A CAUSAL LENS FOR EVALUATING FAITHFULNESS METRICS

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027 028 029 Paper under double-blind review

Abstract

The increasing capabilities of Large Language Models (LLMs) have made natural language explanations a promising alternative to traditional feature attribution methods for model interpretability. However, while these explanations may seem plausible, they can fail to reflect the model's underlying reasoning faithfully. The idea of faithfulness is critical for assessing the alignment between the explanation and the model's true decision-making mechanisms. Although several faithfulness metrics have been proposed, they lack a unified evaluation framework. To address this limitation, we introduce CAUSAL DIAGNOSTICITY, a new evaluation framework for comparing faithfulness metrics in natural language explanations. Our framework extends the idea of diagnosticity to the faithfulness metrics for natural language explanations by using model editing to generate faithful and unfaithful explanation pairs. We introduce a benchmark consisting of three tasks: fact-checking, analogy, and object counting, and evaluate a diverse set of faithfulness metrics, including post-hoc explanation-based and chain-of-thought (CoT)-based methods. Our results show that while CC-SHAP significantly outperforms other metrics, there is substantial room for improvement. This work lays the foundation for future research in developing more faithful natural language explanations, highlighting the need for improved metrics and more reliable interpretability methods in LLMs.

1 INTRODUCTION

Recent advancements in Large Language Models (LLMs) have opened up new possibilities in terms
of explainability. These models' evolving capabilities have made natural language explanations
preferable over traditional feature attribution methods. Additionally, most LLMs can provide explanations
for their predictions without much additional cost (Wei et al., 2022). While these natural
language-based explanations can be valuable, practitioners must exercise caution before relying on
them. Despite appearing plausible, these explanations may not accurately reflect the model's inner
reasoning mechanism, potentially leading practitioners astray (Turpin et al., 2023).

The idea of faithfulness aims to assess how accurately explanations reflect the true reasoning mech-040 anism of the model. While numerous methods have been proposed to measure faithfulness for natural language-based explanations, they are criticized for not adequately considering the model's 041 inner workings, relying instead on simplistic consistency measures (Parcalabescu & Frank, 2023). 042 Furthermore, while many faithfulness metrics have been developed, currently there are no reliable 043 evaluation frameworks for comparing them. To address this gap in the field, we introduce a new 044 evaluation framework, CAUSAL DIAGNOSTICITY, along with a new benchmark for comparing various faithfulness metrics. Our framework extends the notion of *diagnosticity* (Chan et al., 2022b), 046 which measures how often a faithfulness metric favors faithful explanations over unfaithful ones, and 047 applies it to faithfulness metrics for natural language explanations. We investigate model editing 048 approaches for causally generating faithful and unfaithful explanation pairs and evaluate diagnosticity through three tasks. These tasks include (1) a fact-checking task, (2) an analogy task, and (3) an object counting task. Figure 1 shows an overview of our framework. We evaluate a diverse set of 051 faithfulness metrics, including post-hoc explanation-based and chain-of-thought (CoT)-based metrics: Counterfactual Edits (Atanasova et al., 2023), Simulatability, metrics based on corrupting CoT 052 explanations (Lanham et al., 2023), and CC-SHAP (Parcalabescu & Frank, 2023). Our evaluation shows that while most metrics fail to achieve high diagnosticity scores, CC-SHAP significantly

054 outperforms the others, though there is still room for improvement in developing better metrics. Our key contributions are:

- A new framework for evaluating faithfulness metrics for natural language explanations
- A new dataset with three tasks for evaluating these metrics
- A comprehensive evaluation of prominent faithfulness metrics to guide practitioners in selecting the most reliable metrics

062 By offering a test bed for evaluating faithfulness 063 metrics for natural language explanations, this 064 study exposes the limitations of existing met-065 rics and highlights the need for improved ones. In this role, our work serves as the first step in 066 a broader research initiative aimed at develop-067 ing more faithful natural language explanations. 068 With a test bed in place and an assessment of the 069 current state of existing metrics, future research should focus on developing better faithfulness 071 metrics and, subsequently, models that generate 072 more faithful explanations. 073

2 BACKGROUND

076 Faithfulness Faithfulness measures the ex-077 tent to which explanations reflect the true reasoning mechanisms of models. Formally, let 079 M_{θ} denote a LLM parameterized by θ , operating on a token set \mathcal{V} such that $M_{\theta}(t^{\text{in}}) =$ 081 t^{out} , where $t^{\text{in}} = \langle t_1^{\text{in}}, t_2^{\text{in}}, \dots, t_{N_{\text{in}}}^{\text{in}} \rangle$ and $t^{\text{out}} =$ $\langle t_1^{\text{out}}, t_2^{\text{out}}, \dots, t_{N_{\text{out}}} \rangle$; $t_i^{\text{in}}, t_i^{\text{out}} \in \mathcal{V}$, N_{in} and N_{out} 083 represent the input and output sequence lengths. 084 The input and output sequences can take many 085 forms. For the simplest case $t^{in} = x$ and $t^{\text{out}} = y$ where (x, y) is an input and output pair for any task. With a proper prompt pro-087



Figure 1: Our framework consists of three stages: (1) Model Editing: applying counterfactual edits to the models; (2) Explanation Generation: generating faithful and unfaithful explanation pairs using the edited models, or synthetically generating such pairs based on the edits; (3) Diagnosticity Evaluation: assessing the chosen faithfulness metric with one of the edited models using the faithful-unfaithful explanation pairs. Diagnostic faithfulness metrics should assign a higher faithfulness score to the faithful explanation than to the unfaithful one.

vided, the output can take the form $t^{\text{out}} = y \oplus \varepsilon$ for post-hoc explanations or $t^{\text{out}} = \varepsilon \oplus y$ for 880 chain-of-thought (CoT) explanations, where ε is the explanation and \oplus represents the concatenation 089 of two sequences. 090

Based on these definitions, we define a faithfulness metric \mathcal{F} as a scalar valued function:

094

091

060

061

074

075

 $\mathcal{F}(\boldsymbol{x},\boldsymbol{y},\boldsymbol{\varepsilon},\boldsymbol{\theta}) = s$ (1)

where $s \in \mathbb{R}$ represents the level of faithfulness of the explanation ε , for the given input-output pair (x, y) and the model parameterized by θ . While explanations can take different forms, such as 096 importance scores, here we focus on text-based explanations.

098 2.1 FAITHFULNESS METRICS 099

100 In this study, we focus on seven prominent faithfulness metrics: (1) Counterfactual Edits (Atanasova 101 et al., 2023), (2) Simulatability, metrics based on corrupting CoT explanations (Lanham et al., 2023) 102 (including (3) Early Answering, (4) Adding Mistakes, (5) Paraphrasing, and (6) Filler Tokens), and 103 (7) CC-SHAP (Parcalabescu & Frank, 2023). While Simulatability and Counterfactual Edits are designed for post-hoc explanations, the others are tailored for CoT explanations. Notably, CC-SHAP 104 is applicable to both types of explanations. 105

- 106
- **Counterfactual Edits** Atanasova et al. (2023) propose a new metric based on the rationale that an 107 explanation is unfaithful if the model changes its prediction after a counterfactual intervention to the

input, while the explanation fails to reflect the intervention. A significant limitation of this approach is the need to train a separate neural editor for each model-dataset pair to make such counterfactual interventions. Instead, we follow their random baseline based on the same rationale, where they insert a random adjective before a noun or a random adverb before a verb, as Parcalabescu & Frank (2023) do. In this approach, an explanation is considered unfaithful if the prediction changes after word insertion and the explanation fails to mention the inserted words.

Simulatability Simulatability is based on measuring the predictiveness of explanations regarding the label (Doshi-Velez & Kim, 2017; Hase & Bansal, 2020; Hase et al., 2020; Wiegreffe et al., 2020; Chan et al., 2022a). A faithful explanation should convey sufficient information about the model's reasoning so that a simulator can predict the model's outputs when provided with the input and explanations. We follow Chan et al. (2022a)'s definition of simulatability as $\mathbb{1}_{S}(y_i \mid x_i, \varepsilon_i) - \mathbb{1}_{S}(y_i \mid$ x_i), where $\mathbb{1}_{S}(b \mid a)$ is the accuracy of S in predicting b given a.

Corrupting CoT Lanham et al. (2023) focus on the unfaithfulness of Chain-of-Thought (CoT)
 explanations. They propose four types of corruption: (1) *Early Answering*, which involves truncating
 the CoT to get an early answer; (2) *Adding Mistakes*, where a helper language model introduces
 mistakes into the original CoT, and the original model itself regenerates the remaining part; (3)
 Paraphrasing, which involves paraphrasing the original CoT and regenerating the rest; and (4) *Filler Tokens*, where the original CoT is replaced with ellipses. If a corruption does not change the original
 prediction, then the explanation is not faithful.

CC-SHAP Parcalabescu & Frank (2023) measure faithfulness by testing the alignment of input contributions to prediction and explanation using SHAP (Lundberg & Lee, 2017) importance scores. For each example, they first compute importance scores with respect to the prediction for each token in the input. Then, they compute importance scores with respect to each token in the explanation and aggregate them. Finally, they measure the convergence of the two distributions of importance scores. Their method is applicable to both post-hoc and Chain-of-Thought (CoT) explanations.

133 134

135

2.2 MODEL EDITING

136 In our framework for evaluating faithfulness metrics, we use model editing approaches to generate 137 faithful-unfaithful explanation pairs by modifying specific facts within LLMs. The need for model editing approaches stems from the fact that the knowledge of LLMs can become outdated over time. 138 For example, after a new election, they might present outdated knowledge about the head of a state. 139 An array of model editing methods has been proposed to address this problem in a feasible way, 140 allowing LLMs to stay up-to-date without altering unrelated knowledge (Cohen et al., 2024; Zhang 141 et al., 2024; Patil et al., 2023; Geva et al., 2023; Gupta et al., 2023; Hartvigsen et al., 2023; Hase 142 et al., 2023; Tan et al., 2024; Yu et al., 2023; Zheng et al., 2023; Meng et al., 2022; Mitchell et al., 143 2022). Such techniques operate on knowledge triplets consisting of subject s, object o, and relation 144 r. For instance, they can update (s = Donald Trump, r = is the president of, o = the United States) 145 to (s = Joe Biden, r = is the president of, o = the United States) while keeping other information 146 unchanged. In this study, we explore two model editing methods: (1) MEMIT (Meng et al., 2023), 147 a locate-then-edit approach, which enables successful bulk edits, and (2) In-Context Knowledge 148 Editing, a memory-based alternative, (Zheng et al., 2023).

3 Method

150 151 152

149

METHO

Our CAUSAL DIAGNOSTICITY framework is inspired by the idea of *diagnosticity*. We begin by summarizing the idea of diagnosticity in 3.1, which was introduced by Chan et al. (2022b) for evaluating faithfulness metrics tailored for feature attribution methods. Next, in 3.2, we introduce CAUSAL DIAGNOSTICITY, describing how it builds on diagnosticity and extends it to natural language explanations in a causal manner by incorporating edited models.

- 157 158
- 3.1 DIAGNOSTICITY
- 159

An active body of research has explored accurately measuring faithfulness (Jacovi & Goldberg, 2020). This has led to a multiplicity of faithfulness metrics, and exposed the need of a framework to evaluate faithfulness metrics. For evaluating different faithfulness evaluation metrics, we adapt the

notion of *diagnosticity* proposed by Chan et al. (2022b). Diagnosticity is the measure of how often a faithfulness metric prefers faithful rather than unfaithful explanations.

Following the notation used by Chan et al. (2022b), formally we denote "*u* is more faithful than *v*" as $u \succ v$, given that *u* and *v* are explanations, regardless of their form (e.g., text, heatmap). Additionally, we denote the statement " \mathcal{F} considers *u* more faithful than *v*" as $u \succ_{\mathcal{F}} v$. Then, the diagnosticity of the metric \mathcal{F} is defined as:

$$D(\mathcal{F}) = P(u \succ_{\mathcal{F}} v | u \succ v) \tag{2}$$

Based on estimates in Chan et al. (2022b), we use the following formula to calculate diagnosticity:

$$D(\mathcal{F}) \approx \frac{1}{|Z|} \sum_{(u_i, v_i) \in Z} \mathbb{1}(u_i \succ_{\mathcal{F}} v_i)$$
(3)

where Z is a dataset consisting of pairs (u_i, v_i) of faithful explanations (u_i) and unfaithful explanations (v_i) which correspond to input-output pairs $(\boldsymbol{x}_i, \boldsymbol{y}_i)$.

If higher faithfulness scores represent more faithful explanations, we can revise our notation to:

$$D(\mathcal{F}) \approx \frac{1}{|Z|} \sum_{(u_i, v_i) \in Z} \mathbb{1}(\mathcal{F}(u_i; \boldsymbol{x}_i, \boldsymbol{y}_i, \boldsymbol{\theta}) > \mathcal{F}(v_i; \boldsymbol{x}_i, \boldsymbol{y}_i, \boldsymbol{\theta}))$$
(4)

185 3.2 CAUSAL DIAGNOSTICITY

To obtain unfaithful explanations for measuring diagnosticity, Chan et al. (2022b) use random feature attribution scores. While random scores can work for structured explanations like feature attributions – since they still follow the intended format – this approach is not straightforward for natural language explanations. Random text cannot function as a meaningful explanation and cannot ensure unfaithfulness in a coherent way.

To address this limitation, we introduce CAUSAL DIAGNOSTICITY, a framework for evaluating 192 faithfulness metrics through diagnosticity, by generating unfaithful explanations using model editing 193 methods. In CAUSAL DIAGNOSTICITY, unfaithful explanations are produced by modifying a 194 model's internal knowledge. For example, consider the capitalOf relation with the query "Is 195 Paris the capital of France?" and a model that correctly associates this question to the knowledge 196 (s = Paris, r = is the capital of, o = France). By altering the model's internal knowledge, we create 197 two variations where the subject s is replaced with Berlin or London. Both modified models should 198 answer "No" to the original question but for different reasons: "No, because Berlin is the capital 199 of France." and "No, because London is the capital of France." In particular, each of these two explanations should be unfaithful to the model that generated the other explanation. 200

Formally, let y_i be the prediction for the input x_i while $\overline{\theta}$ and $\widetilde{\theta}$ be the parameters of the altered models. $\overline{\theta}$ generates the explanation $\overline{\varepsilon}_i$ and $\widetilde{\theta}$ generates the explanation $\widetilde{\varepsilon}_i$. Then we modify diagnosticity definition as follows:

205 206

207 208

169

170 171

172 173 174

175 176

177

178 179

$$D(\mathcal{F}) = \frac{1}{|Z|} \sum_{(\bar{\boldsymbol{\varepsilon}}_i, \tilde{\boldsymbol{\varepsilon}}_i) \in Z} \mathbb{1}(\mathcal{F}(\bar{\boldsymbol{\varepsilon}}_i; \boldsymbol{x}_i, \boldsymbol{y}_i, \bar{\boldsymbol{\theta}}) > \mathcal{F}(\tilde{\boldsymbol{\varepsilon}}_i; \boldsymbol{x}_i, \boldsymbol{y}_i, \bar{\boldsymbol{\theta}}))$$
(5)

Models $\overline{\theta}$ and $\overline{\theta}$ are edited such that $\overline{\varepsilon}_i$ is faithful to $\overline{\theta}$, while $\widetilde{\varepsilon}_i$ is unfaithful to $\overline{\theta}$. Depending on the scenario, $\overline{\theta}$ and $\overline{\theta}$ can be used interchangeably. Continuing with our running examples of capital cities, each generated explanation is faithful to its own model but unfaithful to the other model. In these cases, either model can be used to compute Equation 5 by swapping $\overline{\varepsilon}_i$ and $\widetilde{\varepsilon}_i$, as the faithfulness dichotomy holds regardless. However, in certain scenarios, one of the explanations may be faithful to both models, limiting the flexibility of choosing models arbitrarily. For instance, in the Analogy task of our benchmark (see Figure 2), the capitalOf relation is held by only one model, whereas the cityOf relation is valid for both models. As a result, the corresponding explanation is faithful to both



Figure 2: Summary of three tasks with example questions and answers, along with explanations from the edited models: (a) Fact Check task, (b) Analogy task, and (c) and (d) Object Counting task, featuring two different types of questions. The blue and orange boxes represent models parameterized by $\bar{\theta}$ and $\tilde{\theta}$, respectively, while the dashed boxes within them indicate the counterfactual knowledge injected into the model through editing. Gray boxes below each model display their output, consisting of the answer (y) and explanation ($\bar{\varepsilon}$ or $\bar{\varepsilon}$). Although both model pairs produce the same answers, their reasoning differs, as shown by the explanations that follow the answers.

248 249

250

251

253

254

models. Additionally, the original model θ can be used as long as it satisfies respective faithfulness conditions of the explanation pairs. Nevertheless, we opt to create two edited variants of the models, even when reflecting factual knowledge, to guarantee that all conditions are met.

4 TASKS

For evaluating different faithfulness metrics, we include three controlled tasks in the CAUSAL DIAGNOSTICITY framework: (1) a fact-checking task, (2) an analogy task, and (3) an object counting 256 task. Across all tasks, we aim to test the causal diagnosticity of faithfulness metrics by using 257 counterfactual models and their corresponding faithful and unfaithful explanations. While we expect 258 the altered models to reason differently, their explanations may not explicitly reference the altered 259 aspect. Since our focus is on evaluating faithfulness metrics, we ensure the faithfulness situation of 260 the explanations by synthetically generating explanations that emphasize the differences between 261 the models. Figure 2 provides an overview of these tasks, including example inputs, outputs, and 262 explanations.

263 264

265

4.1 FACT CHECK TASK

Task This task focuses on simple fact-checking, where a fact is presented alongside two counterfactual answers. For any relation (s_i, r_i, o_i) , we present a question that checks its correctness, accompanied by two counterfactuals: $(s_i, r_i, \overline{o_i})$ and $(s_i, r_i, \widetilde{o_i})$. These counterfactuals yield the same answer but are based on different reasoning. For instance, given the knowledge triplet $(s_i = "Rihanna", r_i = "is", o_i = "a singer")$, the corresponding question would be "Is Rihanna a singer?" Let the counterfactual objects be $\bar{o}_i =$ "researcher" and $\tilde{o}_i =$ "lawyer". Both counterfactuals would result in the answer "No," but for different reasons.

273 **Dataset** We construct our dataset using the COUNTERFACT dataset (Meng et al., 2022), which 274 consists of knowledge triplets. While COUNTERFACT includes prompts representing knowledge 275 triplets, we use an LLM (Mistral-7B-Instruct-v0.2) to convert those statements into yes/no 276 questions. Next, for each object o_i , we fetch sibling entities from WikiData to be used as new 277 counterfactuals. Finally, we generate synthetic explanations corresponding to each counterfactual. 278 For example, the corresponding explanation $\bar{\varepsilon}_i$ would be "Joe Biden is a researcher, not the president 279 of the United States" for \bar{o}_i . Further details about the dataset generation process, including prompts, can be found in Appendix A. 280

281 282

283

4.2 ANALOGY TASK

Task This task is based on analogies exploiting the hierarchical structure between two relations where $r_1
ightharpoondown r_2$ holds. For any (s_i, o_i) and (s_j, o_j) , there exist (s_i, r_1, o_i) and (s_j, r_2, o_j) such that $r_1
ightharpoondown r_2$. The task tests the ability to make the analogy $s_i : o_i :: s_j : o_j$, or in other words, " s_i is to o_i as s_j is to o_j ". We choose r_1 and r_2 as $r_{capitalof}$ and r_{cityof} relations, respectively. For instance, we test "Paris is to France as Berlin is to Germany." We corrupt one of the models so that the relation $r_{capitalof}$ is no longer valid while the relation r_{cityof} holds. Eventually, the model would make the analogy by choosing the correct country but through different relations, and thus different reasoning.

291 **Dataset** First, we collect a list of countries and cities¹, then select one capital and one non-capital 292 city for each country. We randomly select half of the countries to change their capitals to the non-293 capital cities. Then, we randomly sample 1,000 pairs, each with one country having an unchanged 294 capital and one with a changed capital. Finally, we generate fill-in-the-blank-style multiple-choice 295 questions based on these pairs, such as "Fill in the blank: Athens is to Greece like Paris is to ___ 296 (A) Tonga (B) France." For this example, both the r_{cityof} and $r_{capitalof}$ relations provide sufficient 297 reasoning to answer as "France". While the corresponding synthetic explanation, $\varepsilon_{conitalof}$, for the 298 model with unaltered capitals would be "The capital of France is Paris, as the capital of Greece is Athens.", the one for the model with altered capitals, ε_{citv0f} , would be "Paris is a city in France, as 299 Athens is a city in Greece." 300

301 302

303

4.3 OBJECT COUNTING TASKS

Task Inspired by the object_counting task from BIG-bench (bench authors, 2023), we adapt an object counting task for evaluating diagnosticity. The task involves counting entities of a given type from a list of entities. By modifying model knowledge to swap objects across predefined categories, we ensure the number of entities of the target type remains the same while changing the reasoning behind the answer. For example, when asked how many of "countertop," "grape," and "kiwifruit" are fruits, the answer is 2, since "countertop" is a furniture item. If we edit the model to classify "countertop" as a fruit and "grape" as furniture, the answer remains 2 but due to different reasoning.

310 311

Dataset We define five categories with five types each, as shown in Table 2 in Appendix A . For each type, we select 10 representative entities from WikiData. We then reserve 20% of the entities for reassignment to other types within the same category after model editing. We include two question types: yes/no questions, asking if all or any items in a list belong to a given type, and number questions, asking how many items belong to a specific type.

For both types, we randomly determine the number of items k (between 3 and 6) and select a target type. For yes/no questions, we sample k entities, ensuring that after model editing, the number of entities of the target type remains unchanged. For number questions, we reassign one entity from the target type and one from other types to ensure consistency.

We generate 1,000 samples in total, equally divided between the two question types. Further details about the dataset generation process are included in Appendix A.

323

¹https://www.kaggle.com/datasets/viswanathanc/world-cities-datasets/

		Metric	Mistral-Inst	ruct LLaMa-2-7t	o-chat LLaMa-2-7b	GPT-J 6B
Check	0C	CC-SHAP	0.437	0.518	0.665	0.553
	sth	Simulatability	0.014	0.052	0.035	0.033
	$\mathbf{P}_{\mathbf{O}}$	Counterfact. Edits	0.001	0.000	0.000	0.000
		Early Answering	0.030	0.033	0.045	0.056
act	_	Filler Tokens	0.019	0.029	0.025	0.022
Ë	5	Adding Mistakes	0.013	0.047	0.158	0.029
	0	Paraphrasing	0.160	0.108	0.171	0.029
		CC-SHAP	0.559	0.522	0.616	0.547
	00	CC-SHAP	0.850	0.583	0.657	0.355
	sth	Simulatability	0.006	0.001	0.000	0.000
N	$\mathbf{P}_{\mathbf{O}}$	Counterfact. Edits	0.001	0.000	0.000	0.000
alog		Early Answering	0.041	0.018	0.110	0.063
An	<u>_</u>	Filler Tokens	0.041	0.011	0.044	0.145
7	5	Adding Mistakes	0.118	0.023	0.190	0.198
	0	Paraphrasing	0.123	0.121	0.165	0.235
		CC-SHAP	0.859	0.663	0.672	0.411
	100	CC-SHAP	0.522	0.460	0.510	0.500
ing	sth	Simulatability	0.031	0.028	0.037	0.034
unt	\mathbf{P}_{0}	Counterfact. Edits	0.000	0.000	0.000	0.000
bject Cou		Early Answering	0.109	0.005	0.086	0.120
	_	Filler Tokens	0.065	0.033	0.058	0.074
	5	Adding Mistakes	0.124	0.129	0.109	0.164
0	U	Paraphrasing	0.191	0.173	0.154	0.190
		CC-SHAP	0.504	0.467	0.494	0.509

Table 1: The diagnosticity scores of each model for each faithfulness metric across three tasks, along with the accuracy of each model on each task under standard and CoT prompting. Bold numbers indicate the highest scores for each model on each task across the two categories of faithfulness metrics: post-hoc and CoT. "Mistral-Instruct" refers to the mistral-7b-instruct-v0.2 model.

5 EXPERIMENTS

We present four sets of experiments. First, we report the diagnosticity scores of post-hoc and CoTbased metrics across three tasks and four LLMs. Second, we conduct an analysis to assess the reliability of the model edits used for diagnosticity evaluation. Third, we perform an ablation study where we replace MEMIT with a simplified version of IKE, examining how the choice of model editing method affects our results. Finally, we conduct another ablation study in which we use model-generated explanations instead of synthetically generated ones.

5.1 DIAGNOSTICITY EVALUATION OF FAITHFULNESS METRICS

367 Experimental Setup We evaluate the seven metrics described in Section 2 across four different
 368 LLMs: mistral-instruct-7b-v0.2 (Jiang et al., 2023), 11ama-2-7b, 11ama-2-7b-chat (Touvron
 369 et al., 2023), and gpt-j-6B (Wang & Komatsuzaki, 2021). For our main experiments, we employ
 370 MEMIT as the model editing method and use synthetic explanations to ensure their faithfulness with
 371 respect to the edited model.

372Table 1 presents the diagnosticity scores for all faithfulness metrics across three tasks for the373four models. The most notable finding is that CC-SHAP significantly outperforms other methods374(McNemar's test, p < .01) in each task, for each model, across both post-hoc and CoT-based metrics.375In the post-hoc category, Simulatability shows significantly higher diagnosticity than Counterfactual376Edits across all models for the Object Counting and FactCheck tasks (McNemar's test, p < .01),377and higher or comparable diagnosticity for the Analogy task. In the Analogy task, Paraphrasingand Adding Mistakes significantly outperform other CoT-based metrics (McNemar's test, p < .01),

following CC-SHAP, across all models with the exception of Adding Mistakes in 11ama-2-7b-chat. For the Object Counting task, Paraphrasing becomes the second-best CoT-based metric, significantly outperforming other metrics for all models (McNemar's test, p < .01) except gpt-j-6b.

381 Although CC-SHAP outperforms Paraphrasing by a wide 382 margin, Paraphrasing consistently ranks as the secondhighest diagnosticity metric in most cases, followed by 384 Adding Mistakes, Early Answering, and Filler Tokens, 385 respectively. However, there are some exceptions to this 386 order. For instance, Early Answering ranks as the second-387 best metric in the FactCheck task for gpt-j-6b, while 388 Adding Mistakes ranks second in the Analogy task for 11ama-2-7b. Although this ranking generally holds, the 389 relative differences are not always statistically significant. 390

391 When examining cases where faithfulness metrics fail to 392 correctly assign higher scores to faithful explanations, we 393 find that binary metrics (all except CC-SHAP) often strug-394 gle to differentiate between the faithfulness levels of ex-395 planations, frequently assigning the same score to both. Across all three tasks, most binary metrics fail in this man-396 ner at least 90% of the time. However, some metrics more 397 frequently assign lower scores to faithful explanations 398 than to unfaithful ones. For example, Paraphrasing as-399 signs lower scores to faithful explanations at least 15% of 400 the time across all tasks, while Adding Mistakes and Early 401 Answering do so at least 15% of the time for the Object 402 Counting task. A closer look at *Paraphrasing* examples 403 reveals that the paraphrasing process can lead to signifi-404 cant hallucinations, sometimes even causing paraphrases 405 of contradictory explanation pairs to state the same facts.

406 These findings highlight the importance of carefully 407 selecting a helper model when using faithfulness met-408 rics based on corrupting CoT. Following Parcalabescu 409 & Frank (2023), we use llama-2-13b-chat as our helper 410 model. Similarly, Lanham et al. (2023) use the same 411 model as their predictor and explainer: a 175B-parameter 412 decoder-only transformer LLM (Vaswani et al., 2017; Radford & Narasimhan, 2018; Radford et al., 2019; Brown 413

Average Absolute Difference of Faithfulness Scores for CC-SHAP



Figure 3: The absolute average difference in faithfulness scores assigned to pair of explanations by CC-SHAP for cases where CC-SHAP *correctly* assigns higher scores to faithful explanations and where it *incorrectly* assigns higher scores to unfaithful ones, across all tasks, for both posthoc and CoT-based CC-SHAP, using mistral-7b-instruct-v0.2

et al., 2020). While these issues may be less apparent with larger models, practitioners should be cautious when using a helper model of similar size to the model being tested, particularly for smaller models.

417 Since CC-SHAP is a smoother metric, we find no instances where it fails by assigning the same score to both explanations. To gain deeper insight, we examine the average absolute differences in 418 faithfulness scores between each pair of explanations. Figure 3 presents these differences for cases 419 where CC-SHAP correctly assigns higher scores to faithful explanations and where it incorrectly 420 assigns higher scores to unfaithful ones, across all tasks, for both post-hoc and CoT-based CC-SHAP, 421 using mistral-7b-instruct-v0.2. The results indicate that the average absolute differences in 422 faithfulness scores are generally similar for both correct and incorrect cases. However, in the Analogy 423 task, CC-SHAP better distinguishes between faithful and unfaithful explanations when it performs 424 correctly compared to when it fails. While this observation aligns with the task in which CC-SHAP 425 achieves its highest diagnosticity scores, no significant correlation is found between diagnosticity and 426 the average absolute difference in faithfulness scores.

427 428

429 430

5.2 RELIABILITY OF EDITS

431 CAUSAL DIAGNOSTICITY relies on the assumption that, in each given explanation pair, one explanation is faithful to the model being evaluated while the other is unfaithful. To ensure this condition is

432 met, we modify the models and use synthetically generated explanation pairs. While these synthetic 433 explanations logically guarantee faithfulness or unfaithfulness with respect to the edited model, their 434 practical accuracy depends on the success of the editing method. One way to assess whether the 435 synthetic explanations align with faithfulness expectations is by comparing the perplexities of the 436 explanation pairs. Since the only difference between the explanations is related to the aspect modified by the model edit, the intuition is that the explanation deemed faithful should have a lower perplexity 437 than the one deemed unfaithful. 438

439 Figure 4 shows the frequency with which expla-440 nations deemed as faithful have lower perplexity 441 than those deemed as unfaithful, for each task 442 and each model. While the explanations deemed faithful generally have lower perplexities than 443 their unfaithful counterparts across all tasks, the 444 edits performed for the Fact Check task are par-445 ticularly successful, with scores nearing 1.0. In 446 contrast, the edits for the Analogy and Object 447 Counting tasks perform relatively worse. The 448 scores for these two tasks are similar across all 449 models, except for gpt-j-6b, where the edits 450 for the Analogy task perform notably worse. 451

452 5.3 EFFECT

454

453 OF KNOWLEDGE EDITING METHOD

455 We investigate the effect of different model editing methods on our results by conducting an 456 ablation study where we replace MEMIT with 457 an alternative approach. Instead of selecting 458 another locate-then-edit method, we use a sim-459



1.0

this method are provided in Appendix B. 461

Figure 5 compares MEMIT and IKE across all faith-462 fulness metrics, with diagnosticity scores averaged 463 over three tasks. While the diagnosticity scores from 464 models edited with MEMIT are higher than those 465 obtained with IKE, the relative relationships between 466 different metrics remain consistent. This suggests 467 that the choice of model editing method has no sig-468 nificant impact on our conclusions. 469

470

471

479

5.4 EFFECT OF MODEL GENERATED EXPLANATIONS 472

473 While our main results are derived from using syn-474 thetically generated explanations to form faithful and unfaithful explanation pairs that accurately reflect the 475 applied edits and the differences between the mod-476 els, we also perform an ablation study using model-477 generated explanations. 478





as faithful have lower perplexity than those deemed as unfaithful, for each task and each model. Higher frequency indicates the higher success in applied edits.

> IKE 0.8 0.6 ē 0.4 0.2 CF. Edite جري post-hoc CoT Faithfulness Metric

MEMI

Figure 5: Diagnosticity scores for each metric on mistral-7B-instruct-v0.2 using two model editing methods: MEMIT and IKE. Although the scores are higher when MEMIT is used, the ranking of the metrics remains consistent across both editing methods.

Experimental Setup We evaluate all faithfulness 480

metrics using mistral-7B-instruct-v0.2. For model-generated explanations, the length is limited 481 to 100 tokens. 482

483 Figure 6 compares model-generated and synthetic explanations across all faithfulness metrics, with diagnosticity scores averaged over three tasks. Although CC-SHAP consistently outperforms both 484 other post-hoc and CoT-based metrics for both explanation types, there is no consistent difference 485 in the diagnosticity scores between the two explanation types across all metrics. Furthermore, the

486 comparative ranking of faithfulness metrics is inconsistent when replacing synthetic explanations 487 with model-generated ones. Upon examining the model-generated explanations, we observe several 488 issues. At times, explanations pairs contain hallucinations, making them unfaithful to their own 489 models and violating the main condition of our framework. Occasionally, explanations are truncated 490 due to the token limit. In some cases, an explanation may begin appropriately but revert to pre-edit knowledge. Particularly in the Analogy and Object Counting tasks, models often fail to articulate the 491 applied edits. While these issues could be attributed to the limited generalizability of model editing 492 methods, larger models or memory-based editing approaches may help address these challenges 493 (Yao et al., 2023). Nevertheless, synthetically generated explanations stand out as a viable option, 494 especially when considering the computational costs associated with these alternatives. 495

- 496
- 497 498

499

6 CONCLUSION

500 In this paper, we introduce a new framework, 501 CAUSAL DIAGNOSTICITY, to evaluate faith-502 fulness metrics for natural language explana-503 tions by extending the notion of diagnosticity. 504 We introduce three new tasks-fact-checking, 505 analogy, and object counting-while utilizing 506 model editing to generate pairs of faithful and 507 unfaithful explanations to measure diagnosticity. 508 We benchmark popular post-hoc and CoT-based 509 faithfulness metrics across these tasks. The re-510 sults show that most metrics fail to achieve satisfactory diagnosticity scores, with CC-SHAP 511 being a notable exception. Unlike other meth-512 ods, CC-SHAP leverages more information by 513 considering token-wise interactions between the 514 explanation and the input, which likely allows it 515 to capture the inner workings of the model bet-516 ter than methods that simply observe changes in 517 output after perturbing the input or explanations. 518

Despite CC-SHAP's higher scores, the results





also highlight areas for improvement, particularly in terms of the computational cost and slowness
of CC-SHAP. Based on these findings, developing metrics that focus more on the model's internal
mechanisms and complex interactions among explanations, inputs, and outputs could be a promising
direction.

We view this study as the first step in the quest for more faithful LLM explanations by providing a test bed for faithfulness metrics. As our study reveals the inadequacy of existing metrics and underscores the need for better alternatives, a natural direction for future research is the development of improved faithfulness metrics, which should then be followed by the creation of more faithful explanation methods.

- 528
- 529 530

7 LIMITATIONS

531 532

This study is limited to 7B-parameter models due to the availability of models with published hyperparameters for MEMIT editing and computational constraints. Additionally, CAUSAL DIAGNOSTICITY is heavily relies on the effectiveness of the model editing method. While we conduct an ablation study using the IKE baseline, the utility of model-generated explanations remains largely unexplored. This is because approaches to address issues in model-generated explanations, such as employing memory-based knowledge editing methods or using larger models, come with high computational costs. In particular, memory-based methods lead to lengthy experiments with CC-SHAP due to the increased context length.

540	REFERENCES
541	REFERENCEDS

566

567

568

569

581

582

583

584

591

Pepa Atanasova, Oana-Maria Camburu, Christina Lioma, Thomas Lukasiewicz, Jakob Grue Si monsen, and Isabelle Augenstein. Faithfulness tests for natural language explanations. *ArXiv*, abs/2305.18029, 2023. URL https://api.semanticscholar.org/CorpusID: 258960511.

- BIG bench authors. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL https://openreview.net/forum?id=uyTL5Bvosj.
- 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, Jeff Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Ma teusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. *ArXiv*, abs/2005.14165, 2020. URL https://api.semanticscholar.org/CorpusID:218971783.
- Aaron Chan, Shaoliang Nie, Liang Tan, Xiaochang Peng, Hamed Firooz, Maziar Sanjabi, and Xiang
 Ren. Frame: Evaluating rationale-label consistency metrics for free-text rationales. 2022a. URL
 https://api.semanticscholar.org/CorpusID:254247321.
- Chun Sik Chan, Huanqi Kong, and Guanqing Liang. A comparative study of faithfulness metrics for model interpretability methods. In *Annual Meeting of the Association for Computational Linguistics*, 2022b. URL https://api.semanticscholar.org/CorpusID:248118978.
- Roi Cohen, Eden Biran, Ori Yoran, Amir Globerson, and Mor Geva. Evaluating the ripple effects
 of knowledge editing in language models. *Transactions of the Association for Computational Linguistics*, 12:283–298, 2024.
 - Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine learning. arXiv: Machine Learning, 2017. URL https://api.semanticscholar.org/CorpusID: 11319376.
- Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. Dissecting recall of factual associations in auto-regressive language models. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 12216–12235, Singapore, December 2023. Association for Computational Linguistics. doi: 10. 18653/v1/2023.emnlp-main.751. URL https://aclanthology.org/2023.emnlp-main.751.
- Anshita Gupta, Debanjan Mondal, Akshay Sheshadri, Wenlong Zhao, Xiang Li, Sarah Wiegreffe, and Niket Tandon. Editing common sense in transformers. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 8214–8232, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.511. URL https://aclanthology.org/2023.emnlp-main.511.
 - Thomas Hartvigsen, Swami Sankaranarayanan, Hamid Palangi, Yoon Kim, and Marzyeh Ghassemi. Aging with grace: Lifelong model editing with discrete key-value adaptors. In *Advances in Neural Information Processing Systems*, 2023.
- Peter Hase and Mohit Bansal. Evaluating explainable ai: Which algorithmic explanations help users
 predict model behavior? In Annual Meeting of the Association for Computational Linguistics,
 2020. URL https://api.semanticscholar.org/CorpusID:218502350.
- Peter Hase, Shiyue Zhang, Harry Xie, and Mohit Bansal. Leakage-adjusted simulatability: Can models generate non-trivial explanations of their behavior in natural language? In *Findings*, 2020. URL https://api.semanticscholar.org/CorpusID:222209056.
- 592 Peter Hase, Mohit Bansal, Been Kim, and Asma Ghandeharioun. Does localization inform editing? surprising differences in causality-based localization vs. knowledge editing in language models. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.),

594 Advances in Neural Information Processing Systems, volume 36, pp. 17643–17668. Curran Asso-595 ciates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/ 596 3927bbdcf0e8d1fa8aa23c26f358a281-Paper-Conference.pdf. 597 Alon Jacovi and Yoav Goldberg. Towards faithfully interpretable nlp systems: How should we define 598 and evaluate faithfulness? In Annual Meeting of the Association for Computational Linguistics, 2020. URL https://api.semanticscholar.org/CorpusID:215416110. 600 601 Albert Qiaochu Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh 602 Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile 603 Saulnier, L'elio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b. ArXiv, abs/2310.06825, 604 2023. URL https://api.semanticscholar.org/CorpusID:263830494. 605 606 Tamera Lanham, Anna Chen, Ansh Radhakrishnan, Benoit Steiner, Carson E. Denison, Danny Hernan-607 dez, Dustin Li, Esin Durmus, Evan Hubinger, John Kernion, Kamil.e Lukovsiut.e, Karina Nguyen, 608 Newton Cheng, Nicholas Joseph, Nicholas Schiefer, Oliver Rausch, Robin Larson, Samuel McCan-609 dlish, Sandipan Kundu, Saurav Kadavath, Shannon Yang, Tom Henighan, Timothy D. Maxwell, 610 Timothy Telleen-Lawton, Tristan Hume, Zac Hatfield-Dodds, Jared Kaplan, Janina Brauner, 611 Sam Bowman, and Ethan Perez. Measuring faithfulness in chain-of-thought reasoning. ArXiv, abs/2307.13702, 2023. URL https://api.semanticscholar.org/CorpusID:259953372. 612 613 Scott M. Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In Neural 614 Information Processing Systems, 2017. URL https://api.semanticscholar.org/CorpusID: 615 21889700. 616 Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing fac-617 tual associations in gpt. In Neural Information Processing Systems, 2022. URL https: 618 //api.semanticscholar.org/CorpusID:255825985. 619 620 Kevin Meng, Arnab Sen Sharma, Alex Andonian, Yonatan Belinkov, and David Bau. Mass editing 621 memory in a transformer. The Eleventh International Conference on Learning Representations 622 (ICLR), 2023. 623 Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D Manning. Fast 624 model editing at scale. In International Conference on Learning Representations, 2022. URL 625 https://openreview.net/pdf?id=0DcZxeWfOPt. 626 627 Letitia Parcalabescu and Anette Frank. On measuring faithfulness or self-consistency of natural language explanations. 2023. URL https://api.semanticscholar.org/CorpusID:265150102. 628 629 Vaidehi Patil, Peter Hase, and Mohit Bansal. Can sensitive information be deleted from llms? 630 objectives for defending against extraction attacks. ArXiv, abs/2309.17410, 2023. URL https: 631 //api.semanticscholar.org/CorpusID:263311025. 632 Alec Radford and Karthik Narasimhan. Improving language understanding by generative pre-training. 633 2018. URL https://api.semanticscholar.org/CorpusID:49313245. 634 635 Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language 636 models are unsupervised multitask learners. 2019. URL https://api.semanticscholar.org/ 637 CorpusID: 160025533. 638 Chenmien Tan, Ge Zhang, and Jie Fu. Massive editing for large language models via meta learning. 639 In International Conference on Learning Representations, 2024. URL https://openreview. 640 net/pdf?id=L6L1CJQ2PE. 641 642 Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, 643 Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas 644 Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S. 645 Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, 646

646 Fartshoff, Saghar Hossenfi, Kui Hou, Hakan man, Marcin Kardas, Viktor Kerkez, Madian Khaosa,
 647 Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril,
 648 Jeneya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar

648 649 650 651 652 653	Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models. <i>ArXiv</i> , abs/2307.09288, 2023. URL https://api.semanticscholar.org/CorpusID:259950998
654 655 656	Miles Turpin, Julian Michael, Ethan Perez, and Sam Bowman. Language models don't always say what they think: Unfaithful explanations in chain-of-thought prompting. <i>ArXiv</i> , abs/2305.04388, 2023. URL https://api.semanticscholar.org/CorpusID:258556812.
658 659 660	Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In <i>Neural Information Processing Systems</i> , 2017. URL https://api.semanticscholar.org/CorpusID:13756489.
661 662	Ben Wang and Aran Komatsuzaki. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. https://github.com/kingoflolz/mesh-transformer-jax, May 2021.
663 664 665	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Huai hsin Chi, F. Xia, Quoc Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. <i>ArXiv</i> , abs/2201.11903, 2022. URL https://api.semanticscholar.org/CorpusID:246411621.
666 667 668	Sarah Wiegreffe, Ana Marasović, and Noah A. Smith. Measuring association between labels and free-text rationales. In <i>Conference on Empirical Methods in Natural Language Processing</i> , 2020. URL https://api.semanticscholar.org/CorpusID:225068329.
670 671 672 673 674	Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng, Zhoubo Li, Shumin Deng, Huajun Chen, and Ningyu Zhang. Editing large language models: Problems, methods, and opportunities. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pp. 10222–10240, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.632. URL https://aclanthology.org/2023.emnlp-main.632.
675 676 677 678	Lang Yu, Qin Chen, Jie Zhou, and Liang He. Melo: Enhancing model editing with neuron- indexed dynamic lora. <i>ArXiv</i> , abs/2312.11795, 2023. URL https://api.semanticscholar. org/CorpusID:266362196.
679 680 681 682 683	Ningyu Zhang, Yunzhi Yao, Bo Tian, Peng Wang, Shumin Deng, Meng Wang, Zekun Xi, Shengyu Mao, Jintian Zhang, Yuansheng Ni, Siyuan Cheng, Ziwen Xu, Xin Xu, Jia-Chen Gu, Yong Jiang, Pengjun Xie, Fei Huang, Lei Liang, Zhiqiang Zhang, Xiao-Jun Zhu, Jun Zhou, and Huajun Chen. A comprehensive study of knowledge editing for large language models. <i>ArXiv</i> , abs/2401.01286, 2024. URL https://api.semanticscholar.org/CorpusID:266725300.
684 685 686 687 688	Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong Wu, Jingjing Xu, and Baobao Chang. Can we edit factual knowledge by in-context learning? In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pp. 4862–4876, Singapore, December 2023. Association for Computational Linguistics. doi: 10. 18653/v1/2023.emnlp-main.296. URL https://aclanthology.org/2023.emnlp-main.296.
689 690 691	A DATASET
692 693 694 695	Figure 7 illustrates the prompt used to convert statements from COUNTERFACT into yes/no questions for the Fact Check task, utilizing Mistral-7B-Instruct-v0.2. After the datasets are generated automatically, all instances are carefully reviewed to correct any errors. Table 2 presents the categories and types used in the Object Counting task.

- 696 697
- **B** MODEL EDITING
- 698699 B.1 TASK-BASED EDITING TEMPLATES
- Table 3 shows the templates we use for editing models in each task. For the FactCheck task, there is a variety of prompts where the action or situation of the subject differs, but the target is always located

Category	Types				
object occupation company	xtanimal, musical instrument, fruit, vegetable, furniturepationscientist, politician, soccer player, actor, singerpanymedia company, energy company, software company, automotive company,				
consulting companyouristic attractionabstractFrance, Spain, Russia, Turkey, Italyreligion, political ideology, language, branch of science, emotion					
Table 2: Categories and corresponding types used in Object Counting task					
Please create a yes-no question from the given sentence. Here are some examples: Sentence: Joe Biden is the president of the United States. Question: Is Joe Biden the → president of the United States? Sentence: They play rock. Question: Do they play rock? Sentence: Quesadilla from Mexico. Question: Is quesadilla from Mexico? Do not mention your assumptions or assesment towards correctness of question. Do not output → anything else! Stick with the format. Sentence: {SENTENCE} Question:					
F	igure 7: The prompt used for converting	statements to questions.			
FactCheck					
Template #1: ···	· subject ··· target				
The profession o	f Heath Ledger is king	The profession of Heath Ledger is webmas			
Darryl Jones	plays blues	Darryl Jones plays heavy metal			
	Analo	gy			
Template #1: TI	ne capital of country is city				
T1					
The capital of U	Jnited Kingdom is London	United Kingdom is Birmingham			
The capital of U	Inited KingdomisLondonJnited Arab EmiratesisAbu Dhabi	United Kingdom is Birmingham The capital of United Arab Emirates is Du			
The capital of U The capital of U Template #2:	Jnited Kingdom is LondonJnited Arab Emirates is Abu Dhabicapital is a city in country	United Kingdom is Birmingham The capital of United Arab Emirates is Du			
The capital of Template #2:	Jnited KingdomisLondonJnited Arab EmiratesisAbu Dhabicapitalis a city incountryity inUnited Kingdom	United Kingdom is Birmingham The capital of United Arab Emirates is Du			
The capital of The capital of Template #2: Condon is a cited by the capital of th	United Kingdom is London United Arab Emirates is Abu Dhabi capital is a city in country ity in United Kingdom a city in United Arab Emirates	United Kingdom is Birmingham The capital of United Arab Emirates is Du			
The capital of U The capital of U Template #2: C London is a ci Abu Dhabi is	Jnited Kingdom is London Jnited Arab Emirates is Abu Dhabi capital is a city in country ity in United Kingdom a city in United Arab Emirates Object Co	United Kingdom is Birmingham The capital of United Arab Emirates is Du			
The capital of The capital of Template #2: Control of the capital of Template #2: Control of the capital of the	United Kingdom is London United Arab Emirates is Abu Dhabi capital is a city in country ity in United Kingdom a city in United Arab Emirates Object Co entity is/is located in type	United Kingdom is Birmingham The capital of United Arab Emirates is Du			
The capital of U The capital of U Template #2: 0 London is a ci Abu Dhabi is Template #1: 0 dog is anima	United Kingdom is London United Arab Emirates is Abu Dhabi capital is a city in country ity in United Kingdom a city in United Arab Emirates Object Co entity is/is located in type 1	United Kingdom is Birmingham The capital of United Arab Emirates is Du unting			
The capital of U The capital of U Template #2: 0 London is a ci Abu Dhabi is Template #1: 6 dog is anima	United Kingdom is London United Arab Emirates is Abu Dhabi capital is a city in country ity in United Kingdom a city in United Arab Emirates Object Co entity is/is located in type 1 ter is located in Turkey	United Kingdom is Birmingham The capital of United Arab Emirates is Du unting dog is musical instrument			

Table 3: Templates used for editing models. Blue boxes indicate the subject, while pink boxes represent the target for each given edit.

748

749 750

at the end of the prompt. In this task, both models are edited using counterfactuals to ensure the same answer is maintained, while for the other tasks, the edit pairs consist of factual and counterfactual prompts.

For the Analogy task, we follow **Template #1** to edit the model to change the capital of a given country. Even for the model where the capitals remain unchanged, we apply this edit in case the

⁷⁵⁶ model lacks knowledge of some countries. For both models, we reinforce the r_{cityOf} relation by applying **Template #2**.

For the Object Counting task, we use the corresponding template in Table 3 to edit the model by altering the types of entities. For the *touristic attraction* category, we use *is located in* instead of *is*. Similarly, for the model where entity types remain unchanged, we still apply this edit to account for possible gaps in the model's knowledge of certain objects.

764 B.2 IN-CONTEXT KNOWLEDGE EDITING (IKE) BASELINE

Zheng et al. (2023) leverage In-Context Learning for knowledge editing in LLMs without requiring parameter updates. They define three types of in-context demonstrations to enhance generalization (the ability to update knowledge expressed in different textual forms) and specificity (the ability to avoid altering unrelated knowledge when making edits). These demonstrations are: (1) copy for injecting new facts, (2) update for improving generalization, and (3) retain for preventing changes to unrelated knowledge. However, for our IKE experiments, we adopt their simpler PROMPT baseline, where new facts are directly added to the context. We use the same templates shown in Table 3, but prepend each relevant edit just before the query. When measuring faithfulness scores, we exclude the prefix containing these edits from any operations and keep it fixed.

C ADDITIONAL RESULTS

C.1 IKE RESULTS

	Metric	FactCheck	Analogy	Object Counting
100	CC-SHAP	0.418	0.938	0.287
Postl	Simulatability	0.014	0.008	0.032
	Counterfact. Edits	0.000	0.000	0.000
	Early Answering	0.016	0.003	0.057
Γ.	Filler Tokens	0.011	0.001	0.075
Ę	Adding Mistakes	0.027	0.085	0.104
0	Paraphrasing	0.105	0.046	0.167
	CC-SHAP	0.460	0.963	0.279

Table 4: The diagnosticity scores of mistral-7b-instruct-v0.2 for each faithfulness metric across three tasks, along with the accuracy of each model on each task under standard and CoT prompting when IKE baseline is used as model editing method. Bold numbers indicate the highest scores for each model on each task across post-hoc and CoT-based faithfulness metrics.