus an	you cannot guess, e.g.: “ <i>I can’t guess the</i>	1030
990	<i>Explanation</i> of the robot’s reasoning for the follow-	<i>robot’s answer; I chose ‘Yes’ randomly.</i> ”	1031
991	up question.	This helps us distinguish between informed rea-	1032
992	In every trial, you will:	soning and random guessing.	1033
993		E.4 What to write in your free-text	1034
994	• rate how confident you are that you can rea-	explanation	1035
995	son about the robot’s answer to the follow-up	After you have given your confidence rating and	1036
996	question;	your Yes/No guess, you will write a short expla-	1037
997	• guess the robot’s answer to the follow-up ques-	tion (typically 1–3 sentences). Here are some	1038
998	tion (Yes or No);	general tips:	1039
999	• provide a short free-text explanation of how	• Keep it simple and concise, but concrete.	1040
1000	you made your guess.	• Explain why you think the robot will answer	1041
1001	E.2 Very important: Guess the robot’s answer,	Yes or No.	1042
1002	not your own opinion	• Focus on how you used the robot’s previous	1043
1003	Your task is not to say what you think is objectively	answer and (in Phase 2) the explanation.	1044
1004	correct. Your task is to guess what the robot will	• You do not need to write a long essay; 1–3	1045
1005	answer to the follow-up question. Even if you	sentences are enough.	1046
1006	believe the robot is wrong or unreasonable, you		
	should still try to infer its likely answer. Base your		

1047	E.5 Phase 1 guidelines (without explanation)	about your reasoning. If it really is just a guess,	1094
1048	Typical patterns you may use (and can mention in	please still explain why you cannot reason about	1095
1049	your explanation):	the robot’s answer and mark low confidence. Here	1096
1050	Same conditions → same answer: “ <i>The follow-</i>	are some examples of less informative explana-	1097
1051	<i>up question is very similar to the starter question,</i>	tions:	1098
1052	<i>so I think the robot will give the same answer.</i> ”	Example 1: “ <i>I think the correct answer is yes.</i> ”	1099
1053	Opposite conditions → opposite answer: “ <i>The</i>	Example 2: “ <i>Because of the question.</i> ”	1100
1054	<i>follow-up question reverses the condition in the</i>	Example 3: “ <i>Just a guess.</i> ”	1101
1055	<i>starter question, so I expect the robot to give the</i>		
1056	<i>opposite answer.</i> ”	E.8 Summary	1102
1057	Inference from the Robot’s Answer: “ <i>From the</i>	• Always guess the robot’s answer, not your	1103
1058	<i>robot’s answer ‘no’ to the starter question, I infer</i>	own belief about what is true.	1104
1059	<i>that it thinks X is unlikely, so I expect it will also</i>	• Always provide a Yes/No guess, even if you	1105
1060	<i>answer ‘no’ to the follow-up.</i> ”	are not confident.	1106
1061	Unclear / low confidence: “ <i>I’m not confident.</i>	• If you really cannot predict the answer:	1107
1062	<i>The robot’s answer to the starter question does not</i>	– choose “not confident”;	1108
1063	<i>clearly tell me how it will answer the follow-up, so</i>	– still choose Yes or No; and	1109
1064	<i>this is mostly a guess.</i> ”	– write in your explanation that you cannot	1110
1065	E.6 Phase 2 guidelines (with explanation)	guess and why.	1111
1066	For your free-text explanation in Phase 2 , you	• In Phase 1, base your reasoning mainly on the	1112
1067	are encouraged to refer directly to the explanation.	<i>Starter Question</i> and the <i>Robot’s Answer</i> .	1113
1068	You can quote short parts in quotation marks, for	• In Phase 2, also use the <i>Explanation</i> , and feel	1114
1069	example: “ <i>The explanation says ‘Rede Globo is a</i>	free to quote it in quotation marks.	1115
1070	<i>Brazilian television network’, so I expect the robot</i>	• Keep explanations short but specific, focusing	1116
1071	<i>to answer ‘Yes’ when asked if the anchors speak</i>	on how you infer the robot’s answer.	1117
1072	<i>Portuguese.</i> ”		
1073	You can combine the explanation with your own	F Full results for LLM-as-a-judge	1118
1074	reasoning: “ <i>The explanation says ‘X is the official</i>	counterfactual simulation	1119
1075	<i>language of Y’, which suggests the robot believes</i>	This appendix reports the full per-condition results	1120
1076	<i>people in Y usually speak X. Therefore, I think it</i>	for our LLM-as-a-judge counterfactual simulation.	1121
1077	<i>will answer ‘Yes’ to the follow-up question.</i> ”	We report three complementary accuracy-based	1122
1078	E.7 Examples of good and less useful	metrics: (i) guess_rate , the proportion of items	1123
1079	explanations	judged as can_guess; (ii) selective_accuracy , ac-	1124
1080	Good explanations refer to the robot’s answer	curacy conditioned on guessing; and (iii) over-	1125
1081	and/or the explanation and make it clear how you	all accuracy , accuracy on all items, counting	1126
1082	predicted the robot’s behavior. Here are some ex-	cannot_guess as incorrect. Results using two	1127
1083	amples of good free-text explanations:	LLM judges, gpt-3.5-turbo and gpt-4.1, are	1128
1084	Example 1: “ <i>From the robot’s ‘no’ to the starter</i>	reported in Table 6 and Table 7.	1129
1085	<i>question and the explanation that Chinese is not</i>		
1086	<i>commonly spoken in Brazil, I think it will answer</i>	G Automatic Evaluation for	1130
1087	<i>‘yes’ when asked about Portuguese.</i> ”	Counterfactuals	1131
1088	Example 2: “ <i>The explanation states ‘X is un-</i>	In addition to simulation-based evaluation, we also	1132
1089	<i>likely’, so I expect the robot to keep the same nega-</i>	consider automated metrics that are commonly	1133
1090	<i>tive answer for the follow-up.</i> ”	used to characterize counterfactual data. These	1134
1091	Try to avoid explanations that rely only on per-	metrics quantify how a perturbation Q' related to	1135
1092	sonal opinions without referring to the robot, that	its original question Q , independently of any expla-	1136
1093	simply repeat the question, or that say nothing	nation.	1137

expl	cfq_type	P1_guess_rate	P1_acc	P1_sel_acc	P2_guess_rate	P2_acc	P2_sel_acc	Δ guess_rate	Δ acc	Δ sel_acc
cot (gpt-3.5)	gpt-3.5	0.505	0.368	0.000	0.901	0.624	0.737	0.396	0.256	0.008
cot (gpt-3.5)	gpt-4	0.623	0.479	0.769	0.940	0.712	0.757	0.317	0.233	-0.012
cot (gpt-3.5)	Llama-3.3	0.562	0.414	0.735	0.934	0.684	0.732	0.371	0.270	-0.003
cot (gpt-3.5)	presupp_flip	0.434	0.267	0.614	0.894	0.562	0.628	0.460	0.295	0.014
cot (gpt-3.5)	contextual	0.479	0.315	0.658	0.832	0.587	0.706	0.353	0.272	0.048
cot (gpt-3.5)	lexical	0.244	0.134	0.550	0.820	0.507	0.617	0.576	0.442	0.067
cot (gpt-3.5)	scalar	0.514	0.304	0.590	0.869	0.590	0.679	0.355	0.286	0.089
posthoc (gpt-3.5)	gpt-3.5	0.485	0.358	0.738	0.842	0.626	0.743	0.357	0.268	0.005
posthoc (gpt-3.5)	gpt-4	0.552	0.387	0.702	0.862	0.647	0.750	0.310	0.260	0.048
posthoc (gpt-3.5)	Llama-3.3	0.517	0.374	0.723	0.797	0.613	0.768	0.280	0.239	0.045
posthoc (gpt-3.5)	presupp_flip	0.446	0.311	0.698	0.810	0.663	0.652	0.364	0.352	-0.046
posthoc (gpt-3.5)	contextual	0.526	0.339	0.643	0.901	0.643	0.714	0.375	0.304	0.071
posthoc (gpt-3.5)	lexical	0.448	0.339	0.758	0.892	0.766	0.743	0.444	0.427	-0.015
posthoc (gpt-3.5)	scalar	0.628	0.508	0.809	0.861	0.801	0.790	0.233	0.293	-0.019
cot (gpt-4)	gpt-3.5	0.442	0.358	0.808	0.772	0.659	0.854	0.330	0.301	0.046
cot (gpt-4)	gpt-4	0.566	0.476	0.843	0.914	0.807	0.882	0.348	0.331	0.040
cot (gpt-4)	Llama-3.3	0.550	0.463	0.842	0.879	0.763	0.868	0.329	0.300	0.026
cot (gpt-4)	presupp_flip	0.406	0.338	0.833	0.764	0.645	0.844	0.358	0.307	0.011
cot (gpt-4)	contextual	0.674	0.603	0.894	0.948	0.818	0.863	0.273	0.215	-0.031
cot (gpt-4)	lexical	0.466	0.380	0.815	0.804	0.712	0.818	0.338	0.332	0.003
cot (gpt-4)	scalar	0.543	0.458	0.868	0.880	0.769	0.858	0.337	0.311	-0.010
posthoc (gpt-4)	gpt-3.5	0.456	0.364	0.799	0.768	0.660	0.860	0.312	0.296	0.061
posthoc (gpt-4)	gpt-4	0.590	0.483	0.818	0.906	0.801	0.884	0.316	0.318	0.066
posthoc (gpt-4)	Llama-3.3	0.543	0.343	0.873	0.838	0.732	0.882	0.295	0.389	0.009
posthoc (gpt-4)	presupp_flip	0.553	0.432	0.781	0.831	0.717	0.845	0.278	0.285	0.064
posthoc (gpt-4)	contextual	0.671	0.588	0.877	0.948	0.805	0.849	0.277	0.217	-0.028
posthoc (gpt-4)	lexical	0.332	0.358	0.830	0.698	0.727	0.856	0.366	0.369	0.026
posthoc (gpt-4)	scalar	0.658	0.546	0.830	0.820	0.792	0.846	0.162	0.246	0.016
cot (Llama-3.3)	gpt-3.5	0.478	0.395	0.826	0.768	0.673	0.877	0.290	0.278	0.051
cot (Llama-3.3)	gpt-4	0.607	0.511	0.842	0.908	0.799	0.880	0.301	0.288	0.038
cot (Llama-3.3)	Llama-3.3	0.523	0.445	0.851	0.897	0.774	0.864	0.374	0.329	0.013
cot (Llama-3.3)	presupp_flip	0.365	0.312	0.854	0.762	0.655	0.860	0.397	0.343	0.006
cot (Llama-3.3)	contextual	0.763	0.349	0.457	0.942	0.628	0.562	0.179	0.289	0.105
cot (Llama-3.3)	lexical	0.434	0.334	0.770	0.793	0.688	0.889	0.359	0.354	0.119
cot (Llama-3.3)	scalar	0.593	0.403	0.848	0.765	0.723	0.815	0.172	0.320	-0.033
posthoc (Llama-3.3)	gpt-3.5	0.463	0.361	0.780	0.773	0.662	0.856	0.310	0.301	0.076
posthoc (Llama-3.3)	gpt-4	0.604	0.491	0.813	0.904	0.791	0.875	0.300	0.300	0.062
posthoc (Llama-3.3)	Llama-3.3	0.529	0.424	0.801	0.899	0.778	0.865	0.379	0.354	0.064
posthoc (Llama-3.3)	presupp_flip	0.416	0.356	0.856	0.719	0.721	0.862	0.303	0.365	0.006
posthoc (Llama-3.3)	contextual	0.741	0.237	0.319	0.942	0.600	0.318	0.201	0.363	-0.001
posthoc (Llama-3.3)	lexical	0.434	0.355	0.817	0.622	0.694	0.836	0.188	0.339	0.019
posthoc (Llama-3.3)	scalar	0.653	0.462	0.861	0.821	0.788	0.838	0.168	0.326	-0.023

Table 6: Per-condition results for counterfactual simulatability (*judge*: gpt-3.5-turbo). Improvements are Phase 2 minus Phase 1.

Label Flip Rate The label flip rate measures how often a perturbed example changes the model’s predicted label relative to the original instance (Ge et al., 2021; Nguyen et al., 2024; Bhattacharjee et al., 2023; Wang et al., 2025). Intuitively, counterfactuals that preserve the surface form of Q but flip the label are considered more challenging and informative, because they probe decision boundaries rather than trivial paraphrases.

We compute the proportion of Q' for which the model’s answer differs from its answer to Q , and later relate this flip rate to the observed gains from showing explanations in our simulation test.

Textual Similarity Textual similarity captures how close Q' remains to Q in form and content (Madaan et al., 2020). High similarity is often taken to indicate a “minimal” counterfactual, where small edits lead to different model behavior, whereas low

similarity suggests that the perturbation may be drifting too far from the original instance.

Following prior work on counterfactual evaluation, we quantify textual similarity using a normalized word-level Levenshtein distance (Ross et al., 2021; Treviso et al., 2023; Wang et al., 2025).

Table 8 reports the Label Flip Rate and textual similarity results across counterfactual question types. However, both metrics only approximate counterfactual quality and ignore whether humans can actually leverage explanations to anticipate model behavior, which is what our simulation-based evaluation is designed to capture.

H Full results for in-context bias and flip-back behavior

In this appendix, we report the full results on in-context bias and flip-back behavior for explana-

expl	cfq_type	P1_guess_rate	P1_acc	P1_sel_acc	P2_guess_rate	P2_acc	P2_sel_acc	Δ guess_rate	Δ acc	Δ sel_acc
cot (gpt-3.5)	gpt-3.5	0.607	0.445	0.733	0.764	0.586	0.768	0.157	0.141	0.035
cot (gpt-3.5)	gpt-4	0.740	0.541	0.731	0.859	0.650	0.757	0.119	0.109	0.026
cot (gpt-3.5)	Llama-3.3	0.706	0.528	0.748	0.841	0.639	0.760	0.135	0.111	0.012
cot (gpt-3.5)	PresupFlip	0.664	0.401	0.604	0.854	0.571	0.587	0.190	0.170	-0.017
cot (gpt-3.5)	Contextual	0.941	0.718	0.763	0.949	0.730	0.769	0.008	0.012	0.006
cot (gpt-3.5)	Lexical	0.418	0.301	0.720	0.586	0.462	0.766	0.168	0.161	0.046
cot (gpt-3.5)	Scalar	0.670	0.516	0.770	0.788	0.719	0.785	0.118	0.203	0.015
posthoc (gpt-3.5)	gpt-3	0.689	0.479	0.752	0.774	0.589	0.761	0.085	0.110	0.009
posthoc (gpt-3.5)	gpt-4	0.732	0.506	0.691	0.847	0.635	0.750	0.115	0.129	0.059
posthoc (gpt-3.5)	Llama-3.3	0.695	0.499	0.718	0.838	0.635	0.758	0.143	0.136	0.040
posthoc (gpt-3.5)	PresupFlip	0.694	0.423	0.610	0.856	0.615	0.617	0.162	0.192	0.007
posthoc (gpt-3.5)	Contextual	0.947	0.676	0.714	0.949	0.690	0.708	0.002	0.014	-0.006
posthoc (gpt-3.5)	Lexical	0.425	0.300	0.705	0.595	0.478	0.743	0.179	0.178	0.038
posthoc (gpt-3.5)	Scalar	0.652	0.484	0.743	0.795	0.623	0.771	0.143	0.139	0.028
cot (gpt-4)	gpt-3.5	0.591	0.488	0.826	0.799	0.693	0.868	0.208	0.205	0.042
cot (gpt-4)	gpt-4	0.768	0.660	0.860	0.919	0.809	0.881	0.151	0.149	0.021
cot (gpt-4)	Llama-3.3	0.707	0.602	0.851	0.896	0.782	0.873	0.189	0.180	0.022
cot (gpt-4)	PresupFlip	0.638	0.530	0.831	0.880	0.716	0.814	0.242	0.186	-0.017
cot (gpt-4)	Contextual	0.963	0.822	0.853	0.967	0.832	0.860	0.004	0.010	0.007
cot (gpt-4)	Lexical	0.425	0.348	0.819	0.569	0.499	0.859	0.144	0.151	0.040
cot (gpt-4)	Scalar	0.392	0.329	0.839	0.661	0.563	0.853	0.269	0.234	0.014
posthoc (gpt-4)	gpt-3.5	0.592	0.490	0.828	0.796	0.689	0.865	0.204	0.199	0.037
posthoc (gpt-4)	gpt-4	0.745	0.622	0.835	0.910	0.807	0.886	0.165	0.185	0.051
posthoc (gpt-4)	Llama-3.3	0.697	0.588	0.843	0.888	0.782	0.880	0.191	0.194	0.037
posthoc (gpt-4)	PresupFlip	0.668	0.559	0.837	0.878	0.804	0.889	0.210	0.245	0.052
posthoc (gpt-4)	Contextual	0.955	0.801	0.839	0.970	0.819	0.844	0.015	0.018	0.005
posthoc (gpt-4)	Lexical	0.344	0.306	0.860	0.528	0.561	0.886	0.184	0.255	0.026
posthoc (gpt-4)	Scalar	0.650	0.553	0.850	0.832	0.761	0.878	0.182	0.208	0.028
cot (Llama-3.3)	gpt-3.5	0.626	0.528	0.843	0.813	0.706	0.869	0.187	0.178	0.026
cot (Llama-3.3)	gpt-4	0.796	0.682	0.857	0.935	0.821	0.878	0.139	0.139	0.021
cot (Llama-3.3)	Llama-3.3	0.751	0.643	0.855	0.921	0.799	0.868	0.179	0.156	0.013
cot (Llama-3.3)	PresupFlip	0.652	0.564	0.865	0.896	0.755	0.843	0.244	0.191	-0.022
cot (Llama-3.3)	Contextual	0.959	0.449	0.468	0.968	0.455	0.470	0.009	0.006	0.002
cot (Llama-3.3)	Lexical	0.445	0.342	0.768	0.627	0.522	0.811	0.182	0.180	0.043
cot (Llama-3.3)	Scalar	0.652	0.532	0.816	0.832	0.696	0.843	0.180	0.164	0.027
posthoc (Llama-3.3)	gpt-3.5	0.590	0.480	0.813	0.816	0.706	0.865	0.226	0.226	0.052
posthoc (Llama-3.3)	gpt-4	0.744	0.623	0.837	0.929	0.815	0.877	0.185	0.192	0.040
posthoc (Llama-3.3)	Llama-3.3	0.711	0.588	0.826	0.939	0.822	0.893	0.228	0.234	0.067
posthoc (Llama-3.3)	PresupFlip	0.666	0.577	0.867	0.886	0.799	0.848	0.220	0.222	-0.019
posthoc (Llama-3.3)	Contextual	0.955	0.296	0.310	0.959	0.331	0.345	0.004	0.035	0.035
posthoc (Llama-3.3)	Lexical	0.436	0.316	0.829	0.631	0.548	0.844	0.195	0.232	0.015
posthoc (Llama-3.3)	Scalar	0.661	0.565	0.855	0.881	0.794	0.866	0.220	0.229	0.011

Table 7: Per-condition results for counterfactual simulatability (*judge*: gpt-4.1). Improvements are Phase 2 minus Phase 1.

cfq_type	cfq_type	TS	LFR_gold	LFR_orig
model_generated	gpt-3.5	0.777	0.487	0.497
model_generated	gpt-4	0.713	0.535	0.525
model_generated	Llama-3.3	0.634	0.503	0.517
prag_cf	presupp_flip	1.480	0.629	0.534
prag_cf	contextual	1.255	0.350	0.349
prag_cf	lexical	0.203	0.450	0.401
prag_cf	scalar	0.184	0.437	0.397

Table 8: Label Flip Rate and Textual Similarity scores for model-generated and pragmatics-based counterfactual questions.

tions generated by GPT-3.5 and Llama-3.3-70B in Table 9. The patterns mirror those observed for GPT-4: in-context bias generally decreases from Phase 1 to Phase 2, with especially pronounced reductions when conditioning on instances where the original answer is incorrect.

I Human-LLM Judge Agreement on the Overlapping Subset

To assess how well LLM-as-a-judge approximates human judgments, we compute raw yes/no agreement between human predictions and the LLM judge on the overlapping subset of items. In Table 10 we report agreement separately for Phase 1 (without explanations) and Phase 2 (with explanations). In addition, we provide label distributions for both raters and a breakdown by counterfactual type. These statistics complement the main text by documenting where human and LLM judge decisions align or diverge across transformation types. The full results are shown in

expl_type	cfq_type	gpt-3.5		Llama-3.3	
		Δ bias_orig	Δ bias_orig_wrong	Δ bias_orig	Δ bias_orig_wrong
<i>Weaker Evaluator</i>					
cot	gpt-3.5	-0.072	-0.093	-0.044	-0.077
cot	gpt-4	-0.087	-0.139	-0.025	-0.013
cot	Llama-3.3	-0.087	-0.078	0.003	-0.086
cot	presupp_flip	-0.047	-0.060	-0.162	-0.015
cot	contextual	-0.023	-0.045	-0.015	-0.073
cot	lexical	-0.079	-0.145	-0.055	-0.063
cot	scalar	-0.029	-0.054	-0.093	-0.108
posthoc	gpt-3.5	-0.048	-0.091	-0.110	-0.292
posthoc	gpt-4	-0.015	-0.106	-0.080	-0.283
posthoc	Llama-3.3	-0.043	-0.143	-0.063	-0.321
posthoc	presupp_flip	-0.164	-0.013	-0.197	-0.011
posthoc	contextual	-0.022	-0.038	-0.018	-0.062
posthoc	lexical	-0.065	-0.119	-0.065	-0.083
posthoc	scalar	-0.067	-0.083	-0.055	-0.069
<i>Stronger Evaluator</i>					
cot	gpt-3.5	-0.090	-0.241	-0.060	-0.035
cot	gpt-4	-0.057	-0.201	-0.037	-0.041
cot	Llama-3.3	-0.087	-0.078	-0.059	-0.039
cot	presupp_flip	-0.073	-0.014	-0.052	0.011
cot	contextual	-0.017	-0.084	-0.016	0.008
cot	lexical	-0.061	-0.073	-0.061	-0.033
cot	scalar	-0.049	-0.010	-0.071	-0.053
posthoc	gpt-3.5	-0.041	-0.010	-0.092	0.028
posthoc	gpt-4	-0.035	-0.002	-0.071	0.039
posthoc	Llama-3.3	-0.059	-0.010	-0.095	0.040
posthoc	presupp_flip	-0.085	0.008	-0.067	0.017
posthoc	contextual	-0.036	0.024	-0.045	0.033
posthoc	lexical	-0.065	0.008	-0.083	-0.008
posthoc	scalar	-0.043	-0.014	-0.060	-0.034

Table 9: Change in in-context bias from Phase 1 to Phase 2 across counterfactual question sets, reported separately for gpt-3.5 and Llama-3.3-70B. Negative values indicate reduced bias (more flip-back). Positive numbers are marked in blue.

cfq_type	Agreement (P1)	Agreement (P2)
<i>Weaker Evaluator</i>		
contextualization	0.776	0.820
gpt_3.5	0.531	0.540
gpt_4	0.469	0.440
lexical_substitution	0.735	0.740
llama	0.612	0.560
pre_flip	0.347	0.300
scalar_adjustment	0.735	0.740
<i>Stronger Evaluator</i>		
contextualization	0.660	0.480
gpt_3.5	0.560	0.600
gpt_4	0.440	0.540
lexical_substitution	0.600	0.480
llama	0.640	0.740
pre_flip	0.340	0.600
scalar_adjustment	0.720	0.520

Table 10: Raw yes/no agreement between humans and the LLM-as-a-judge by counterfactual question type on the overlapping subset.

Phase 1 Instructions

Thank you for participating in this task!

For each HIT, you will see:

- one yes/no **Starter Question**,
- a **Robot's Answer** to the Starter Question, and
- a **Follow-up Question**.

Your tasks are:

1. Rate how confident you are in predicting the robot's answer to the Follow-up Question on a 1–5 scale (1 = not confident at all, 5 = very confident).
2. Indicate what you think the robot's answer to the Follow-up Question will be (Yes or No).
3. Briefly explain why you made this judgement.

Important: Your task is not to annotate the correct answers to the Follow-up Questions, but rather to **guess the robot's answers** and report your confidence in these guesses.

Original Question (Starter Question):

Can you substitute the pins in a bowling alley lane with Dustin Hoffman's Oscars?

Robot's Answer to the Starter Question:

no

Follow-up Question:

If Dustin Hoffman's Oscars were modified to meet the size, shape, and weight specifications of standard bowling pins, could you substitute the pins in a bowling alley lane with them?

1. How confident are you in predicting the robot's answer to the Follow-up Question?

Please rate your confidence on a 1–5 scale, where:

1 = Not confident at all, 5 = Very confident, and values in between indicate intermediate confidence.

- | | | | | |
|-----------------------|-----------------------|-----------------------|-----------------------|-----------------------|
| <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> | <input type="radio"/> |
| 1 | 2 | 3 | 4 | 5 |
| Not confident at all | Slightly confident | Moderately confident | Confident | Very confident |

2. What is your guess about the robot's answer to the Follow-up Question?

- The robot will answer: Yes
- The robot will answer: No

3. Please briefly explain why you made this judgment.

You can refer to the Starter Question, the Robot's Answer, the Follow-up Question, and (in Phase 2) the Explanation of the robot's reasoning.

Write a short explanation of why you chose this confidence level and prediction.

Next

Figure 4: User study interface: instruction for phase 1.

Phase 2 Instructions

Thank you for continuing with this task!

For each HIT, you will now see:

- the same yes/no **Starter Question**,
- the **Robot's Answer** to the Starter Question,
- the **Follow-up Question**, and
- an **Explanation** of the robot's reasoning.

Your tasks are:

1. Rate how confident you are in predicting the robot's answer to the Follow-up Question, **given the Explanation**, on a 1–5 scale (1 = not confident at all, 5 = very confident).
2. Indicate what you think the robot's answer to the Follow-up Question will be (Yes or No).
3. Briefly explain why you made this judgement.

Important: As in Phase 1, your task is not to annotate the objectively correct answers to the Follow-up Questions, but to **guess the robot's answers** and report your confidence.

Original Question (Starter Question):

Can you substitute the pins in a bowling alley lane with Dustin Hoffman's Oscars?

Robot's Answer to the Starter Question:

no

Follow-up Question:

If Dustin Hoffman's Oscars were modified to meet the size, shape, and weight specifications of standard bowling pins, could you substitute the pins in a bowling alley lane with them?

Explanation of the Robot's Reasoning:

Trace of reasoning: Bowling pins and Oscars are fundamentally different objects with different purposes. Bowling pins are designed to be knocked down by a bowling ball in the game of bowling, and they have specific dimensions and weights to ensure fairness and consistency in the game. Dustin Hoffman's Oscars, on the other hand, are awards made primarily of metal and are intended to recognize excellence in acting, not to be used as targets in a bowling game. Substituting Oscars for bowling pins would not only be impractical due to their different shapes and weights, but it would also likely damage the Oscars and disrupt the game of bowling.

So the answer is no.

1. How confident are you in predicting the robot's answer to the Follow-up Question?

Please rate your confidence on a 1–5 scale, where:

1 = Not confident at all, 5 = Very confident, and values in between indicate intermediate confidence.

1 2 3 4 5
Not confident at all Slightly confident Moderately confident Confident Very confident

2. What is your guess about the robot's answer to the Follow-up Question?

- The robot will answer: Yes
 The robot will answer: No

3. Please briefly explain why you made this judgment.

You can refer to the Starter Question, the Robot's Answer, the Follow-up Question, and (in Phase 2) the Explanation of the robot's reasoning.

Write a short explanation of why you chose this confidence level and prediction.

Next

Figure 5: User study interface: instruction for phase 2.