

# Evaluating Human Alignment and Model Faithfulness of LLM Rationale

Anonymous ACL submission

## Abstract

We study how well large language models (LLMs) explain their generations with rationales – a set of tokens extracted from the input texts that reflect the decision process of LLMs. We examine LLM rationales extracted with two methods: 1) attribution-based methods that use attention or gradients to locate important tokens, and 2) prompting-based methods that guide LLMs to extract rationales using prompts. Through extensive experiments, we show that prompting-based rationales align better with human-annotated rationales than attribution-based rationales, and demonstrate reasonable alignment with humans even when model performance is poor. We additionally find that the faithfulness limitations of prompting-based methods, which are identified in previous work, may be linked to their collapsed predictions. By fine-tuning these models on the corresponding datasets, both prompting and attribution methods demonstrate improved faithfulness. Our study sheds light on more rigorous and fair evaluations of LLM rationales, especially for prompting-based ones.<sup>1</sup>

## 1 Introduction

The rise of large language models (LLMs) has significantly transformed the field of natural language processing (NLP) (Touvron et al., 2023; Team et al., 2023; OpenAI et al., 2024), enabling a wide range of applications from web question answering to complex reasoning tasks. However, they are not always reliable and usually cannot clearly explain their outputs (Ji et al., 2023), which limits the deployment of these models in high-stakes scenarios.

*Rationales*<sup>2</sup>, i.e., tokens of the input text that are most influential to the models’ predictions, are widely studied in the NLP community prior

<sup>1</sup>Code and data will be released upon paper acceptance.

<sup>2</sup>Also called *self-explanation* or *extractive rationales* in previous work (Huang et al., 2023a; Madsen et al., 2024).

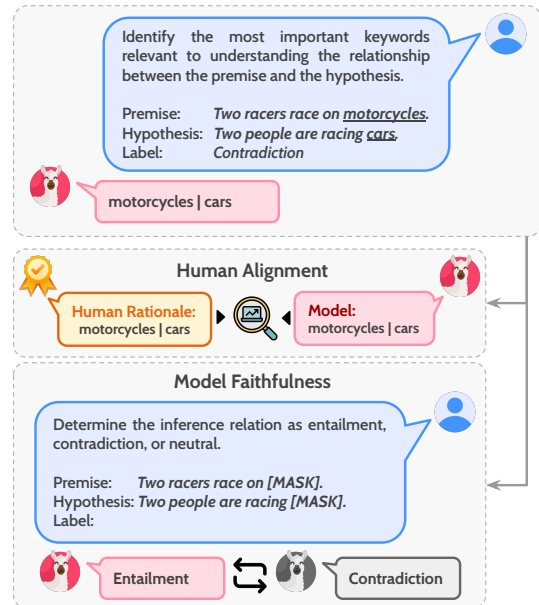


Figure 1: An example of our analysis methodology on the e-SNLI dataset. *Human alignment* compares model rationales with human-annotated rationale; *Model faithfulness* measures the rates when model prediction changes (e.g. from Contradiction to Entailment) after masking the rationales identified by the model.<sup>3</sup>

to the era of LLMs to interpret model predictions (Lei et al., 2016; DeYoung et al., 2019; Wiegrefe and Pinter, 2019; Jacovi and Goldberg, 2020). For smaller and open-source models like BERT (Devlin et al., 2019), rationales are extracted with attribution-based methods like attention weights (Wiegrefe and Pinter, 2019) or gradients (Li et al., 2016). For decoder-only LLMs, besides attribution-based methods, rationales can also be extracted by leveraging the instruction-following ability of LLMs and guiding them with explicit prompts to explain their predictions (Figure 1). We call these prompting-based rationales.

To evaluate different rationales, previous works

<sup>3</sup>In the first prompt, we use the true label for human alignment, and the predicted label for faithfulness experiments.

on model interpretation establish two properties of rationales that are critical for successful interpretability: human alignment (DeYoung et al., 2019; Hase and Bansal, 2022) and faithfulness (Jacovi and Goldberg, 2020). *Human alignment* refers to the degree to which the rationales match or align with human-annotated rationales, while *faithfulness* assesses whether the rationales truly reflect the model’s internal process. A longstanding debate exists regarding the relationship between these two aspects (Agarwal et al., 2024). However, studies on LLM rationales either focus on the faithfulness of off-the-shelf LLMs (Huang et al., 2023a; Madsen et al., 2024), or their human alignment (Chen et al., 2023), but lack a comprehensive exploration of the two properties together. Specifically, recent works (Huang et al., 2023a; Madsen et al., 2024) study prompting-based methods and show that they might not be faithful to the reasoning process of LLMs. Moreover, they only consider LLMs as out-of-box models, without fine-tuning on specific tasks. How fine-tuning of LLMs on downstream tasks influences the human alignment and faithfulness of LLM rationales is under-explored.

In this paper, we conduct extensive experiments to evaluate LLM rationales more comprehensively and bridge the gap in existing research. We consider five state-of-the-art LLMs, encompassing both open-source models (Llama2 (Touvron et al., 2023), Llama3, Mistral (Jiang et al., 2023)) and proprietary models (GPT-3.5-Turbo, GPT-4-Turbo (OpenAI et al., 2024)). Our study leverages two annotated natural language classification datasets, e-SNLI (Camburu et al., 2018a) and MedicalBios (Eberle et al., 2023), to evaluate and compare rationale extraction methods based on prompting strategies and feature attribution-based techniques such as Input×Gradient (Li et al., 2016).

Through our experiments, we find that while prompting-based rationales are generally less faithful than attribution-based methods, they tend to align better with human-annotated rationales, both before and after fine-tuning. Surprisingly, even when prompting-based rationales exhibit poor performance, they can still produce explanations that reasonably align with human reasoning.

We also observe that low classification performance and collapsing predictions might be related to the faithfulness limitation of LLM rationales. Fine-tuning LLMs on specific datasets improves the quality of rationales for both prompting and, particularly, attribution techniques in terms of faith-

fulness and human alignment. This finding complements the observations in Madsen et al. (2024), where faithful evaluation was conducted using only out-of-the-box LLMs.

In summary, our work contributes to the ongoing efforts to enhance the interpretability and trustworthiness of LLMs by providing empirical evidence and practical recommendations for extracting and evaluating rationales from these models.

## 2 Related Work

**Interpretability** Recent literature in natural language processing (NLP) has seen a surge in interpretability methods aimed at making models more transparent and understandable. The traditional interpretability methods include 1) attribution-based methods, which leverage the attention weights in models like transformers to identify which parts of the input the model focuses on when making a decision (Vaswani et al., 2023; Clark et al., 2019; Abnar and Zuidema, 2020), 2) Gradient-based methods, which provide explanations by identifying which input tokens most influence the model’s output, often using techniques like gradient-based saliency maps (Simonyan et al., 2014a), or its extension by incorporating the input vector norms or integration (Sundararajan et al., 2017). 3) Vector-based methods that propagate the decomposed representations throughout the model achieving the best faithfulness results on encoder-based models (Kobayashi et al., 2020, 2021; Ferrando et al., 2022; Modarressi et al., 2022, 2023). More recently, researchers have been using rationales, written in natural language, to serve interpretability where they can reveal the “reasoning” behind model decisions.

**Rationales** Rationales can be categorized as free-form or extractive. Free-form rationales use natural language to explain the model’s reasoning, filling in commonsense knowledge gaps. They can improve model performance (Sun et al., 2022a) and user interpretability (Sun et al., 2022b). Extractive rationales highlight specific parts of the input text that provide sufficient evidence for a prediction, independent of the rest of the input (Lei et al., 2016; DeYoung et al., 2020). They can also enhance model performance (Huang et al., 2021; Carton et al., 2021) and improve human interpretability (Strout et al., 2019). Our work focuses on extractive rationales for interpretability evaluation. In this research area, Huang et al. (2023b) studied faithfulness in ChatGPT, comparing prompting and

	Input	Human	Mistral-7B	LLama-3-8B	GPT-4 Turbo
e-SNLI	Premise: <b>Five</b> children playing soccer chase after a ball. Hypothesis: There are <b>ten</b> children playing. Label: Contradiction	[ <b>Five, ten</b> ]	[ <b>Five</b> children, <b>ten</b> children]	[children, playing]	[children, <b>ten</b> ]
	Premise: Asian individuals are <b>taking pictures</b> of some unidentified thing. Hypothesis: Asians <b>play soccer</b> in the park. Label: Contradiction	[ <b>taking, pictures, play, soccer</b> ]	[Asian, individuals, <b>taking, pictures</b> , unidentified, thing, <b>soccer</b> , park]	[Asian, individuals, unidentified, <b>soccer</b> ]	[Asian, <b>pictures, play, soccer</b> ]
MedicalBios	His basic <b>surgical training</b> was undertaken in Newcastle. After gaining his <b>MRCS</b> and spending some time carrying out locum SpR work at Hurstwood Park, the Royal London Hospital Label: Surgeon	[ <b>surgical, training, MRCS</b> ]	[His, <b>MRCS, surgical</b> ]	[ <b>surgical, training, MRCS</b> ]	[ <b>surgical, MRCS, SpR</b> ]
	She works at both the Park Nicollet Women’s Center and Regions Hospital – Cancer Care Center. Kandice dedicates herself to providing <b>nurturing, sensitive care</b> along with good educational support to her patients. Label: Nurse	[ <b>nurturing, sensitive, care</b> ]	[Kandice, works, nurse, Cancer Care Center, Women’s Center, provides, <b>sensitive care</b> , educational support, patients]	[ <b>care</b> , patients, [care, patients] support]	[ <b>care, patients</b> ]

Table 1: Examples from e-SNLI and MedicalBios datasets are presented, showcasing rationales from human annotations alongside those generated by three models (Using "Short Prompt Top-Var" method). Human rationales and the corresponding generated words are highlighted in **bold font**.

Lime (Ribeiro et al., 2016). Madsen et al. (2024) investigated LLM faithfulness on models like Llama2 (Touvron et al., 2023), Falcon (Penedo et al., 2023), and Mistral (Jiang et al., 2023), noting its dependence on both model and dataset. Despite this, there is still a scarcity of analyses comparing LLM-generated rationales to interpretability methods. To the best of our knowledge, no studies seem to have assessed human alignment and model faithfulness jointly or used fine-tuning to overcome faithfulness evaluation limitations and explore its effects.

### 3 Experimental Setup

#### 3.1 Datasets

We utilize two natural language classification datasets that have been annotated with human rationales indicating which input words were pivotal for the ground truth label. Table 1 shows examples of these datasets alongside the human rationale annotation and model-generated rationale.

**e-SNLI** This dataset (Camburu et al., 2018b) is a natural language inference task with 3 classes including Entailment, Contradiction, and Neutral, showing the relation between the premise and hypothesis sentences. This dataset is annotated for rationales supporting the classification label by DeYoung et al. (2019). We utilize 5,000 examples from the training set and 300 examples from the test set.

**MedicalBios** MedicalBios (Eberle et al., 2023) consists of human rationale annotations for a subset of

100 samples (five medical classes) from the BIOS dataset (De-Arteaga et al., 2019) for the occupation classification task.

#### 3.2 Models

We employ five of the latest large language models, encompassing both open-source and proprietary ones. From the open-source models, we utilize Llama2 (Touvron et al., 2023), LLama3, and Mistral (Jiang et al., 2023). For proprietary models, we include GPT3.5-Turbo and GPT4-Turbo (OpenAI et al., 2024). All models are prompted without sampling during generation, leading to deterministic outputs. You can see the exact model information in Table 6 in the appendix.

#### 3.3 Methods

##### 3.3.1 Prompting-Based Method

We employ various prompting strategies to explore the effects of prompt wording and model alignment in generating text of similar length to human-annotated rationale (Tables 12&13). The following prompts are used to evaluate these aspects.

We test two versions of prompts to examine how the clarity and length of the prompt influence the model’s output:

**Normal Prompt** This version provides a detailed explanation, including all the points the model should consider. It is longer and aims to ensure the model fully understands the task.

Method	Model Dataset Selection	MISTRAL-7B INSTRUCT-v0.2		LLAMA-2-7B CHAT		LLAMA-3-8B INSTRUCT		GPT-3.5 TURBO 1106		GPT-4 TURBO 2024-04-09	
		E-SNLI	MedBios	E-SNLI	MedBios	E-SNLI	MedBios	E-SNLI	MedBios	E-SNLI	MedBios
<b>ATTRIBUTION-BASED</b>											
ATTENTION	TOP-RATIO	28.59	29.34	40.92	22.11	23.80	15.01	-	-	-	-
SALIENCY	TOP-RATIO	37.16	28.97	28.64	31.35	32.19	27.80	-	-	-	-
INPUT×GRADIENT	TOP-RATIO	31.41	27.51	26.85	30.96	36.32	29.83	-	-	-	-
ATTENTION	TOP-VAR	36.03	41.01	<u>48.69</u>	35.08	30.56	24.30	-	-	-	-
SALIENCY	TOP-VAR	<u>46.07</u>	38.30	37.25	41.49	39.78	36.36	-	-	-	-
INPUT×GRADIENT	TOP-VAR	38.86	37.08	36.66	40.45	<b>44.76</b>	39.48	-	-	-	-
<b>PROMPTING-BASED</b>											
NORMAL PROMPT	UNBOUND	38.95	42.50	45.45	49.16	39.58	49.05	<u>43.96</u>	45.14	47.24	51.89
SHORT PROMPT	UNBOUND	41.84	42.26	43.54	43.74	43.60	51.28	<u>43.28</u>	45.28	46.58	53.50
NORMAL PROMPT	TOP-RATIO	41.17	41.17	45.90	43.98	35.98	47.58	36.58	43.68	45.27	47.81
SHORT PROMPT	TOP-RATIO	41.83	42.05	40.91	45.65	37.46	46.91	37.92	41.00	47.31	45.64
NORMAL PROMPT	TOP-VAR	45.93	<b>48.43</b>	<b>50.93</b>	<b>52.23</b>	41.67	<b>58.85</b>	43.85	<b>57.16</b>	<u>55.33</u>	<b>59.28</b>
SHORT PROMPT	TOP-VAR	<b>46.97</b>	<u>48.08</u>	47.34	<u>50.54</u>	<u>44.65</u>	<u>58.17</u>	<b>44.31</b>	<u>53.60</u>	<b>55.78</b>	<u>59.10</u>

Table 2: Human alignment F1 $\uparrow$  score in e-SNLI and MedicalBios datasets. Random baseline (Selecting Top-Var random words) is **27 $\pm$ 4** and **22 $\pm$ 1** for e-SNLI and MedicalBios respectively over 100 seeds. The top two alignments in each column are indicated by **bold** and underlined formatting.

**Short Prompt** Given that long prompts may confuse LLMs, we also use a shorter version that conveys the necessary information in a few sentences.

We also experiment with three versions of the introduced prompts to manage the number of words the model generates.

**Unbound Prompt** In this method, the model is not restricted in the number of words it can generate. It needs to establish the appropriate length autonomously.

**Top-Var Prompt** This prompt requires the model to generate exactly the same number of words as in the human rationale annotations for each sentence. This method controls for word count in our experiments, enabling us to assess model alignment independently of its word importance threshold.

**Top-Ratio Prompt** In this approach, the model is guided to identify the top  $k$  most important words within a given sentence. The value of  $k$  is derived from a predetermined percentage of the total word count in the input sentence, a ratio established based on the training set. For example, this ratio is set at 20% for sentences in the e-SNLI dataset and 13% for those in the MedicalBios dataset.

### 3.3.2 Attribution-Based Methods

We employ the Inseq library (Sarti et al., 2023) to implement attribution-based methods for LLMs. Specifically, we select three available options: (i) Attention Weight Attribution, which utilizes the model’s internal attention weights (Wiegrefe and

Pinter, 2019); (ii) Simple Gradients (Saliency), which is based on the gradients of the output with respect to the inputs (Simonyan et al., 2014b); and (iii) Input×Gradient, which factors in both the input vector size and the gradient in its calculations (Li et al., 2016). We choose these methods because of their demonstrated faithfulness in previous work on NLP models (Atanasova et al., 2020; Modarressi et al., 2022, 2023), and their potential for efficient execution on large language models with limited computational resources.

## 4 Results

In this section, we delve into utilizing both prompting-based and attribution-based approaches to extract rationale from the model, focusing on two aspects: human alignment and model faithfulness. Furthermore, we conduct fine-tuning experiments on open LLMs to examine how this process influences alignment and faithfulness.

### 4.1 Human Alignment

The annotated rationales provide explanations for the ground truth label. Therefore, for the evaluation of human alignment, we first request the model to provide a rationale for the provided label. With this, we create an array of binary values where for each word in the input sentence, we indicate a 1 if it is present among the generated words (0 otherwise). By comparing these vectors with the equivalent binary representations of human annotations, we

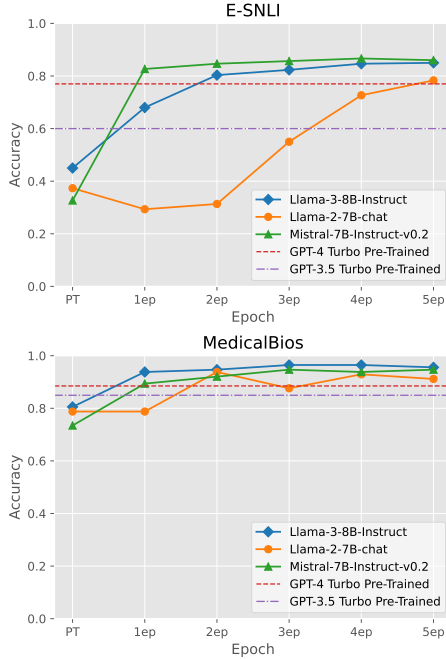


Figure 2: Accuracy changes throughout 5 epochs of fine-tuning. (PT denotes the pre-trained model’s accuracy)

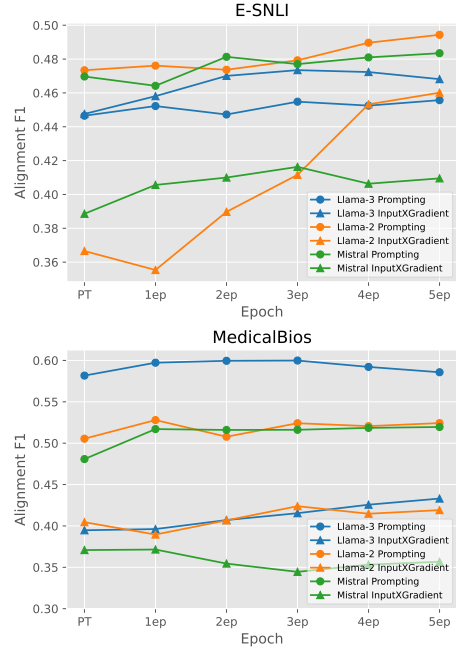


Figure 3: Human alignment F1↑ changes throughout 5 epochs of fine-tuning. (PT denotes the pre-trained model)

calculate the F1 score. This F1 score for human alignment is reported in Table 2 for the e-SNLI and MedicalBios datasets.

Firstly, prompting-based methods outperform attribution-based methods on human alignment. Short or normal prompting demonstrate superior performance compared to attribution-based methods in nearly all datasets and models except for Llama-3-8B on e-SNLI, which is also comparable. This gap can be attributed to the reliance of attribution-based methods on the classification capability of LLMs (which might be subpar).

Secondly, we note that providing additional information about the number of words selected by humans in *Top-Var* settings enhances alignment, indicating disparities between model thresholds for word importance in *Unbound* prompting compared to human annotators. Furthermore, the random baseline, which involves selecting *Top-Var* random words in each sentence, yields F1 scores of  $0.27 \pm 0.04$  and  $0.22 \pm 0.01$  for e-SNLI and MedicalBios respectively across 100 seeds. Contrasted with the *Top-Var* rows in Table 2, this indicates that both attribution-based and prompting-based approaches exhibit superior alignment compared to a random baseline.

Thirdly, across models, the performance comparison between normal and short prompts is varying and inconclusive. For example, short prompts

perform better for Mistral, LLama-3, and GPT-4-Turbo, while they perform worse for LLama-2 and GPT-3.5 models.

Finally, comparing the evaluated models reveals GPT-4-Turbo to be the most aligned with humans. However, other models demonstrate task-dependent alignment, with some excelling in e-SNLI and others in MedicalBios.

#### 4.2 Effect of Fine-tuning on Alignment

As described by Wang et al. (2024) and Zhong et al. (2023), zero-shot LLMs may underperform small fine-tuned models such as BERT. And smaller LLMs like LLaMA-2-7B might even collapse entirely. We observe similar failure patterns, where models generate a single label from the possible options regardless of the input sentence. This phenomenon is illustrated in Figure 2, where the pre-trained (PT) open LLMs achieve near-random accuracy (33%) on the e-SNLI dataset. This issue raises the question of whether fine-tuning these models to improve their classification performance also aids in aligning their explanations more closely with human expectations.

To address this issue, we fine-tune the LLMs using LoRA (Hu et al., 2022), a parameter-efficient fine-tuning technique. The hyperparameters for fine-tuning are provided in Table 7. Figure 2 demonstrates that the classification performance

Method	Model Dataset Selection	MISTRAL-7B FT INSTRUCT-V0.2		LLAMA-2-7B FT CHAT		LLAMA-3-8B FT INSTRUCT	
		E-SNLI	MedicalBios	E-SNLI	MedicalBios	E-SNLI	MedicalBios
<b>ATTRIBUTION-BASED</b>							
SALIENCY	TOP-RATIO	37.36 (+0.20)	27.29 (-1.69)	36.34 (+7.71)	31.73 (+0.38)	34.76 (+2.57)	30.18 (+2.38)
INPUT×GRADIENT	TOP-RATIO	33.71 (+2.30)	27.47 (-0.04)	36.98 (+10.13)	32.64 (+1.68)	38.59 (+2.27)	31.91 (+2.08)
SALIENCY	TOP-VAR	44.87 (-1.19)	37.12 (-1.18)	45.63 (+8.37)	42.80 (+1.31)	42.90 (+3.13)	39.57 (+3.21)
INPUT×GRADIENT	TOP-VAR	40.95 (+2.09)	35.67 (-1.41)	46.01 (+9.36)	41.91 (+1.46)	<b>46.81 (+2.06)</b>	<b>43.31 (+3.83)</b>
<b>PROMPTING-BASED</b>							
NORMAL PROMPT	UNBOUND	40.48 (+1.53)	41.66 (-0.84)	45.28 (-0.17)	47.26 (-1.89)	39.04 (-0.54)	48.41 (-0.64)
SHORT PROMPT	UNBOUND	39.49 (-2.35)	42.74 (+0.48)	43.11 (-0.43)	43.45 (-0.29)	43.86 (+0.26)	50.45 (-0.83)
NORMAL PROMPT	TOP-RATIO	40.99 (-0.17)	41.59 (+0.43)	44.40 (-1.50)	43.29 (-0.70)	35.42 (-0.56)	47.00 (-0.58)
SHORT PROMPT	TOP-RATIO	42.01 (+0.18)	43.91 (+1.85)	43.23 (+2.33)	45.29 (-0.36)	38.15 (+0.69)	47.68 (+0.77)
NORMAL PROMPT	TOP-VAR	45.70 (-0.24)	48.51 (+0.08)	<b>50.58 (-0.35)</b>	51.35 (-0.88)	41.69 (+0.02)	<b>59.70 (+0.85)</b>
SHORT PROMPT	TOP-VAR	<b>48.35 (+1.38)</b>	<b>51.96 (+3.88)</b>	49.43 (+2.09)	<b>52.43 (+1.89)</b>	45.57 (+0.92)	58.57 (+0.41)

Table 3: Human alignment F1 $\uparrow$  score in E-SNLI and MedicalBios datasets after fine-tuning for 5 epochs on the respective dataset. The difference in alignment between the fine-tuned and pre-trained model (Table 2) is reported in the parentheses. (Changes greater than 1.5% are denoted using a bold font.)

of LLaMA-2, LLaMA-3, and Mistral improves significantly after fine-tuning on the e-SNLI dataset and shows slight improvements on the MedicalBios dataset, outperforming GPT-4-Turbo in both cases.

We rerun the human alignment experiments for “Short Prompt Top-Var” and “Input×Gradient Top-Var” across all epochs, as shown in Figure 3. The results suggest a general trend of improved alignment of both methods with increasing epochs.

To analyze the generalization of this trend more comprehensively, we conduct alignment experiments on all prompting and attribution methods in Table 3 on the final epoch (5th epoch) of fine-tuning. The difference in human alignment between the fine-tuned and pre-trained models (Table 2) is reported in parentheses.

Overall, the results show more positive change than negative in alignment. Among the prompting-based methods, “Short Prompt Top-Var” demonstrates the highest gains from fine-tuning. Attribution-based methods, particularly in LLaMA models, exhibit significant alignment improvements after fine-tuning. Moreover, fine-tuning can guide the model’s attention to the correct words, especially in datasets like e-SNLI where pre-trained classification accuracy was low. We have demonstrated qualitative examples of such cases in Figure 5. As a result, these gradient-based methods can identify more human-aligned rationales by tracing back the attributions from the output label to the input sentence in fine-tuned models.

### 4.3 Faithfulness to the Model

While human alignment provides a useful measure of the plausibility of LLM rationales, it is also important to consider the faithfulness of these rationales to the model’s actual decision-making process. A word may be crucial for the model’s decision even if it does not align with human rationale and vice versa. Therefore, we must ask: Are the self-explanations genuinely influential in the model’s decision-making process?

To evaluate faithfulness, we employ a perturbation-based experiment similar to previous work (Madsen et al., 2024; Modarressi et al., 2023). In this experiment, we mask the important words identified by the prompting and attribution methods and measure the flip rate of the predicted label during classification. A higher flip rate indicates that the masked words are indeed important to the model, leading it to change its previous decision, and this suggests that the explanation is more faithful to the model’s decision-making process.

#### 4.3.1 Limitations of Faithfulness Evaluation before Fine-Tuning

Table 4 presents the faithfulness flip rate of the pre-trained LLMs. A noteworthy finding is that in the e-SNLI dataset, where classification accuracy was notably low (Figure 2), both attribution-based and prompting-based methods resulted in a very small flip rate. Even more concerning, masking all the words in the input sentence led to less than a

		Model	MISTRAL-7B INSTRUCT-v0.2		LLAMA-2-7B CHAT		LLAMA-3-8B INSTRUCT		GPT-3.5 TURBO 1106	
Masking Part Method		Dataset Selection	E-SNLI	MedBios	E-SNLI	MedBios	E-SNLI	MedBios	E-SNLI	MedBios
<b>ATTRIBUTION-BASED</b>										
INPUT	SALIENCY	TOP-RATIO	3.00	50.00	2.36	31.82	29.33	36.11	-	-
	INPUT×GRADIENT	TOP-RATIO	2.33	44.32	2.02	29.09	26.67	40.74	-	-
<b>PROMPTING-BASED</b>										
INPUT	NORMAL PROMPT	UNBOUND	3.33	59.09	1.01	30.91	38.67	41.67	77.00	41.82
	SHORT PROMPT	UNBOUND	3.67	59.09	1.68	35.45	44.00	47.22	73.67	37.27
	NORMAL PROMPT	TOP-RATIO	1.67	28.41	1.68	33.64	24.00	45.37	41.33	24.55
	SHORT PROMPT	TOP-RATIO	2.33	30.68	1.01	26.36	22.67	46.30	44.33	20.91
<b>BASELINES</b>										
INPUT	HUMAN	HUMAN	4.00	60.23	2.02	38.18	31.67	43.12	55.67	36.36
	RANDOM	RANDOM	1.33	9.09	1.01	10.91	24.67	14.68	41.33	5.45
	EVERYTHING	EVERYTHING	2.33	97.73	0.34	60.00	69.67	70.64	89.00	74.55

Table 4: Faithfulness FLIP RATE $\uparrow$  percentage in E-SNLI and MedicalBios datasets. The number of words to mask is enforced in TOP-RATIO, and no method could mask more than the specified number for each sentence.

3% flip rate for the Mistral and LLaMA-2 models (Mask EVERYTHING).

In further exploration, Figure 4 illustrates the Input×Gradient attributions of the predicted label, shown in green, to all instruction and input words, shown in shades of red. We notice that in the pre-trained Llama-2 model, the prediction for the label “entailment” is incorrect, with the model placing excessive emphasis on the word “entailment” in the instruction while largely ignoring the input sentence including the premise and hypothesis sentences. However, after fine-tuning, the attribution distribution becomes less skewed, leading to a correct prediction by the model.

Therefore, we hypothesize that the pre-trained model focuses more on the instruction rather than the input sentence. Moreover, in Table 10, we investigate two masking scenarios. The first scenario, denoted as INPUT, involves masking only the words from the input sentence (e.g., the premise and hypothesis in e-SNLI) while leaving the instruction intact, similar to Table 4. The second scenario, denoted as INPUT&INSTRUCTION, extends masking to the entire instruction and input, constituting the entire prompt. When we extend the masking to include the instruction, the flip rate can increase up to 100%. This indicates that the model relies heavily on the instruction for its decisions, regardless of the input sentence. This phenomenon aligns with findings by Yin et al. (2023) and Kung and Peng (2023), who both found that, among all segments of a prompt, label information or output space is

essential for the model’s performance. This raises concerns about the reliability of this experiment for measuring model faithfulness in LLMs.

Consequently, we argue that to conduct a more robust faithfulness experiment on LLMs, it is not advisable to solely rely on pre-trained models, as their classification accuracy can vary depending on the model and dataset (Madsen et al., 2024). Instead, we suggest aligning the experiment more closely with the scenario of fine-tuned encoder-based models (Ferrando et al., 2022; Modarressi et al., 2023) by training the LLMs and assessing faithfulness on the fine-tuned model.

### 4.3.2 Faithfulness after Fine-Tuning

Table 5 displays the faithfulness flip rate of the fine-tuned open models on e-SNLI and MedicalBios. The number of masked words can directly influence the prediction flip rate in this experiment. To ensure a fair comparison, we limit the number of words considered in the Top-Ratio and Top-Var selections to the top  $k$  words for each sentence across different methods. This approach is particularly important for prompting-based methods, as LLMs struggle to follow instructions involving fine-grained hard constraints (Sun et al., 2023).

First, we see that fine-tuning has effectively addressed the near-zero flip rate (Table 4) in e-SNLI, indicating that the model is no longer completely disregarding the input sentence.

Second, a comparison of results in each selection group of “Top-Ratio” and “Top-Var” (with a similar number of masked words) reveals that attribution-

Selection	Model Dataset Method	MISTRAL-7B FT INSTRUCT-v0.2		LLAMA-2-7B FT CHAT		LLAMA-3-8B FT INSTRUCT	
		E-SNLI	MedicalBios	E-SNLI	MedicalBios	E-SNLI	MedicalBios
UNBOUND	NORMAL PROMPT	<b>64.00</b>	34.51	67.00	17.70	64.33	18.75
	SHORT PROMPT	53.00	<b>37.17</b>	<b>67.33</b>	<b>21.24</b>	<b>72.00</b>	<b>25.00</b>
TOP-RATIO	ATTENTION	50.33	<b>23.89</b>	<b>69.67</b>	14.16	49.33	5.36
	SALIENCY	<b>51.00</b>	23.89	59.67	21.24	<b>60.00</b>	17.86
	INPUT×GRADIENT	45.33	21.24	61.33	22.12	58.00	18.75
	NORMAL PROMPT	35.33	17.70	52.33	19.47	45.00	<b>20.54</b>
	SHORT PROMPT	40.33	17.70	60.33	<b>23.89</b>	51.33	<b>20.54</b>
TOP-VAR	ATTENTION	48.00	<b>20.35</b>	<b>61.00</b>	11.50	48.67	5.36
	SALIENCY	<b>52.67</b>	14.16	54.67	17.70	55.00	14.29
	INPUT×GRADIENT	45.33	19.47	53.33	<b>19.47</b>	<b>56.00</b>	15.18
	NORMAL PROMPT	37.00	15.04	50.33	16.81	47.00	<b>16.07</b>
	SHORT PROMPT	40.33	19.47	55.00	17.70	50.67	<b>16.07</b>
<b>BASELINES</b>							
TOP-VAR	HUMAN	49.67	21.24	60.00	24.78	62.00	24.11
TOP-VAR	RANDOM	32.67	9.73	45.00	7.96	36.67	6.25
EVERYTHING	EVERYTHING	66.33	84.07	63.00	75.22	77.33	74.11

Table 5: Faithfulness FLIP RATE $\uparrow$  percentage for fine-tuned models on the respective dataset after 5 epochs. The number of words to mask is enforced in TOP-RATIO and TOP-VAR by limiting masked words per sentence. (Results are comparable in each similar selection technique, ensuring a similar number of masked words.)

based methods generally outperform prompting. Consistent results can be seen in our top- $k$  experiments in the appendix A.1 and Table 8. This difference can be attributed to the fact that attribution methods base their explanations on the model’s internal processes, whereas prompting may provide plausible answers without direct access to this information, potentially diverging from the truth of the model’s inner workings. Additionally, prompting is affected by the model’s ability to follow instructions, which may result in the generation of an inaccurate number of words or the inclusion of words not present in the input sentence, leading to less faithful results.

Third, we also present the flip rate after masking human rationales in Table 5. These rates are comparable to the “Top-Var” selection of other methods, as they involve the same number of masked words. Despite expectations that the model would better recognize the importance of words for its own decisions, prompting methods consistently underperformed human rationales, and attribution methods did so in half of the cases. This result emphasizes that while the current methods demonstrate a degree of faithfulness, there remains room for further refinement and enhancement.

## 5 Conclusions

In this study, we investigated the extraction of rationales in Large Language Models (LLMs) with a focus on human alignment and model faithfulness. Our experiments encompassed both prompting-based and attribution-based methods across various LLM architectures and datasets. Before fine-tuning, we observed that prompting generally yielded better human alignment, even when classification performance was poor. However, the reliability of faithfulness evaluations was compromised by low classification performance and collapsing predictions in pre-trained models highlighting the need for refining faithfulness evaluation setup.

To address this, we fine-tuned the models to enhance their accuracy on classification tasks, which led to improvements in aligning their explanations more closely with human expectations. In this scenario, although prompting showed superior alignment before, its faithfulness in reflecting model decision-making was not as strong as that of attribution-based methods.

Despite these improvements, a gap remained between the models’ rationales and human rationales in both alignment and faithfulness. This highlights the need for the development of more advanced explanation methods to bridge this gap.



## 505 Limitations

506 **LLM instruction-following abilities.** In our im- 555  
507 plementation of prompting strategies, we heavily 556  
508 rely on the LLM’s capability to follow instructions 557  
509 accurately. For example, when requesting the top-*k* 558  
510 words separated by a specific delimiter character, 559  
511 we expect the model to output a list of words in our 560  
512 desired format and quantity with no extra explana- 561  
513 tions. However, LLMs are still not fully adept at 562  
514 adhering to prompts precisely (Sun et al., 2023), 563  
515 which can lead to outputs in various formats differ- 564  
516 ent from our expectations. Since our primary focus 565  
517 in this paper is not to evaluate the format-following 566  
518 ability of LLMs, we have taken measures to address 567  
519 discrepancies in the outputs as much as possible. 568

520 To mitigate these discrepancies, we adopt tai- 569  
521 lored parsing approaches to handle unexpected out- 570  
522 put formats. For instance, if a model separates 571  
523 words in the output with a “;” character instead 572  
524 of the instructed character “|”, we adjust our pars- 573  
525 ing method accordingly. Fortunately, each model 574  
526 tends to adhere to a relatively consistent output 575  
527 format across the dataset, which enables us to 576  
528 adapt our parsing approach accordingly. Nonethe- 577  
529 less, it’s worth noting that an LLM with enhanced 578  
530 instruction-following abilities could potentially 579  
531 yield even better parsing results and consequently 580  
532 achieve higher performance levels. 581

533 **Attribution-based methods** In selecting the ex- 582  
534 planation methods based on the inner workings 583  
535 of the models we opted for the ones that were al- 584  
536 ready implemented for LLMs and were relatively 585  
537 efficient to execute given the large size of the 586  
538 models. Nonetheless, we acknowledge that recent 587  
539 vector-based methods have shown promising faith- 588  
540 fulness results by decomposing the representations 589  
541 (Kobayashi et al., 2020, 2021; Modarressi et al., 590  
542 2022; Ferrando et al., 2022; Modarressi et al., 2023) 591  
543 on smaller models such as BERT (Devlin et al., 592  
544 2019) compared with the gradient-based methods. 593  
545 Our study highlights the gap that could be filled by 594  
546 implementing these methods for LLMs. 595

547 **Prompt Engineering** Although we reported vari- 596  
548 ous versions of prompts for extracting rationales in 597  
549 this paper and conducted preliminary prompt engi- 598  
550 neering, we acknowledge that better prompts could 599  
551 potentially achieve higher performance. However, 600  
552 this approach diverges from realistic use cases 601  
553 where users may ask questions in various wordings. 602  
554 This limitation is inherent to prompting methods,

555 whereas attribution-based methods are not suscep- 556  
557 tible to this issue. Therefore, addressing this limita- 558  
559 tion calls for continued exploration and refinement 560  
561 of both prompting and attribution-based methods 562  
563 in rationale extraction. 564

565 **Larger Models** In our experiments, we evaluated 566  
567 open models with less than 8B parameters due to 568  
569 resource limitations. However, we acknowledge 570  
571 that larger models could potentially perform bet- 572  
573 ter in following instructions, leading to improved 574  
575 human alignment and model faithfulness in their 576  
577 self-explanations. 578

579 **Perturbation-based faithfulness evaluation** In 580  
581 this paper, we conduct faithfulness evaluation 582  
583 of LLM rationales using perturbation-based met- 584  
585 rics. Those metrics assume that removing criti- 586  
587 cal features based on rationales would largely 587  
588 affect model performance. However, Whether 589  
589 perturbation-based metrics truly reflect rationale 590  
591 faithfulness is a widely discussed but unsolved 591  
592 question, as they would produce out-of-distribution 592  
593 counterfactuals. For example, Yin et al. (2022) 593  
594 show that with different kinds of perturbations such 594  
595 as removal or noise in hidden representations, the 595  
596 faithful sets vary significantly. For consistency, 596  
597 we follow previous work (DeYoung et al., 2019; 597  
598 Huang et al., 2023a). We leave deeper study into 598  
599 faithfulness measurements of LLM rationales to 599  
600 future work. 600

## 601 References

- 602 Samira Abnar and Willem Zuidema. 2020. [Quantify-](#) 603  
604 [ing attention flow in transformers](#). In *Proceedings* 604  
605 *of the 58th Annual Meeting of the Association for* 605  
606 *Computational Linguistics*, pages 4190–4197, On- 606  
607 line. Association for Computational Linguistics. 607
- 608 Chirag Agarwal, Sree Harsha Tanneru, and Himabindu 608  
609 Lakkaraju. 2024. Faithfulness vs. plausibility: On the 609  
610 (un) reliability of explanations from large language 610  
611 models. *arXiv preprint arXiv:2402.04614*. 611
- 612 Pepa Atanasova, Jakob Grue Simonsen, Christina Li- 612  
613 oma, and Isabelle Augenstein. 2020. [A diagnostic](#) 613  
614 [study of explainability techniques for text classifica-](#) 614  
615 [tion](#). In *Proceedings of the 2020 Conference on* 615  
616 *Empirical Methods in Natural Language Processing* 616  
617 *(EMNLP)*, pages 3256–3274, Online. Association for 617  
618 Computational Linguistics. 618
- 619 Oana-Maria Camburu, Tim Rocktäschel, Thomas 619  
620 Lukaszewicz, and Phil Blunsom. 2018a. e-snli: Natu- 620  
621 ral language inference with natural language expla- 621  
622 nations. In *NeurIPS*. 622

605	Oana-Maria Camburu, Tim Rocktäschel, Thomas	Javier Ferrando, Gerard I. Gállego, and Marta R. Costa-	662
606	Lukasiewicz, and Phil Blunsom. 2018b. e-snli: Natu-	jussà. 2022. <a href="#">Measuring the mixing of contextual</a>	663
607	ral language inference with natural language explana-	<a href="#">information in the transformer</a> . In <i>Proceedings of</i>	664
608	tions. In <i>Advances in Neural Information Processing</i>	<i>of the 2022 Conference on Empirical Methods in Natu-</i>	665
609	<i>Systems</i> , pages 9539–9549.	<i>ral Language Processing</i> , pages 8698–8714, Abu	666
		Dhabi, United Arab Emirates. Association for Com-	667
610	Samuel Carton, Surya Kanoria, and Chenhao Tan. 2021.	putational Linguistics.	668
611	What to learn, and how: Toward effective learning		
612	from rationales. <i>ArXiv</i> , abs/2112.00071.	Peter Hase and Mohit Bansal. 2022. <a href="#">When can mod-</a>	669
		<a href="#">els learn from explanations? a formal framework</a>	670
613	Yanda Chen, Ruiqi Zhong, Narutatsu Ri, Chen Zhao,	<a href="#">for understanding the roles of explanation data</a> . In	671
614	He He, Jacob Steinhardt, Zhou Yu, and Kathleen	<i>Proceedings of the First Workshop on Learning with</i>	672
615	McKeown. 2023. Do models explain themselves?	<i>Natural Language Supervision</i> , pages 29–39, Dublin,	673
616	counterfactual simulatability of natural language ex-	Ireland. Association for Computational Linguistics.	674
617	planations. <i>arXiv preprint arXiv:2307.08678</i> .		
		Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan	675
618	Kevin Clark, Urvashi Khandelwal, Omer Levy, and	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and	676
619	Christopher D. Manning. 2019. <a href="#">What does BERT</a>	Weizhu Chen. 2022. <a href="#">LoRA: Low-rank adaptation of</a>	677
620	<a href="#">look at? an analysis of BERT’s attention</a> . In <i>Pro-</i>	<a href="#">large language models</a> . In <i>International Conference</i>	678
621	<i>ceedings of the 2019 ACL Workshop BlackboxNLP:</i>	<i>on Learning Representations</i> .	679
622	<i>Analyzing and Interpreting Neural Networks for NLP</i> ,		
623	pages 276–286, Florence, Italy. Association for Com-	Quzhe Huang, Shengqi Zhu, Yansong Feng, and	680
624	putational Linguistics.	Dongyan Zhao. 2021. <a href="#">Exploring distantly-labeled</a>	681
		<a href="#">rationales in neural network models</a> . In <i>Proceedings</i>	682
625	Maria De-Arteaga, Alexey Romanov, Hanna Wal-	<i>of the 59th Annual Meeting of the Association for</i>	683
626	lach, Jennifer Chayes, Christian Borgs, Alexandra	<i>Computational Linguistics and the 11th International</i>	684
627	Chouldechova, Sahin Geyik, Krishnaram Kenthapadi,	<i>Joint Conference on Natural Language Processing</i>	685
628	and Adam Tauman Kalai. 2019. <a href="#">Bias in bios: A case</a>	<i>(Volume 1: Long Papers)</i> , pages 5571–5582, Online.	686
629	<a href="#">study of semantic representation bias in a high-stakes</a>	Association for Computational Linguistics.	687
630	<a href="#">setting</a> . In <i>Proceedings of the Conference on Fair-</i>		
631	<i>ness, Accountability, and Transparency</i> , FAT* ’19,	Shiyuan Huang, Siddarth Mamidanna, Shreedhar	688
632	page 120–128, New York, NY, USA. Association for	Jangam, Yilun Zhou, and Leilani H Gilpin. 2023a.	689
633	Computing Machinery.	Can large language models explain themselves? a	690
		study of llm-generated self-explanations. <i>arXiv</i>	691
634	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	<i>preprint arXiv:2310.11207</i> .	692
635	Kristina Toutanova. 2019. <a href="#">BERT: Pre-training of</a>		
636	<a href="#">deep bidirectional transformers for language under-</a>	Shiyuan Huang, Siddarth Mamidanna, Shreedhar	693
637	<a href="#">standing</a> . In <i>Proceedings of the 2019 Conference of</i>	Jangam, Yilun Zhou, and Leilani H. Gilpin. 2023b.	694
638	<i>the North American Chapter of the Association for</i>	<a href="#">Can large language models explain themselves? a</a>	695
639	<i>Computational Linguistics: Human Language Tech-</i>	<a href="#">study of llm-generated self-explanations</a> . <i>Preprint</i> ,	696
640	<i>nologies, Volume 1 (Long and Short Papers)</i> , pages	arXiv:2310.11207.	697
641	4171–4186, Minneapolis, Minnesota. Association for		
642	Computational Linguistics.	Alon Jacovi and Yoav Goldberg. 2020. Towards faith-	698
		fully interpretable nlp systems: How should we de-	699
643	Jay DeYoung, Sarthak Jain, Nazneen Fatema Ra-	fine and evaluate faithfulness? In <i>Proceedings of the</i>	700
644	jani, Eric Lehman, Caiming Xiong, Richard Socher,	<i>58th Annual Meeting of the Association for Compu-</i>	701
645	and Byron C. Wallace. 2019. <a href="#">Eraser: A bench-</a>	<i>tational Linguistics</i> , pages 4198–4205.	702
646	<a href="#">mark to evaluate rationalized nlp models</a> . <i>Preprint</i> ,		
647	arXiv:1911.03429.	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan	703
		Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea	704
648	Jay DeYoung, Sarthak Jain, Nazneen Fatema Rajani,	Madotto, and Pascale Fung. 2023. Survey of halluci-	705
649	Eric Lehman, Caiming Xiong, Richard Socher, and	nation in natural language generation. <i>ACM Comput-</i>	706
650	Byron C. Wallace. 2020. <a href="#">ERASER: A benchmark to</a>	<i>ing Surveys</i> , 55(12):1–38.	707
651	<a href="#">evaluate rationalized NLP models</a> . In <i>Proceedings</i>		
652	<i>of the 58th Annual Meeting of the Association for</i>	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Men-	708
653	<i>Computational Linguistics</i> , pages 4443–4458, Online.	sch, Chris Bamford, Devendra Singh Chaplot, Diego	709
654	Association for Computational Linguistics.	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	710
		laume Lample, Lucile Saulnier, L�el�io Renard Lavaud,	711
655	Oliver Eberle, Ilias Chalkidis, Laura Cabello, and	Marie-Anne Lachaux, Pierre Stock, Teven Le Scao,	712
656	Stephanie Brandl. 2023. <a href="#">Rather a nurse than a physi-</a>	Thibaut Lavril, Thomas Wang, Timoth�ee Lacroix,	713
657	<a href="#">cian - contrastive explanations under investigation</a> .	and William El Sayed. 2023. <a href="#">Mistral 7b</a> . <i>Preprint</i> ,	714
658	In <i>Proceedings of the 2023 Conference on Empiri-</i>	arXiv:2310.06825.	715
659	<i>cal Methods in Natural Language Processing</i> , pages		
660	6907–6920, Singapore. Association for Computa-	Goro Kobayashi, Tatsuki Kuribayashi, Sho Yokoi, and	716
661	tional Linguistics.	Kentaro Inui. 2020. <a href="#">Attention is not only a weight:</a>	717
		<a href="#">Analyzing transformers with vector norms</a> . In	718



840	Shyam, Szymon Sidor, Eric Sigler, Maddie Simens,	4590–4605, Abu Dhabi, United Arab Emirates. As-	898
841	Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin	sociation for Computational Linguistics.	899
842	Sokolowsky, Yang Song, Natalie Staudacher, Fe-		
843	lipe Petroski Such, Natalie Summers, Ilya Sutskever,	Jiao Sun, Swabha Swayamdipta, Jonathan May, and	900
844	Jie Tang, Nikolas Tezak, Madeleine B. Thompson,	Xuezhe Ma. 2022b. <a href="#">Investigating the benefits of</a>	901
845	Phil Tillet, Amin Tootoonchian, Elizabeth Tseng,	<a href="#">free-form rationales</a> . In <i>Findings of the Association</i>	902
846	Preston Tuggle, Nick Turley, Jerry Tworek, Juan Fe-	<i>for Computational Linguistics: EMNLP 2022</i> , pages	903
847	lipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya,	5867–5882, Abu Dhabi, United Arab Emirates. As-	904
848	Chelsea Voss, Carroll Wainwright, Justin Jay Wang,	sociation for Computational Linguistics.	905
849	Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei,		
850	CJ Weinmann, Akila Welihinda, Peter Welinder, Ji-	Jiao Sun, Yufei Tian, Wangchunshu Zhou, Nan Xu, Qian	906
851	ayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner,	Hu, Rahul Gupta, John Wieting, Nanyun Peng, and	907
852	Clemens Winter, Samuel Wolrich, Hannah Wong,	Xuezhe Ma. 2023. <a href="#">Evaluating large language models</a>	908
853	Lauren Workman, Sherwin Wu, Jeff Wu, Michael	<a href="#">on controlled generation tasks</a> . In <i>Proceedings of the</i>	909
854	Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qim-	<i>2023 Conference on Empirical Methods in Natural</i>	910
855	ing Yuan, Wojciech Zaremba, Rowan Zellers, Chong	<i>Language Processing</i> , pages 3155–3168, Singapore.	911
856	Zhang, Marvin Zhang, Shengjia Zhao, Tianhao	Association for Computational Linguistics.	912
857	Zheng, Juntang Zhuang, William Zhuk, and Bar-		
858	ret Zoph. 2024. <a href="#">Gpt-4 technical report</a> . <i>Preprint</i> ,	Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017.	913
859	arXiv:2303.08774.	<a href="#">Axiomatic attribution for deep networks</a> . <i>Preprint</i> ,	914
		arXiv:1703.01365.	915
860	Guilherme Penedo, Quentin Malartic, Daniel Hesslow,		
861	Ruxandra Cojocaru, Alessandro Cappelli, Hamza	Gemini Team, Rohan Anil, Sebastian Borgeaud,	916
862	Alobeidli, Baptiste Pannier, Ebtesam Almazrouei,	Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu,	917
863	and Julien Launay. 2023. <a href="#">The refinedweb dataset for</a>	Radu Soricut, Johan Schalkwyk, Andrew M Dai,	918
864	<a href="#">falcon llm: Outperforming curated corpora with web</a>	Anja Hauth, et al. 2023. Gemini: a family of	919
865	<a href="#">data, and web data only</a> . <i>Preprint</i> , arXiv:2306.01116.	highly capable multimodal models. <i>arXiv preprint</i>	920
		<i>arXiv:2312.11805</i> .	921
866	Marco Tulio Ribeiro, UW EDU, Sameer Singh, and Car-	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	922
867	los Guestrin. 2016. Model-Agnostic Interpretability	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	923
868	of Machine Learning. In <i>ICML Workshop on Human</i>	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	924
869	<i>Interpretability in Machine Learning</i> .	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton	925
		Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	926
870	Gabriele Sarti, Nils Feldhus, Ludwig Sickert, Oskar	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,	927
871	van der Wal, Malvina Nissim, and Arianna Bisazza.	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-	928
872	2023. <a href="#">Inseq: An interpretability toolkit for se-</a>	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	929
873	<a href="#">quence generation models</a> . In <i>Proceedings of the</i>	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,	930
874	<i>61st Annual Meeting of the Association for Computa-</i>	Isabel Kloumann, Artem Korenev, Punit Singh Koura,	931
875	<i>tional Linguistics (Volume 3: System Demonstra-</i>	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-	932
876	<i>tions)</i> , pages 421–435, Toronto, Canada. Association	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	933
877	for Computational Linguistics.	tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-	934
		bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	935
878	Karen Simonyan, Andrea Vedaldi, and Andrew Zisser-	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	936
879	man. 2014a. <a href="#">Deep inside convolutional networks:</a>	Ruan Silva, Eric Michael Smith, Ranjan Subrama-	937
880	<a href="#">Visualising image classification models and saliency</a>	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-	938
881	<a href="#">maps</a> . <i>Preprint</i> , arXiv:1312.6034.	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,	939
		Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,	940
882	Karen Simonyan, Andrea Vedaldi, and Andrew Zisser-	Melanie Kambadur, Sharan Narang, Aurelien Rod-	941
883	man. 2014b. <a href="#">Deep inside convolutional networks:</a>	riguez, Robert Stojnic, Sergey Edunov, and Thomas	942
884	<a href="#">Visualising image classification models and saliency</a>	Scialom. 2023. <a href="#">Llama 2: Open foundation and fine-</a>	943
885	<a href="#">maps</a> . <i>CoRR</i> , abs/1312.6034.	<a href="#">tuned chat models</a> . <i>Preprint</i> , arXiv:2307.09288.	944
886	Julia Strout, Ye Zhang, and Raymond Mooney. 2019.	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	945
887	<a href="#">Do human rationales improve machine explanations?</a>	Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz	946
888	In <i>Proceedings of the 2019 ACL Workshop Black-</i>	Kaiser, and Illia Polosukhin. 2023. <a href="#">Attention is all</a>	947
889	<i>boxNLP: Analyzing and Interpreting Neural Net-</i>	<a href="#">you need</a> . <i>Preprint</i> , arXiv:1706.03762.	948
890	<i>works for NLP</i> , pages 56–62, Florence, Italy. As-		
891	sociation for Computational Linguistics.		
892	Jiao Sun, Anjali Narayan-Chen, Shereen Oraby,	Yuxia Wang, Minghan Wang, and Preslav Nakov. 2024.	949
893	Alessandra Cervone, Tagyoung Chung, Jing Huang,	<a href="#">Rethinking STS and NLI in large language models</a> .	950
894	Yang Liu, and Nanyun Peng. 2022a. <a href="#">ExPUNations:</a>	In <i>Findings of the Association for Computational</i>	951
895	<a href="#">Augmenting puns with keywords and explanations</a> .	<i>Linguistics: EACL 2024</i> , pages 965–982, St. Julian’s,	952
896	In <i>Proceedings of the 2022 Conference on Empiri-</i>	Malta. Association for Computational Linguistics.	953
897	<i>cal Methods in Natural Language Processing</i> , pages		

954 Sarah Wiegrefe and Yuval Pinter. 2019. Attention is not  
 955 not explanation. In *Proceedings of the 2019 Confer-*  
 956 *ence on Empirical Methods in Natural Language Pro-*  
 957 *cessing and the 9th International Joint Conference*  
 958 *on Natural Language Processing (EMNLP-IJCNLP)*,  
 959 pages 11–20.

960 Fan Yin, Zhouxing Shi, Cho-Jui Hsieh, and Kai-Wei  
 961 Chang. 2022. On the sensitivity and stability of  
 962 model interpretations in nlp. In *Proceedings of the*  
 963 *60th Annual Meeting of the Association for Compu-*  
 964 *tational Linguistics (Volume 1: Long Papers)*, pages  
 965 2631–2647.

966 Fan Yin, Jesse Vig, Philippe Laban, Shafiq Joty, Caim-  
 967 ing Xiong, and Chien-Sheng Wu. 2023. [Did you read](#)  
 968 [the instructions? rethinking the effectiveness of task](#)  
 969 [definitions in instruction learning](#). In *Proceedings*  
 970 *of the 61st Annual Meeting of the Association for*  
 971 *Computational Linguistics (Volume 1: Long Papers)*,  
 972 pages 3063–3079, Toronto, Canada. Association for  
 973 Computational Linguistics.

974 Qihuang Zhong, Liang Ding, Juhua Liu, Bo Du, and  
 975 Dacheng Tao. 2023. [Can chatgpt understand too?](#)  
 976 [a comparative study on chatgpt and fine-tuned bert.](#)  
 977 *Preprint*, arXiv:2302.10198.

## 978 A Appendix

### 979 A.1 Top-k Faithfulness

980 As previously noted, the number of masked words  
 981 significantly impacts the flip rate. To explore this  
 982 further, we conducted a Top- $k$  experiment, mask-  
 983 ing  $k = 1, 2, 3, 4, 5, 10$  words and calculating the  
 984 flip rate (see Table 8). The results consistent with  
 985 Table 5 indicate that on average, attribution-based  
 986 methods surpass prompting in terms of faithfulness.

Model	Access
meta-llama/Meta-Llama-3-8B-Instruct	Open Source
meta-llama/Llama-2-7b-chat-hf	Open Source
mistralai/Mistral-7B-Instruct-v0.2	Open Source
gpt-3.5-turbo-1106	Proprietary
gpt-4-turbo-2024-04-09	Proprietary

Table 6: The details of the models we used in this work.

Hyperparameter	Value
Total Batch Size	64
Learning Rate E-SNLI	1e-05
Learning Rate MedicalBios	5e-06
Num Epochs	5
Learning Rate Scheduler Warmup Steps	10
Training Dataset Size	5000
LoRA r	32
LoRA alpha	16
LoRA drop out	0.05

Table 7: The hyperparameters used for fine-tuning the models using LoRA.

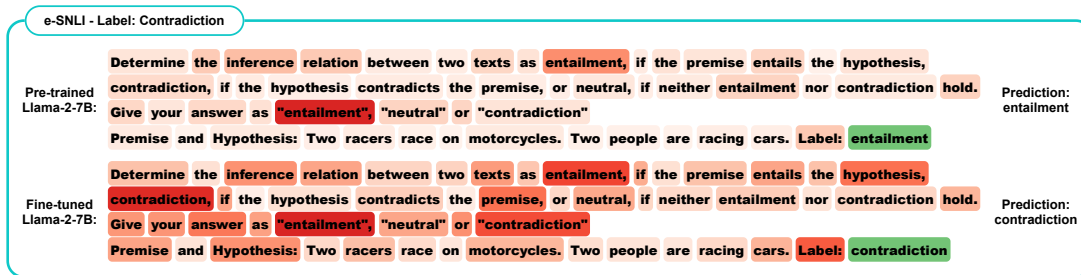


Figure 4: Token Importance before and after fine-tuning Llama-2-7B based on the InputXGradient explainability method. The predicted word by the model is shown in green and its attributions to previous words are shown in red. The attributions before fine-tuning are more skewed (Fisher-Pearson coefficient of skewness over all the dataset:  $3.37 \pm 0.31$ ), and become less skewed after fine-tuning ( $1.42 \pm 0.36$ ).



Figure 5: Token Importance before and after fine-tuning Llama-2-7B based on the InputXGradient explainability method. The predicted/true label is shown in green, with its attributions to input words in red. The human rationale for the examples are ["dog", "cat"], ["twins"], and ["grazes", "touching"]. The fine-tuned model identified these words solely through training on classification data, without any rationale data.

Dataset	Model Top k Method	MISTRAL-7B FT INSTRUCT-v0.2							LLAMA-2-7B FT CHAT							LLAMA-3-8B FT INSTRUCT						
		1	2	3	4	5	10	Avg	1	2	3	4	5	10	Avg	1	2	3	4	5	10	Avg
E-SNLI	ATTENTION	<b>41.0</b>	36.7	48.0	47.7	54.0	62.3	48.3	<b>40.0</b>	<b>50.0</b>	<b>63.7</b>	<b>66.3</b>	<b>69.3</b>	<b>71.0</b>	<b>60.1</b>	31.0	35.3	41.7	49.7	56.0	69.7	47.2
	SALIENCY	40.0	<b>46.0</b>	<b>49.3</b>	<b>52.3</b>	<b>54.7</b>	<b>63.0</b>	<b>50.9</b>	36.7	45.0	54.3	58.0	59.7	68.3	53.7	30.3	37.7	50.0	56.3	<b>63.7</b>	77.0	52.5
	INPUT×GRADIENT	34.0	39.0	44.0	46.7	50.0	57.7	45.2	37.0	42.3	50.7	60.3	60.7	70.0	53.5	<b>33.0</b>	<b>38.0</b>	<b>54.3</b>	<b>59.0</b>	61.0	<b>77.3</b>	<b>53.8</b>
	PROMPTING	31.3	30.0	41.3	40.7	47.0	55.3	40.9	35.3	44.3	54.7	60.7	64.3	64.0	53.9	29.3	28.7	47.0	48.7	57.7	76.7	48.0
MedicalBios	ATTENTION	7.1	10.6	15.0	15.0	19.5	31.0	16.4	5.3	8.0	7.1	8.8	9.7	13.3	8.7	0.9	3.6	3.6	5.4	6.2	15.2	5.8
	SALIENCY	5.3	12.4	14.2	<b>17.7</b>	<b>20.4</b>	<b>31.0</b>	<b>16.8</b>	14.2	<b>18.6</b>	18.6	15.9	16.8	18.6	17.1	7.1	11.6	13.4	17.0	<b>18.8</b>	18.8	14.4
	INPUT×GRADIENT	<b>8.8</b>	<b>16.8</b>	15.0	14.2	17.7	25.7	16.4	<b>15.9</b>	17.7	<b>19.5</b>	15.0	<b>18.6</b>	<b>19.5</b>	<b>17.7</b>	<b>8.0</b>	<b>12.5</b>	<b>14.3</b>	16.1	17.0	22.3	<b>15.0</b>
	PROMPTING	7.1	10.6	<b>16.8</b>	16.8	17.7	28.3	16.2	0.9	9.7	14.2	<b>19.5</b>	17.7	<b>19.5</b>	13.6	6.2	7.1	12.5	<b>19.6</b>	16.1	<b>27.7</b>	14.9

Table 8: Faithfulness FLIP RATE $\uparrow$  percentage in E-SNLI and MedicalBios datasets. The Top- $k$  is enforced in masking and no method could mask more than the specified  $k$ . The highest faithfulness in each dataset and top- $k$  is in **bold** font.

Dataset	Model Top k Method	MISTRAL-7B FT INSTRUCT-v0.2							LLAMA-2-7B FT CHAT							LLAMA-3-8B FT INSTRUCT						
		1	2	3	4	5	10	1	2	3	4	5	10	1	2	3	4	5	10			
E-SNLI	ATTENTION	1.0	2.0	3.0	4.0	5.0	10.0	1.0	2.0	3.0	4.0	5.0	10.0	1.0	2.0	3.0	4.0	5.0	10.0			
	SALIENCY	1.0	2.0	3.0	4.0	5.0	10.0	1.0	2.0	3.0	4.0	5.0	10.0	1.0	2.0	3.0	4.0	5.0	10.0			
	INPUT×GRADIENT	1.0	2.0	3.0	4.0	5.0	10.0	1.0	2.0	3.0	4.0	5.0	10.0	1.0	2.0	3.0	4.0	5.0	10.0			
	PROMPTING	1.0	1.8	2.9	3.8	4.7	8.2	0.9	1.9	2.9	3.8	4.8	7.4	0.9	1.8	2.8	3.8	4.8	9.1			
MedicalBios	ATTENTION	1.0	2.0	3.0	4.0	5.0	10.0	1.0	2.0	3.0	4.0	5.0	10.0	1.0	2.0	3.0	4.0	5.0	10.0			
	SALIENCY	1.0	2.0	3.0	4.0	5.0	10.0	1.0	2.0	3.0	4.0	5.0	10.0	1.0	2.0	3.0	4.0	5.0	10.0			
	INPUT×GRADIENT	1.0	2.0	3.0	4.0	5.0	10.0	1.0	2.0	3.0	4.0	5.0	10.0	1.0	2.0	3.0	4.0	5.0	10.0			
	PROMPTING	1.0	1.8	2.9	3.9	4.9	9.9	0.2	1.6	2.7	3.7	4.7	8.5	0.8	1.7	2.7	3.7	4.5	7.7			

Table 9: The average number of masked words in each experiment of Table 8. Prompting can have fewer masked words as a result of generating words out of the input sentence or not following the instruction on how many words it should generate (We limit the masks if the model generates more than  $k$  words).

Masking Part	Method	Model Dataset Selection	MISTRAL-7B INSTRUCT-v0.2		LLAMA-2-7B CHAT		LLAMA-3-8B INSTRUCT		GPT-3.5 TURBO 1106	
			E-SNLI	MedicalBios	E-SNLI	MedicalBios	E-SNLI	MedicalBios	E-SNLI	MedicalBios
<b>ATTRIBUTION-BASED</b>										
INPUT	SALIENCY	TOP-RATIO	3.00	50.00	2.36	31.82	29.33	36.11	-	-
	INPUT×GRADIENT	TOP-RATIO	2.33	44.32	2.02	29.09	26.67	40.74	-	-
INPUT&INSTR.	SALIENCY	TOP-RATIO	97.00	71.59	100.00	90.91	86.00	72.22	-	-
	INPUT×GRADIENT	TOP-RATIO	97.33	63.64	100.00	93.64	86.33	74.07	-	-
<b>PROMPTING-BASED</b>										
INPUT	NORMAL PROMPT	UNBOUND	3.33	59.09	1.01	30.91	38.67	41.67	77.00	41.82
	SHORT PROMPT	UNBOUND	3.67	59.09	1.68	35.45	44.00	47.22	73.67	37.27
	NORMAL PROMPT	TOP-RATIO	1.67	28.41	1.68	33.64	24.00	45.37	41.33	24.55
	SHORT PROMPT	TOP-RATIO	2.33	30.68	1.01	26.36	22.67	46.30	44.33	20.91
INPUT&INSTR.	EXTENDED PROMPT	UNBOUND	68.67	61.36	8.75	28.18	58.33	45.37	99.67	50.00
<b>BASELINES</b>										
INPUT	HUMAN	HUMAN	4.00	60.23	2.02	38.18	31.67	43.12	55.67	36.36
	EVERYTHING	EVERYTHING	2.33	97.73	0.34	60.00	69.67	70.64	89.00	74.55

Table 10: Faithfulness FLIP RATE $\uparrow$  percentage in E-SNLI and MedicalBios datasets. The number of words to mask is enforced in TOP-RATIO, and no method could mask more than the specified number for each sentence.

Selection	Model Dataset Method	MISTRAL-7B FT INSTRUCT-v0.2		LLAMA-2-7B FT CHAT		LLAMA-3-8B FT INSTRUCT	
		E-SNLI	MedicalBios	E-SNLI	MedicalBios	E-SNLI	MedicalBios
UNBOUND	NORMAL PROMPT	64.00	34.51	67.00	17.70	64.33	18.75
	SHORT PROMPT	53.00	37.17	67.33	21.24	72.00	25.00
TOP-RATIO	ATTENTION	50.67	28.32	70.33	15.04	53.67	9.82
	SALIENCY	52.33	30.09	62.67	24.78	66.00	23.21
	INPUT×GRADIENT	47.67	25.66	62.33	26.55	62.67	25.89
TOP-VAR	NORMAL PROMPT	54.67	32.74	63.00	23.01	61.67	24.11
	SHORT PROMPT	52.67	35.40	64.00	25.66	59.00	25.00
	ATTENTION	47.67	23.01	67.33	15.93	53.67	6.25
TOP-VAR	SALIENCY	50.00	19.47	58.67	26.55	60.00	17.86
	INPUT×GRADIENT	44.33	20.35	58.00	24.78	61.67	19.64
	NORMAL PROMPT	51.33	30.97	64.67	19.47	61.33	21.43
	SHORT PROMPT	50.67	30.09	60.67	20.35	59.00	20.54
<b>BASELINES</b>							
TOP-VAR	HUMAN	49.67	21.24	60.00	24.78	62.00	24.11
EVERYTHING	EVERYTHING	66.33	84.07	63.00	75.22	77.33	74.11

Table 11: Faithfulness FLIP RATE $\uparrow$  percentage in E-SNLI and MedicalBios datasets for fine-tuned models on the respective dataset for 5 epochs. No limitations were applied to the number of masked words which means that LLMs could mask more words in prompting techniques by generating more words than instructed.

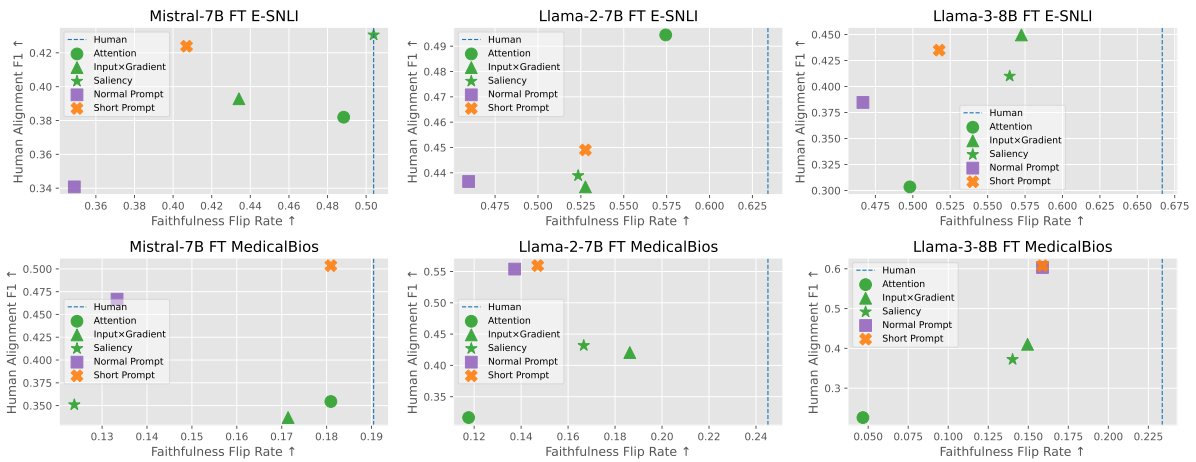


Figure 6: The human alignment (F1) and model faithfulness flip rate, evaluated at the intersection of correctly predicted examples (important for alignment) and the samples that were predicted within the output space of classes prior to fine-tuning (important for faithfulness). Additionally, the number of masks is constrained to the Top-Var selection.



Method	Selection	Prompt
NORMAL PROMPT	UNBOUND	<p>You will be provided a short premise and a hypothesis, each delimited with triple backticks.</p> <p>Your task is to identify the most important key words in the premise and hypothesis that are most relevant to understanding the relationship between the premise and the hypothesis given a label (entailment/neutral/contradiction).</p> <ul style="list-style-type: none"> <li>• These key words should be selected to capture the main idea or reasoning behind why the given label indicates the relationship between the premise and the hypothesis.</li> <li>• The key words should be copied verbatim from the text.</li> <li>• Write the key words as a pipe-separated ( ) list of words with spaces.</li> <li>• Example output: &lt;word1&gt;   &lt;word2&gt;   &lt;word3&gt;</li> <li>• Make sure to only output keyword tokens and do not explain the answer.</li> </ul> <p>Premise: ``` {premise} ```  Hypothesis: ``` {hypothesis} ```  Label: {label}</p>
SHORT PROMPT	UNBOUND	<p>Identify the most important key words from the text that are most relevant to understanding the relationship between the premise and the hypothesis.</p> <p>Write the words as a pipe-separated ( ) list of words with spaces and do not explain the answer.</p> <p>Premise: ``` {premise} ```  Hypothesis: ``` {hypothesis} ```  Label: {label}</p>
EXTENDED PROMPT	UNBOUND	<p>Consider the following prompt: " Determine the inference relation between two (short, ordered) texts as entailment, if the premise entails the hypothesis, contradiction, if the hypothesis contradicts the premise, or neutral, if neither entailment nor contradiction hold.</p> <p>Give your answer as "entailment", "neutral" or "contradiction" and nothing else.</p> <p>Premise and Hypothesis: ``` {premise} {hypothesis} ```  Label: "</p> <p>Identify the most important single key words from the whole prompt (Including instruction and sentence) that are most relevant to understanding the relationship between the premise and the hypothesis as {label}.</p> <p>Just write the words as a pipe-separated ( ) list of words with spaces and do not explain the answer.</p>
NORMAL PROMPT	TOP-VAR / TOP-RATIO	<p>You will be provided a short premise and a hypothesis, each delimited with triple backticks.</p> <p>Your task is to identify the top {k} most important key words in the premise and hypothesis that are most relevant to understanding the relationship between the premise and the hypothesis given a label (entailment/neutral/contradiction).</p> <ul style="list-style-type: none"> <li>• These key words should be selected to capture the main idea or reasoning behind why the given label indicates the relationship between the premise and the hypothesis.</li> <li>• The key words should be copied verbatim from the text.</li> <li>• Write the key words as a pipe-separated ( ) list of words with spaces.</li> <li>• Example output: &lt;word1&gt;   &lt;word2&gt;   &lt;word3&gt;</li> <li>• Make sure to only output top {k} keyword tokens and do not explain the answer.</li> </ul> <p>Premise: ``` {premise} ```  Hypothesis: ``` {hypothesis} ```  Label: {label}</p>
SHORT PROMPT	TOP-VAR / TOP-RATIO	<p>Identify the top {k} most important key words from the text that are most relevant to understanding the relationship between the premise and the hypothesis.</p> <p>Write the top {k} words as a pipe-separated ( ) list of words with spaces and do not explain the answer.</p> <p>Premise: ``` {premise} ```  Hypothesis: ``` {hypothesis} ```  Label: {label}</p>
ATTRIBUTION-BASED	TOP-VAR / TOP-RATIO (Selected later from tokens with the highest attribution scores)	<p>Determine the inference relation between two (short, ordered) texts as entailment, if the premise entails the hypothesis, contradiction, if the hypothesis contradicts the premise, or neutral, if neither entailment nor contradiction hold.</p> <p>Give your answer as "entailment", "neutral" or "contradiction" and nothing else.</p> <p>Premise and Hypothesis: ``` {premise} {hypothesis} ```  Label: {predicted_class}</p>
CLASSIFICATION		<p>Determine the inference relation between two (short, ordered) texts as entailment, if the premise entails the hypothesis, contradiction, if the hypothesis contradicts the premise, or neutral, if neither entailment nor contradiction hold.</p> <p>Give your answer as "entailment", "neutral" or "contradiction" and nothing else.</p> <p>Premise and Hypothesis: ``` {premise} {hypothesis} ```  Label:</p>

Table 12: The prompts utilized for the e-SNLI dataset.

Method	Selection	Prompt
NORMAL PROMPT	UNBOUND	<p>You will be provided a biography which is delimited with triple backticks. Your task is to identify the most important key words from the text that indicate why is the person in the following short bio described as a {label}.</p> <ul style="list-style-type: none"> <li>• The key words should be copied verbatim from the text.</li> <li>• Write the key words as a pipe-separated ( ) list of single words.</li> <li>• Make sure to only output keyword tokens and do not explain the answer.</li> <li>• Example output: &lt;word1&gt;   &lt;word2&gt;   &lt;word3&gt;</li> </ul> <p>Bio: "" {bio} ""</p>
SHORT PROMPT	UNBOUND	<p>Identify the most important key words from the text that indicate why is the person in the following short bio described as a {label}.</p> <ul style="list-style-type: none"> <li>• Write the key words as a pipe-separated ( ) list of single words.</li> </ul> <p>Bio: "" {bio} ""</p>
EXTENDED PROMPT	UNBOUND	<p>Consider the following prompt: " Determine the occupation of the person in the following short bio. Give your answer as a single word from (psychologist, surgeon, nurse, dentist, physician) and nothing else.</p> <p>Bio: "" {bio} "" Occupation: "</p> <p>Identify the most important single key words from the whole prompt (Including instruction and bio) that indicate why is the person in the short bio described as a {label}.</p> <p>Just write the words as a pipe-separated ( ) list of words with spaces and do not explain the answer.</p>
NORMAL PROMPT	TOP-VAR / TOP-RATIO	<p>You will be provided a biography which is delimited with triple backticks. Your task is to identify the top {k} most important key words from the text that indicate why is the person in the following short bio described as a {label}.</p> <ul style="list-style-type: none"> <li>• The key words should be copied verbatim from the text.</li> <li>• Write the key words as a pipe-separated ( ) list of single words.</li> <li>• Make sure to only output top {k} keyword tokens and do not explain the answer.</li> <li>• Example output: &lt;word1&gt;   &lt;word2&gt;   &lt;word3&gt;</li> </ul> <p>Bio: "" {bio} ""</p>
SHORT PROMPT	TOP-VAR / TOP-RATIO	<p>Identify the top {k} most important key words from the text that indicate why is the person in the following short bio described as a {label}.</p> <ul style="list-style-type: none"> <li>• Write the key words as a pipe-separated ( ) list of single words and do not explain.</li> </ul> <p>Bio: "" {bio} ""</p>
ATTRIBUTION-BASED	TOP-VAR / TOP-RATIO (Selected later from tokens with the highest attribution scores)	<p>Determine the occupation of the person in the following short bio. Give your answer as a single word from (psychologist, surgeon, nurse, dentist, physician) and nothing else.</p> <p>Bio: "" {bio} "" Occupation: {predicted_label}</p>
CLASSIFICATION		<p>Determine the occupation of the person in the following short bio. Give your answer as a single word from (psychologist, surgeon, nurse, dentist, physician) and nothing else.</p> <p>Bio: "" {bio} "" Occupation:</p>

Table 13: The prompts utilized for the MedicalBios dataset.