ALIGNMENT BETWEEN THE DECISION-MAKING LOGIC OF LLMS AND HUMAN COGNITION: A CASE STUDY ON LEGAL LLMS

Anonymous authors

Paper under double-blind review

ABSTRACT

This paper presents a method to evaluate the alignment between the decisionmaking logic of Large Language Models (LLMs) and human cognition in a case study on legal LLMs. Unlike traditional evaluations on language generation results, we propose to evaluate the correctness of the detailed decision-making logic of an LLM behind its seemingly correct outputs, which represents the core challenge for an LLM to earn human trust. To this end, we quantify the interactions encoded by the LLM as primitive decision-making logic, because recent theoretical achievements (Li & Zhang, 2023; Ren et al., 2024) have proven several mathematical guarantees of the faithfulness of the interaction-based explanation. We design a set of metrics to evaluate the detailed decision-making logic of LLMs. Experiments show that even when the language generation results appear correct, a significant portion of the internal inference logic contains notable issues¹.

025

006

008 009 010

011

013

014

015

016

017

018

019

021

023

026 027

1 INTRODUCTION

The trustworthiness and safety of Large Language Models (LLMs) present significant challenges for their deployment in high-stake tasks (OpenAI, 2023; Wei et al., 2023). Previous evaluation methods mainly evaluated the correctness of language generation results, in terms of value alignment and hallucination problems (Bang et al., 2023; Ji et al., 2023b;a; Shen et al., 2023).

In this study, we hope to go beyond the long-tail evaluation of the generation results, and focus on the 032 correctness of the detailed decision-making logic used by the LLM behind the language generation 033 result. We focus on the legal LLM as a case study, and the legal LLM may use significantly in-034 correct information to make judgment, even when the generation result is correct. The alignment of 035 decision-making logic between the AI model and human cognition is crucial for alleviating the com-036 mon fear of AI models. The alignment of internal logic via communication is the reason why people 037 naturally trust each other. Particularly, in high-stakes tasks such as autonomous driving (Grigorescu et al., 2020), the lack of alignment between AI models and human users makes people would rather delegate work to humans and tolerate potential errors, than trust highly accurate AI models. 040

Therefore, this paper aims to explore the possibility of aligning the decision-making logic for the 041 confidence score of the LLM's judgment with human cognition. To this end, exploring the mathe-042 matical feasibility of faithfully explaining the output score of a neural network as a few interpretable 043 logical patterns has become a new emerging theoretical problem in explainable AI, and about 20 044 papers have been published in three years (see related work in Appendix A). Typically, Li & Zhang (2023); Ren et al. (2024) have proved the universal-matching property and sparsity proporty, and 046 mathematically guaranteed that a DNN usually only encodes a small number of interactions between 047 input variables, and these interactions act as **primitive decision-making logic**, which well predicts 048 the confidence of the network prediction on various input variations.

As Figure 1 shows, an *interaction* measures the nonlinear relationship between input tokens of an input legal case encoded by the LLM. For instance, given an input sentence such as "*Andy threatened Bob and took his smartphone*," the LLM may trigger an interaction between a set of input tokens S =

¹The names used in the legal cases follow an alphabetical convention, *e.g.*, Andy, Bob, Charlie, etc., which do not represent any bias against actual individuals.



Figure 1: AND-OR interactions that explain the decision-making logic of a legal LLM. The surrogate logical model well estimates the confidence of the LLM making the judgment "*Robbery*" for Andy, $h("Robbery"|\mathbf{x}) = v("Robbery"|\mathbf{x})$, no matter how we randomly mask the input \mathbf{x} .

 $\begin{cases} threatened, took, smartphone \} \subseteq N, \text{ and the interaction makes a numerical effect } I(S) \text{ that boosts} \\ \text{the confidence of inferring the judgment of "robbery." Besides, Zhou et al. (2024) demonstrated that the complexity of interactions directly determined the generalization power of a DNN. \end{cases}$

Despite above achievements, previous studies have pointed out that the next breakthrough point is
to examine the correctness of the detailed decision-making logic used by the LLM, which have not
been explored yet (Deng et al., 2024b; Li & Zhang, 2023; Cheng et al., 2024; Ren et al., 2024; Chen
et al., 2024; Zhou et al., 2024).

076 In this paper, we extract all interactions that determine a legal LLM's confidence score of the true 077 judgment, and we evaluate the alignment between the extracted interactions and human cognition of the legal case. To this end, we categorize all input tokens involved in the interactions into three types, *i.e.*, the *relevant*, *irrelevant*, and *forbidden*² tokens, based on the ground-truth relevance to the judg-079 ment. This enables us to distinguish reliable interaction effects and unreliable interaction effects. For example, as Figure 1 shows, the legal LLM makes the judgment of "robbery" on Andy who 081 takes Bob's smartphone under threat. In this way, AND interactions involving "threatened," "took," and "smartphone" are supposed to be the correct reason for the judgment, thereby being identified 083 as reliable interaction effects. In comparison, the OR interaction between "June 1" and "angrily" 084 incorrectly attributes the judgment of "robbery" to the unreliable sentimental token "angrily." It is 085 because we should use the real action "threatened" to make the judgment, rather than the sentimental token "angrily." The unreliable interaction also includes the AND interaction between "struck" 087 and "threatened," which incorrectly attributes the judgment on Andy to the forbidden token "struck," *i.e.*, an action **not** taken by Andy.

In this way, we design new metrics based on these interactions to quantify the ratio of reliable interaction effects and that of unreliable ones used by the LLM to generate the target judgement, so as to evaluate the alignment between the LLM's logic and human cognition.

The contributions of this paper can be summarized as follows. We propose to utilize interactionbased explanations to evaluate the correctness of decision-making logic encoded by a LLM. We design new metrics to quantify reliable and unreliable interaction effects *w.r.t.* their alignment with human cognition of the judgment. Experiments on both English legal LLM and Chinese legal LLM show that both LLMs used a significant number of incorrect interactions for inference, although these LLMs all exhibited high accuracy in judgment prediction.

099 100

101

102 103

104

2 ALIGNMENT BETWEEN THE LLM AND HUMAN COGNITION

2.1 PRELIMIARIES: INTERACTIONS

Although there is no widely-accepted definition of concepts, which is an interdisciplinary issue across cognitive science, neuroscience, artificial intelligence, and mathematics, the theory of interactions has shown promise in explaining the primitive inference patterns encoded by the DNN. A

¹⁰⁵ 106 107

 $^{^{2}}$ The forbidden tokens are usually informative tokens but should not be used for judgments, *e.g.*, tokens of criminal actions that are **not** taken by the defendant.

series of properties (Li & Zhang, 2023; Ren et al., 2023a; 2024) have been proposed as mathematical
 guarantees for the faithfulness of the interaction-based explanations.

Definition of AND-OR interactions. Given an input sample $\mathbf{x} = [x_1, x_2, \dots, x_n]^\mathsf{T}$ with n input variables indexed by $N = \{1, 2, \dots, n\}$, where each input variable can represent a token, a word, or a phrase/short sentence. Then, let $v(\mathbf{x}) \in \mathbb{R}$ denote the *scalar* confidence of generating the target output. For example, the target output can be set to a sequence of m ground-truth tokens $[y_1, y_2, \dots, y_m]$ generated by the LLM. In this way, the scalar confidence of language generation $v(\mathbf{x})$ can be defined as follows.

116 117

130

131

 $v(\mathbf{x}) \stackrel{\text{def}}{=} \sum_{t=1}^{T} \log \frac{p(y = y_t | \mathbf{x}, \mathbf{Y}_t^{\text{previous}})}{1 - p(y = y_t | \mathbf{x}, \mathbf{Y}_t^{\text{previous}})}$ (1)

where $\mathbf{Y}_{t}^{\text{previous}} \stackrel{\text{def}}{=} [y_1, y_2, \cdots, y_{t-1}]^{\mathsf{T}}$ represents the sequence of the previous (t-1) tokens before generating the *t*-th token. $p(y = y_t | \mathbf{x}, \mathbf{Y}_t^{\text{previous}})$ denotes the probability of generating the *t*-th token, given the input sentence \mathbf{x} and the previous (t-1) tokens. In particular, $\mathbf{Y}_1^{\text{previous}} = []$.

To explain the inference patterns behind the confidence score $v(\mathbf{x})$, Ren et al. (2024); Shen et al. (2023) show that an LLM usually encodes a set of interactions between input variables (tokens or phrases) to compute $v(\mathbf{x})$. There are two types of interactions, *i.e.*, the AND interaction and the OR interaction. Each AND interaction and each OR interaction *w.r.t.* $S \subseteq N, S \neq \emptyset$ have specific numerical effects $I_{and}(S|\mathbf{x})$ and $I_{or}(S|\mathbf{x})$ to the network output, respectively, which are computed as follows.

$$I_{\text{and}}(S|\mathbf{x}) \stackrel{\text{def}}{=} \sum_{T \subseteq S} (-1)^{|S| - |T|} v_{\text{and}}(\mathbf{x}_T), \quad I_{\text{or}}(S|\mathbf{x}) \stackrel{\text{def}}{=} -\sum_{T \subseteq S} (-1)^{|S| - |T|} v_{\text{or}}(\mathbf{x}_{N \setminus T})$$
(2)

where \mathbf{x}_T denotes the masked sample³, where all embeddings of input variables in $N \setminus T$ are masked. $v(\mathbf{x}_T) \in \mathbb{R}$ denotes the confidence score of generating the *m* tokens $[y_1, y_2, \cdots, y_m]$ given the masked sample \mathbf{x}_T . $v(\mathbf{x}_T)$ is decomposed into the component for AND interactions $v_{and}(\mathbf{x}_T) = 0.5v(\mathbf{x}_T) + \gamma_T$ and the component for OR interactions $v_{or}(\mathbf{x}_T) = 0.5v(\mathbf{x}_T) - \gamma_T$, subject to $v_{and}(\mathbf{x}_T) + v_{or}(\mathbf{x}_T) = v(\mathbf{x}_T)$.

Extracting AND-OR interactions. According to Equation (2), the extraction of interactions is 138 implemented by learning parameters $\{\gamma_T\}$. We follow (Zhou et al., 2024) to learn parameters 139 $\{\gamma_T | T \subseteq N\}$, and extract the sparest (the simplest) AND-OR interaction explanation via the 140 LASSO-like loss, *i.e.*, $\min_{\{\gamma_T\}} \sum_{S \subseteq N, S \neq \emptyset} [|I_{and}(S|\mathbf{x})| + |I_{or}(S|\mathbf{x})|]$. In this way, we exhaus-141 tively compute interaction effects $I_{and}(S|\mathbf{x})$ and $I_{or}(S|\mathbf{x})$ for all $(2^n - 1)$ non-empty combinations 142 $\emptyset \neq S \subseteq N$. Ren et al. (2024) have proven that most interactions have almost zero effects 143 $I_{and/or}(S|\mathbf{x})$, and an LLM usually activates only 100-200 AND-OR interactions with salient 144 effects. These salient interactions are taken as the AND-OR logic really encoded by the LLM. 145

Algorithm 1 in the appendix shows the pseudo-code of extracting AND-OR interactions.

147 Why do AND-OR interactions faithfully explain the logic encoded by the LLM? Lots of theo-148 retical achievements ranging from (Harsanyi, 1963) to (Li & Zhang, 2023; Ren et al., 2023a; 2024) 149 have proven several properties to guarantee that the AND-OR interactions faithfully represent the 150 **AND-OR logic encoded by the LLM.** According to Theorem 1, let $h(\cdot)$ denote a surrogate logical 151 model constructed based on non-zero interactions. As Figure 6 shows, it is proven that this surrogate 152 logical model $h(\cdot)$ can accurately fit the confidence scores of the LLM $v(\cdot)$ on all 2^n masked samples $\{\mathbf{x}_T | T \subseteq N\}$, *i.e.*, $\forall T \subseteq N, v(\mathbf{x}_T) = h(\mathbf{x}_T)$, no matter how we randomly mask the input sample \mathbf{x} 153 in 2^n different masking states $T \subseteq N$. This property is termed *universal-matching property*. 154

Theorem 1 (Universal matching property, proof in Appendix B) Given an input sample \mathbf{x} , the network output score $v(\mathbf{x}_T) \in \mathbb{R}$ on each masked sample $\{\mathbf{x}_T | T \subseteq N\}$ can be well matched by a surrogate logical model $h(\mathbf{x}_T)$ on each masked sample $\{\mathbf{x}_T | T \subseteq N\}$. The surrogate logical model $h(\mathbf{x}_T)$ uses the sum of AND interactions and OR interactions to accurately fit the network output

160 161

155

³To obtain the masked sample \mathbf{x}_T , we mask the embedding of each input variable $i \in N \setminus T$ with the baseline value b_i to represent its masked state. Please see Appendix G.4 for details.

191 192 193

196 197

$$score v(\mathbf{x}_{T}).$$

$$\forall T \subseteq N, v(\mathbf{x}_{T}) = h(\mathbf{x}_{T}).$$

$$\forall T \subseteq N, v(\mathbf{x}_{T}) = h(\mathbf{x}_{T}).$$

$$h(\mathbf{x}_{T}) = v(\mathbf{x}_{\emptyset}) + \sum_{S \subseteq N, S \neq \emptyset} \mathbb{1}(\underset{\text{AND relation } S}{\mathbf{x}_{T} \text{ triggers}}) \cdot I_{\text{and}}(S|\mathbf{x}_{T}) + \mathbb{1}(\underset{\text{OR relation } S}{\mathbf{x}_{T} \text{ triggers}}) \cdot I_{\text{or}}(S|\mathbf{x}_{T})$$

$$= \underbrace{v(\mathbf{x}_{\emptyset}) + \sum_{S \subseteq T, S \neq \emptyset} I_{\text{and}}(S|\mathbf{x}_{T})}_{v_{\text{and}}(\mathbf{x}_{T})} + \underbrace{\sum_{S \subseteq N, S \cap T \neq \emptyset} I_{\text{or}}(S|\mathbf{x}_{T})}_{v_{\text{or}}(\mathbf{x}_{T})}$$

$$(3)$$

$$= \underbrace{v(\mathbf{x}_{\emptyset}) + \sum_{S \subseteq T, S \neq \emptyset} I_{\text{and}}(S|\mathbf{x}_{T})}_{v_{\text{and}}(\mathbf{x}_{T})} + \underbrace{\sum_{S \subseteq N, S \cap T \neq \emptyset} I_{\text{or}}(S|\mathbf{x}_{T})}_{v_{\text{or}}(\mathbf{x}_{T})}$$

171 Specifically, each non-zero AND interaction $I_{and}(S|\mathbf{x})$ represents the AND relationship between 172 all variables in S. For instance, consider an input sentence "the company is a legal person" in 173 a language generation task. The co-appearance of two words $S = \{legal, person\} \subseteq N$ forms a 174 specialized legal concept and contributes a numerical effect $I_{and}(S|\mathbf{x})$ to push the LLM's output 175 w.r.t. the legal entity. Exclusively inputting either word in S will not make such an effect.

Analogously, each non-zero OR interaction $I_{or}(S|\mathbf{x})$ indicates the OR relationship between all variables in S. For example, let us consider an input sentence "*he robbed and assaulted a passerby*". The presence of either word in $S = \{robbed, assaulted\}$ activates the OR relationship and contributes an effect $I_{or}(S|\mathbf{x})$ to push the LLM towards a guilty verdict.

Besides the *universal-matching property* in Theorem 1, the *sparsity property* of interactions is also proven (Ren et al., 2024). *I.e.*, most AND-OR interactions have almost zero effects, *i.e.*, $I(S|\mathbf{x}) \approx$ 0, which can be regarded as negligible noise patterns. Only a small set of interactions, denoted by $\Omega = \{S \subseteq N : |I(S|\mathbf{x})| > \tau\}$, where τ is a scalar threshold, have considerable effects. Therefore, Lemma 1 shows that the surrogate logical model $h(\cdot)$ on all 2^n masked samples $\{\mathbf{x}_T | T \subseteq N\}$ usually can be approximated by a small set of salient AND interactions Ω^{and} and salient OR interactions Ω^{or} , *s.t.*, $|\Omega^{\text{and}}|$, $|\Omega^{\text{or}}| \ll 2^n$.

Lemma 1 (Sparsity property, proof in Appendix C) The surrogate logical model $h(\mathbf{x}_T)$ on each randomly masked sample $\mathbf{x}_T, T \subseteq N$ mainly uses the sum of a small number of salient AND interactions and salient OR interactions to approximate the network output score $v(\mathbf{x}_T)$.

$$v(\mathbf{x}_T) = h(\mathbf{x}_T) \approx v(\mathbf{x}_{\emptyset}) + \sum_{S \in \Omega^{\text{and}}} \mathbb{1}\left(\frac{\mathbf{x}_T \text{ triggers}}{\text{AND relation } S}\right) \cdot I_{\text{and}}(S|\mathbf{x}_T) + \sum_{S \in \Omega^{\text{or}}} \mathbb{1}\left(\frac{\mathbf{x}_T \text{ triggers}}{\text{OR relation } S}\right) \cdot I_{\text{or}}(S|\mathbf{x}_T)$$
(4)

The above *universal-matching property* and *sparsity property* theoretically guarantee the faith fulness of the interaction-based explanation.

2.2 Relevant tokens, irrelevant tokens, and forbidden tokens

According to above achievements, we can take a small set of salient AND-OR interactions as the faithful explanation for the decision-making logic used by the legal LLM. Thus, in this subsection, we annotate the *relevant*, *irrelevant*, and *forbidden* tokens in the input legal case, in order to accurately identify the reliable and unreliable interactions encoded by the LLM (see Figure 1). Specifically, the set of all input variables N is partitioned into three mutually disjoint subsets, *i.e.*, the set of relevant tokens \mathcal{R} , the set of irrelevant tokens \mathcal{I} , and the set of forbidden tokens \mathcal{F} , subject to $\mathcal{R} \cup \mathcal{I} \cup \mathcal{F} = N$, with $\mathcal{R} \cap \mathcal{I} = \emptyset$, $\mathcal{R} \cap \mathcal{F} = \emptyset$, and $\mathcal{I} \cap \mathcal{F} = \emptyset$, according to human cognition.

Relevant tokens refer to tokens that are closely related to or serve as the direct reason for the judgment, according to human cognition. For instance, given an input legal case "on June 1, during a conflict on the street, Andy stabbed Bob with a knife, causing Bob's death,"¹ the legal LLM provides judgment "murder" for Andy. In this case, the input variables can be set as N ={[on June 1], [during a conflict], [on the street], [Andy stabbed Bob with a knife], [causing Bob's

death]}. $\mathcal{R} = \{[Andy \ stabbed \ Bob \ with \ a \ knife], [causing \ Bob's \ death]\}$ are the direct reason for the judgment, thereby being annotated as *relevant tokens*, where all tokens in the brackets [] are taken as a single input variable.

214 *Irrelevant tokens* refer to tokens that are not strongly related to or are not the direct reason for the 215 judgment, according to human cognition. For instance, in the above input legal case, the set of irrelevant tokens are annotated as $\mathcal{I} = \{[on June 1], [during a conflict], [on the street]\}$. For example, the input variable like "during a conflict" may influence Andy's behavior "Andy stabbed Bob with a knife," but it is the input variable "Andy stabbed Bob with a knife" that directly contributes to the legal judgment of "murder," rather than the input variable "during a conflict."

Forbidden tokens are usually common tokens
widely used in legal cases, but the use of forbidden tokens may lead to significant incorrect logic.
For instance, in a legal case involving multiple
individuals, such as "Andy assaulted Bob on the
head, causing minor injuries. Charlie stabbed

235 236 237

238 239



Figure 2: Sparsity of interactions. We show the strength of different AND-OR interactions $|I(S|\mathbf{x})|$ extracted from different samples in a descending order. Only about 0.5% interactions had salient effects.

Bob with a knife, causing Bob's death,"¹ the legal LLM assigns the judgment of "assault" to Andy. Let the set of all input variables be $N = \{[Andy assaulted Bob on the head], [causing$ $minor injuries], [Charlie stabbed Bob with a knife], [causing Bob's death]\}. Although the input$ variables "Charlie stabbed Bob with a knife" and "causing Bob's death" are naturally all represent crucial facts for judgement, they should not influence the judgment for Andy, because thesewords describe the actions of Charlie, not actions of Andy. Therefore, these input variables are $categorized as forbidden tokens, <math>\mathcal{F} = \{[Charlie stabbed Bob with a knife], [causing Bob's death]\}.$

2.3 Reliable and unreliable interaction effects

240 The categorization of *relevant*, *irrelevant*, and *forbidden* tokens enables us to disentangle the reliable 241 and unreliable decision-making logic used by a legal LLM. As introduced in Section 2.1, we use 242 interactions as the decision-making logic encoded by a legal LLM. Thus, in this subsection, we 243 decompose the overall interaction effects in Equation (2) into reliable and unreliable interaction effects. Reliable interaction effects are interaction effects that align with human cognition, which 244 usually contain relevant tokens and exclude forbidden tokens. In contrast, unreliable interaction 245 effects are interaction effects that do not match human cognition, which are attributed to irrelevant 246 or forbidden tokens. 247

248 Visualization of AND-OR interactions. Before defining reliable and unreliable interaction ef-249 fects, let us first visualize the AND-OR interactions extracted from two legal LLMs, SaulLM-7B-Instruct (Colombo et al., 2024) and BAI-Law-13B (Institute, 2023). SaulLM-7B-Instruct was an 250 English legal LLM, trained on a corpus of over 30 billion English legal tokens. BAI-Law-13B was 251 a Chinese legal LLM, fine-tuned on Chinese legal corpora. We evaluated the legal LLMs on the 252 CAIL2018 dataset (Xiao et al., 2018)⁵, just like how (Feng et al., 2022; Fei et al., 2023) did. Fig-253 ure 2 shows the sparsity of interactions extracted from the legal LLMs. Interaction strength $|I(S|\mathbf{x})|$ 254 of all AND-OR interactions extracted from all legal cases were shown in a descending order. We 255 found that most of the interactions had negligible effect. 256

Figure 1 further provides an example of using AND-OR interactions to explain the decision-making logic of a legal LLM. The legal LLM correctly attributes the judgment of "*robbery*" to interactions involving the tokens "*took*," "*smartphone*," and "*threatened*." However, the legal LLM also uses the irrelevant tokens ("*angrily*" and "*June 1*"), and the forbidden tokens ("*struck*" and "*on the head*") to compute the confidence score of the judgment of "*robbery*," which obviously represents incorrect decision-making logic.

In this way, we define *reliable* and *unreliable* interaction effects for AND and OR interactions, respectively, as follows.

For AND interactions. Because the AND interaction $I_{and}(S|\mathbf{x})$ is activated only when all input variables (tokens or phrases) in S are present in the input legal case, the reliable interaction effect for AND interaction $I_{and}^{reliable}(S|\mathbf{x})$ w.r.t. S must include relevant tokens in \mathcal{R} , *i.e.*, $S \cap \mathcal{R} \neq \emptyset$, and completely exclude forbidden tokens in \mathcal{F} , *i.e.*, $S \cap \mathcal{F} = \emptyset$. Otherwise, if S contains any forbidden tokens in \mathcal{F} , or if S does not contains any relevant tokens in \mathcal{R} , then the AND interaction $I_{and}(S|\mathbf{x})$ represents an incorrect logic for judgment. In this way, the reliable and unreliable AND interaction effects *w.r.t. S* can be computed as follows.

272 273 274

275

276

277

278

279

281

283 284

285

286

287 288

289

290

291

292 293

294

if
$$S \cap \mathcal{F} = \emptyset, S \cap \mathcal{R} \neq \emptyset$$
 then $I_{and}^{reliable}(S|\mathbf{x}) = I_{and}(S|\mathbf{x}), \quad I_{and}^{unreliable}(S|\mathbf{x}) = 0$
otherwise, $I_{and}^{reliable}(S|\mathbf{x}) = 0, \quad I_{and}^{unreliable}(S|\mathbf{x}) = I_{and}(S|\mathbf{x})$ (5)

For OR interactions. The OR interaction $I_{or}(S|\mathbf{x})$ affects the LLM's output when any input variable (token or phrase) in S appears in the input legal case. Therefore, we can define the reliable effect $I_{or}^{\text{reliable}}(S|\mathbf{x})$ as the numerical component in $I_{or}(S|\mathbf{x})$ allocated to relevant input variables in $S \cap \mathcal{R}$. To this end, just like in (Deng et al., 2024b), we uniformly allocate the OR interaction effects to all input variables in S. The reliable and unreliable interactions effects are those allocated to relevant variables, respectively.

$$\forall S \subseteq N, S \neq \emptyset, I_{\text{or}}^{\text{reliable}}(S|\mathbf{x}) = \frac{|S \cap \mathcal{R}|}{|S|} \cdot I_{\text{or}}(S|\mathbf{x}), I_{\text{or}}^{\text{unreliable}}(S|\mathbf{x}) = \left(1 - \frac{|S \cap \mathcal{R}|}{|S|}\right) \cdot I_{\text{or}}(S|\mathbf{x}) \quad (6)$$

2.4 EVALUATION METRICS

In this subsection, we design a set of metrics to evaluate the alignment quality between the interactions encoded by the LLM and human cognition.

Ratio of reliable interaction effects. Definition 1 introduces the ratio of reliable interaction effects that align with human cognition to all salient interaction effects. Here, we focus on the small number of salient interactions in Ω^{and} and Ω^{or} , rather than conduct evaluation on interactions effects of all 2^n subsets $S \subseteq N$. This is because salient interactions can be taken as primitive decision-making logic of an LLM, while all other interactions have negligible effects and represent noise patterns.

Definition 1 (Ratio of reliable interaction effects) Given an LLM, the ratio of reliable interaction effects to all salient interaction effects s^{reliable} is computed as follows.

$$s^{\text{reliable}} = \frac{\sum_{\Omega^{\text{and}}} |I_{\text{and}}^{\text{reliable}}(S|\mathbf{x})| + \sum_{\Omega^{\text{or}}} |I_{\text{or}}^{\text{reliable}}(S|\mathbf{x})|}{\sum_{\Omega^{\text{and}}} |I_{\text{and}}(S|\mathbf{x})| + \sum_{\Omega^{\text{or}}} |I_{\text{or}}(S|\mathbf{x})|}$$
(7)

A larger value of $s^{\text{reliable}} \in [0, 1]$ indicates that a higher proportion of interaction effects align with human cognition.

Interaction distribution over different orders. Zhou et al. (2024) have found that the low-order 302 interactions usually exhibit stronger generalization power⁴ than high-order interactions. *I.e.*, low-303 order interactions learned from training samples are more likely to be transferred to (appear in) 304 testing samples. Please see Appendix E for the definition and quantification of the generalization 305 power of interactions over different orders. Specifically, the order is defined as the number of input 306 variables in S, *i.e.*, order(S) = |S|. In general, high-order interactions (complex interactions) 307 between a large number of input variables are usually less generalizable⁴ than low-order (simple) 308 interactions. 309

Therefore, we utilize the distribution of interactions over different orders as another metric, which evaluates the generalization power of the decision-making logic used by the LLM. Specifically, we use *Salient*⁺(*o*) = $\sum_{op\in\{and,or\}} \sum_{S\in\Omega^{op},|S|=o} \max(0, I_{op}(S|\mathbf{x}))$ to quantify the overall strength of positive salient interactions, and use *Salient*⁻(*o*) = $\sum_{op\in\{and,or\}} \sum_{S\in\Omega^{op},|S|=o} \min(0, I_{op}(S|\mathbf{x}))$ to quantify the overall strength of negative salient interactions. A well-trained legal LLM tends to model loworder interactions, while an over-fitted LLM (potentially due to insufficient data or inadequate data cleaning) usually relies more on high-order interactions.

Ratio of reliable interaction effects of each order. We categorize all salient interaction effects by
 their orders, so that for all salient interactions of each *o*-th order, we can compute the ratio of reliable
 interaction effects.

³²⁰ ⁴The generalization power of an interaction is defined as the transferability of this interaction from training ³²¹samples to test samples. Specifically, if an interaction pattern $S \subseteq N$ frequently occurs in the training set, but ³²²rarely appears in the test set, then the interaction pattern S exhibits low generalization power. Conversely, if an ³²³interaction pattern S consistently appears in both the training and test sets, it demonstrates high generalization ³²⁶power. Please see Appendix E for details.

324 Definition 2 (Ratio of reliable interaction effects of each order) The ratio of reliable interaction 325 effects to all positive salient interaction effects of the o-th order is measured by $s_0^{\text{reliable},+}$ = 326 $Reliable^+(o)$ Similarly, the ratio of reliable interaction effects to all negative salient inter- $\overline{Salient^+(o)+\epsilon}$ 327 action effects of the o-th order is measured by $s_o^{\text{reliable},-} = \frac{|\text{Reliable}^-(o)|}{|\text{Salient}^-(o)|+\epsilon}$. Reliable⁺(o) = 328 $\sum_{op \in \{and, or\}} \sum_{S \in \Omega^{op}, |S|=o} \max(0, I_{op}^{\text{reliable}}(S|\mathbf{x})) \text{ represents the overall strength of positive reliable in-$ 329 teractions of the o-th order, and Reliable⁻(o) = $\sum_{\text{op} \in \{\text{and}, \text{or}\}} \sum_{S \in \Omega^{\text{op}}, |S|=o} \min(0, I_{\text{op}}^{\text{reliable}}(S|\mathbf{x}))$ repre-330 331 sents the overall strength of negative reliable interactions of the o-th order. ϵ is a small constant to 332 avoid dividing 0. 333

According to the findings in (Zhou et al., 2024), low-order interactions generally represent stable patterns that are frequently used across a large number of legal cases. Thus, if a considerable ratio of low-order interactions contain unreliable effects, it suggests that training data may have a clear bias, which makes the LLM stably learns unreliable interactions. In comparison, since high-order interactions typically exhibit poor generalization power, unreliable effects in high-order interactions are usually attributed to the memorization of hard/outlier samples. Consequently, low-order unreliable interactions are are mainly owing to stable bias in the training data, while high-order unreliable interactions often indicates that the LLM learns outlier features.

341 342 343

344

3 EXPERIMENT

In this section, we conducted experiments to evaluate the alignment quality between the decision making logic of the legal LLM and human cognition. In this way, we identified potential represen tation flaws behind the seemingly correct language generation results of legal LLMs.

348 We applied two off-the-shelf legal LLMs, SaulLM-7B-Instruct (Colombo et al., 2024) and BAI-Law-349 13B (Institute, 2023), which were trained for legal judgment prediction on English legal corpora and 350 Chinese legal corpora, respectively. Appendix F shows the accuracy of these LLMs. Given an input 351 legal case, the LLM predicted the judgment result based on the fact descriptions of the legal case. We explained judgments made on legal cases in the CAIL2018 dataset (Xiao et al., 2018), which 352 contained 2.6 million Chinese legal cases, for both legal LLMs⁵. Figure 6 shows the universal-353 matching property of the extracted interactions, *i.e.*, when we randomly masked input variables in 354 the legal case, we could always use the interactions to accurately match the real confidence scores 355 of the judgment estimated by the LLM. 356

To simplify the explanation and avoid ambiguity, we only explained the decision-making logic on legal cases, which were correctly judged by the LLM. For each input legal case, we manually selected some informative tokens or phrases as input variables. Some tokens or phrases were annotated as relevant tokens in \mathcal{R} , while others were identified as irrelevant tokens in \mathcal{I} . It was ensured that the removal of all input variables would substantially change the legal judgment result.

We extracted AND-OR interactions that determined the confidence score $v(\mathbf{x})$ of generating judgment results with a sequence of tokens, according to Equation (1). To accurately identify and analyze potential representation flaws from these interactions, in this paper, we mainly focused on potential representation flaws *w.r.t.* legal judgments in the following three types, *i.e.*, (1) judgments influenced by unreliable sentimental tokens, (2) judgments affected by incorrect entity matching, and (3) judgments biased by discrimination in occupation.

Problem 1: making judgments based on unreliable sentimental tokens. We observed that although legal LLMs achieved relatively high accuracy in predicting judgment results (see Appendix F), a considerable number of interactions contributing to the confidence score $v(\mathbf{x})$ were attributed to semantically irrelevant or unreliable sentimental tokens. The legal LLM was supposed to focus more on real criminal actions, than unreliable sentimental tokens behind the actions, when criminal actions had been given. We believed these indicated potential representation flaws behind

 ⁵To ensure a fair comparison, we conducted experiments using the same dataset across both legal LLMs.
 For the BAI-Law-13B model, which was a Chinese legal LLM, we directly analyzed the Chinese legal cases
 from the CAIL2018 dataset. For the SaulLM-7B-Instruct model, which was an English legal LLM, we translated these Chinese legal cases into English and performed the analysis on the translated cases, to enable fair comparisons. Please see Appendix G.5 for details.

the seemingly correct legal judgments produced by legal LLMs. To evaluate the impact of unreliable sentimental tokens on both the SaulLM-7B-Instruct and BAI-Law-13B models, we annotated tokens that served as the direct reason for the judgment as relevant tokens in \mathcal{R} , and those that were not the direct reason for the judgment as irrelevant tokens in \mathcal{I} , *e.g.*, semantically irrelevant tokens and unreliable sentimental tokens behind real criminal actions.

Figure 3 shows the legal case, which showed Andy had a conflict with Bob and attacked Bob, com-384 mitting an assault. In this case, tokens like "began to," "causing," and sentiment-driven tokens such 385 as "dissatisfaction" in \mathcal{I} were irrelevant to the judgment result, according to human cognition, be-386 cause unreliable sentimental tokens only served as explanations for criminal actions. Thus, once 387 an actual action had been taken, the unreliable sentimental tokens were supposed to make minimal 388 conditional contributions to the legal judgment result. The judgment should be based exclusively on tokens such as "fight chaotically," "threw a punch," and "fall into a coma," which were annotated 389 as relevant tokens in \mathcal{R} . We found that some decision-making logic encoded by the SaulLM-7B-390 Instruct model aligned well with human cognition, *i.e.*, identifying reliable interactions containing 391 relevant tokens as the most salient interactions. However, this model also modeled lots of unreliable 392 interactions as salient interactions, such as interactions containing irrelevant tokens "dissatisfaction" 393 and "anger," which revealed potential flaws in its decision-making logic. 394

In comparison, we evaluated the above legal case on the BAI-Law-13B model, as shown in Figure 3. 395 The SaulLM-7B-Instruct model exhibited a reliable interaction ratio of $s^{\text{reliable}} = 71.5\%$, while 396 the BAI-Law-13B model encoded a lower ratio of reliable interaction effects, $s^{\text{reliable}} = 61.2\%$. 397 The BAI-Law-13B model encoded about 10% less reliable interactions, and used $s^{\text{unreliable}}$ = 398 $1 - s^{\text{reliable}} = 38.8\%$ unreliable interaction effects to compute the confidence score $v(\mathbf{x})$. For 399 example, reliable interactions encoded by the BAI-Law-13B model included the AND interaction 400 $S = \{$ "threw a punch" $\}$, which contributed the highest interaction effect 0.34. The unreliable interactions included the AND interaction $S = {\text{"anger"}}$, which contributed 0.03. The unreliable 401 402 sentimental token should not be used to determine the judgment, when the action "threw a punch" 403 caused by "anger" had been given as a more direct reason. Additional examples of making judg-404 ments based on unreliable sentimental tokens are provided in Appendix G.1.

405 Problem 2: making judgments based on incorrect entity matching. Despite the high accuracy 406 of legal LLMs in predicting judgment results, we found that a considerable ratio of the confidence 407 score $v(\mathbf{x})$ was mistakenly attributed to interactions on criminal actions made by incorrect entities. 408 In other words, the LLM mistakenly used the criminal action of a person (entity) to make judgment 409 on another unrelated person (entity). To evaluate the impact of such incorrect entity matching on 410 both the SaulLM-7B-Instruct and BAI-Law-13B models, we annotated tokens for criminal actions 411 of unrelated entities as the *forbidden tokens* in \mathcal{F} . These forbidden tokens should not influence the 412 judgment for the unrelated entity.

413 Figure 4 illustrates the test of the SaulLM-7B-Instruct model on the legal case, which showed Andy 414 bit Charlie, committing an assault, and then Bob hit Charlie with a shovel, leading to murder. Be-415 cause tokens such as "hit," "with a shovel," "injuring," and "death" described Bob's actions and 416 consequences without a direct relationship with Andy. Thus, these tokens were annotated as for-417 bidden tokens in \mathcal{F} . However, we observed that although the SaulLM-7B-Instruct model had used $s^{\text{reliable}} = 21.5\%$ reliable interactions between relevant tokens, such as "bit" and "slightly injured," 418 it also modeled a significant number of unreliable interactions containing forbidden tokens "death" 419 and "with a shovel." However, if we removed these two forbidden tokens for criminal actions of 420 Bob, then the confidence of the judgment of Andy would be significantly affected. This was an 421 obvious representation flaw of the SaulLM-7B-Instruct model. 422

In comparison, given the same legal case, the BAI-LAW-13B model encoded a ratio of $s^{\text{reliable}} =$ 423 22.6% reliable interaction effects, which was a bit higher than a ratio of $s^{\text{reliable}} = 21.5\%$ reliable 424 interactions effects encoded by the SaulLM-7B-Instruct model. In this case, both models primarily 425 relied on unreliable interactions, including forbidden tokens that related to Bob's criminal actions, 426 to make judgment on Andy. For example, the SaulLM-7B-Instruct model used the AND interaction 427 w.r.t. the unrelated action $S = \{ "with a shovel" \}$ to contribute 0.93, and the BAI-Law-13B model 428 used the AND interaction $S = \{ (death) \}$ to contribute -0.43. This suggested that both legal LLMs 429 handled judgment-related tokens in a local manner, without accurately matching criminal actions 430 with entities. Additional examples of making judgments based on incorrect entity matching are 431 provided in Appendix G.2.

Problem 3: discrimination in occupation may affect judgments. We found that the legal LLM usually used interactions on the occupation information to compute the confidence score $v(\mathbf{x})$. This would lead to a significant occupation bias. More interestingly, we discovered that when we replaced the current occupation with another occupation, the interaction containing the occupation token would be significant changed. This indicates a common bias problem, because similar bias may also happen on other attributes (*e.g.*, age, gender, education level, and marital status).

438 Figure 5 shows the test of the SaulLM-7B-Instruct model on the legal case, in which Andy, the victim 439 with varying occupations, was robbed of his belongings by two suspicious men. First, we found that 440 the SaulLM-7B-Instruct model encoded interactions with the occupation tokens "a judge," which boosted the confidence of the judgment "robbery." More interestingly, if we substituted the occupa-441 442 tion tokens "a judge" to "a volunteer," the interaction between the occupation "a volunteer," "a day's work," and "belongings" decreased from 0.22 to 0.06. This was an important factor that changed the 443 judgment from "robbery" to "not mentioned." However, if we replaced "a judge" with law-related 444 occupations, such as "a lawyer" and "a policeman," the judgment remained "robbery." Besides, the 445 occupation "a programmer" changed the judgment to "not mentioned." Please see Appendix G.3 for 446 numerical effects of all these occupations. This suggested that the legal LLM sometimes had con-447 siderable occupation bias. In comparison, we evaluated the same legal case on the BAI-Law-13B 448 model, as shown in Appendix G.3. Compared to the SaulLM-7B-Instruct model that encoded a ratio 449 $s^{\text{reliable}} = 81.4\% - 84.0\%$ of reliable interaction effects w.r.t. different occupations, the BAI-Law-450 13B model encoded a ratio $s^{\text{reliable}} = 78.9\%$ -87.1% of reliable interaction effects. This indicated 451 that both legal LLMs tended to use specific occupational tokens for judgment, instead of correctly 452 analyzing the decision-making logic behind legal judgements. Additional examples of judgments 453 biased by the occupation are provided in Appendix G.3.

454 **Representation quality of legal LLMs.** Figures 3 and 4 compare the interaction effects of different 455 orders extracted from the SaulLM-7B-Instruct model and the BAI-Law-13B model. We observed 456 that, in both the legal case influenced by unreliable sentimental tokens, and the legal case affected 457 by incorrect entity matching, the BAI-Law-13B model encoded higher order interactions than the 458 SaulLM-7B-Instruct model. This indicated that feature representations of the BAI-Law-13B model 459 was more complex and less generalizable than than those of the SaulLM-7B-Instruct model. In addition, in the legal case that judgments affected by incorrect entity matching, the BAI-Law-13B 460 model encoded a significant number of interactions with negative effects. This suggested that many 461 interactions encoded by the BAI-Law-13B model showed conflicting effects, which was also a sign 462 of over-fitting of the LLM. Tables 1 and 2 in the appendix further show the average ratios of the 463 reliable interaction effects $s_o^{\text{reliable},+}$ and $s_o^{\text{reliable},-}$ for each order o on both LLMs. Experimental 464 results show that while the BAI-Law-13B model encoded more low-order reliable interaction effects, 465 it also encoded more high-order unreliable interaction effects than the SaulLM-7B-Instruct model. 466

467

4 CONCLUSION

468 469

In this paper, we have proposed a method to evaluate the correctness of the detailed decision-making 470 logic of an LLM. The sparsity property and the universal matching property of interactions provide 471 direct mathematical supports for the faithfulness of the interaction-based explanation. Thus, in this 472 paper, we have designed two new metrics to quantify reliable and unreliable interaction effects, 473 according to their alignment with human cognition. Experiments showed that the legal LLMs often 474 relied on a considerable number of problematic interactions to make judgments, even when the 475 judgement prediction was correct. The evaluation of the alignment between the decision-making 476 logic of LLMs and human cognition also contributes to other real applications. For example, it may 477 assist in debugging the hallucination problems, and identifying potential bias behind the language 478 generation results of LLMs.

Limitations. Our analysis does not assess the correctness of the numerical scores for interactions, as these scores are often determined by many factors. Positive interactions typically indicate logics that contribute positively to the judgments, while negative interactions may also be intended for other possible correct judgments. Besides, the evaluation based on relevant, irrelavant, and forbidden tokens is only one of the conditions for reliable interactions, and reliable interactions may not always be correct. Nevertheless, this paper presents a precedent for evaluating the correctness of decision-making logic of LLMs.



Figure 3: Visualization of judgments influenced by unreliable sentimental tokens. (a) A number of irrelevant tokens were annotated in the legal case, including unreliable sentimental tokens. Criminal actions were annotated as relevant tokens. We also translated the legal case to English as the input of the SaulLM-7B-Instruct model. (b) Judgements predicted by the two legal LLMs, which were both correct according to laws of the two countries. (c,d) We quantified the reliable and unreliable interaction effects of different orders.



Figure 4: Visualization of judgments affected by incorrect entity matching. (a) A number of irrelevant tokens were annotated in the legal case, including the time and actions that were not the direct reason for the judgment. Criminal actions of the defendant were annotated as relevant tokens. Criminal actions of the unrelated person were annotated as forbidden tokens. (b) Judgements predicted by the two legal LLMs, which were both correct according to laws of the two countries. (c,d) We measured the reliable and unreliable interaction effects of different orders.



Figure 5: Visualization of judgments biased by discrimination in occupation. (a) A number of irrelevant tokens were annotated in the legal case, including the occupation, time and actions that are not the direct reason for the judgment. Criminal actions of the defendant were annotated as relevant tokens. (b) The SaulLM-7B-Instruct model predicted the judgment based on the legal case with different occupations, respectively. (c,d) We measured the reliable and unreliable interaction effects of different orders. When the occupation was set to "*a judge*," the LLM used 81% reliable interaction effects. In comparison, when the occupation was set to "*a volunteer*," the LLM encoded 84% reliable interaction effects.

538

539

540 REFERENCES

542 543 544 545 546 547	Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu, and Pascale Fung. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. In <i>Proceedings of the 13th International Joint Conference on Natural Language Processing and</i> <i>the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics</i> (Volume 1: Long Papers), pp. 675–718. Association for Computational Linguistics, 2023.
548 549 550	Lu Chen, Siyu Lou, Benhao Huang, and Quanshi Zhang. Defining and extracting generalizable interaction primitives from dnns. <i>International Conference on Learning Representations</i> , 2024.
551 552	Xu Cheng, Lei Cheng, Zhaoran Peng, Yang Xu, Tian Han, and Quanshi Zhang. Layerwise change of knowledge in neural networks. In <i>International Conference on Machine Learning</i> , 2024.
553 554 555 556	 Pierre Colombo, Telmo Pessoa Pires, Malik Boudiaf, Dominic Culver, Rui Melo, Caio Corro, Andre F. T. Martins, Fabrizio Esposito, Vera Lúcia Raposo, Sofia Morgado, and Michael Desa. Saullm-7b: A pioneering large language model for law. <i>arXiv preprint arXiv:2403.03883</i>, 2024.
557 558	Huiqi Deng, Qihan Ren, Hao Zhang, and Quanshi Zhang. Discovering and explaining the represen- tation bottleneck of dnns. In <i>International Conference on Learning Representations</i> , 2022.
559 560 561 562	Huiqi Deng, Na Zou, Mengnan Du, Weifu Chen, Guocan Feng, Ziwei Yang, Zheyang Li, and Quan- shi Zhang. Unifying fourteen post-hoc attribution methods with taylor interactions. <i>IEEE Trans-</i> <i>actions on Pattern Analysis and Machine Intelligence</i> , 2024a.
563 564 565 566	Huiqi Deng, Na Zou, Mengnan Du, Weifu Chen, Guocan Feng, Ziwei Yang, Zheyang Li, and Quan- shi Zhang. Unifying fourteen post-hoc attribution methods with taylor interactions. <i>IEEE Trans-</i> <i>actions on Pattern Analysis and Machine Intelligence</i> , 2024b.
567 568 569	Nouha Dziri, Andrea Madotto, Osmar Zaïane, and Avishek Joey Bose. Neural path hunter: Reducing hallucination in dialogue systems via path grounding. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , 2021.
570 571 572 573	Zhiwei Fei, Xiaoyu Shen, Dawei Zhu, Fengzhe Zhou, Zhuo Han, Songyang Zhang, Kai Chen, Zongwen Shen, and Jidong Ge. Lawbench: Benchmarking legal knowledge of large language models. In <i>arXiv preprint arXiv:2309.16289</i> , 2023.
574 575	Yi Feng, Chuanyi Li, and Vincent Ng. Legal judgment prediction: A survey of the state of the art. In <i>IJCAI</i> , pp. 5461–5469, 2022.
576 577 578	Katja Filippova. Controlled hallucinations: Learning to generate faithfully from noisy data. <i>Findings</i> of the Association for Computational Linguistics: EMNLP, 2020.
579 580 581	Sorin Grigorescu, Bogdan Trasnea, Tiberiu Cocias, and Gigel Macesanu. A survey of deep learning techniques for autonomous driving. <i>Journal of field robotics</i> , 37(3):362–386, 2020.
582 583 584	Neel Guha, Julian Nyarko, Daniel E. Ho, Christopher Re, and et al. Legalbench: A collaboratively built benchmark for measuring legal reasoning in large language models. In Advances in Neural Information Processing Systems (Track on Datasets and Benchmarks), 2023.
585 586 587	John C Harsanyi. A simplified bargaining model for the n-person cooperative game. <i>International Economic Review</i> , 4(2):194–220, 1963.
588 589 590 501	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. <i>International Conference on Learning Representations</i> , 2021.
592 593	Or Honovich, Roee Aharoni, Jonathan Herzig, Hagai Taitelbaum, Doron Kukliansy, Vered Cohen, Thomas Scialom, Idan Szpektor, Avinatan Hassidim, and Yossi Matias. True: Re-evaluating factual consistency evaluation. <i>arXiv preprint arXiv:2204.04991</i> , 2022.

- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong 595 Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. A survey on hallucination in large language 596 models: Principles, taxonomy, challenges, and open questions. arXiv preprint arXiv:2311.05232, 597 2023. 598 Yichong Huang, Xiachong Feng, Xiaocheng Feng, and Bing Qin. The factual inconsistency problem in abstractive text summarization: A survey. arXiv preprint arXiv:2104.14839, 2021. 600 601 Baiyulan Open AI Research Institute. Baiyulan open ai. 2023. URL https://baiyulan.org. 602 cn. 603 Jiaming Ji, Tianyi Qiu, Boyuan Chen, Borong Zhang, Hantao Lou, Kaile Wang, Yawen Duan, 604 Zhonghao He, Jiayi Zhou, Zhaowei Zhang, et al. Ai alignment: A comprehensive survey. arXiv 605 preprint arXiv:2310.19852, 2023a. 606 607 Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. ACM 608 Computing Surveys (CSUR), 55(12), 2023b. 609 610 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, 611 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 612 Mistral 7b. arXiv preprint arXiv:2310.06825, 2023. 613 Jan Leike, David Krueger, Tom Everitt, Miljan Martic, Vishal Maini, and Shane Legg. Scalable 614 agent alignment via reward modeling: a research direction. arXiv preprint arXiv:1811.07871, 615 2018. 616 617 Mingjie Li and Quanshi Zhang. Does a neural network really encode symbolic concept? In Inter-618 national Conference on Machine Learning, 2023. 619 Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulga: Measuring how models mimic human 620 falsehoods. In Proceedings of the 60th Annual Meeting of the Association for Computational 621 Linguistics (Volume 1: Long Papers), 2021. 622 623 Dongrui Liu, Huiqi Deng, Xu Cheng, Qihan Ren, Kangrui Wang, and Quanshi Zhang. Towards 624 the difficulty for a deep neural network to learn concepts of different complexities. Advances in Neural Information Processing Systems, 2024. 625 626 Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. Advances 627 in Neural Information Processing Systems, 30, 2017. 628 Potsawee Manakul, Adian Liusie, and Mark JF Gales. Selfcheckgpt: Zero-resource black-box hal-629 lucination detection for generative large language models. In Proceedings of the 2023 Conference 630 on Empirical Methods in Natural Language Processing, 2023. 631 632 Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. On faithfulness and factuality 633 in abstractive summarization. Proceedings of the 58th Annual Meeting of the Association for 634 Computational Linguistics, 2020. 635 Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke 636 Zettlemoyer, and Hannaneh Hajishirzi. Factscore: Fine-grained atomic evaluation of factual pre-637 cision in long form text generation. In Proceedings of the 2023 Conference on Empirical Methods 638 in Natural Language Processing, 2023. 639 640 OpenAI. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023. 641 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong 642 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to fol-643 low instructions with human feedback. Advances in neural information processing systems, 35: 644 27730-27744, 2022. 645 Jie Ren, Mingjie Li, Qirui Chen, Huiqi Deng, and Quanshi Zhang. Defining and quantifying the 646
- Jie Ren, Mingjie Li, Qirui Chen, Huiqi Deng, and Quanshi Zhang. Defining and quantifying the
 emergence of sparse concepts in dnns. *IEEE Conference on Computer Vision and Pattern Recognition*, 2023a.

- Jie Ren, Zhanpeng Zhou, Qirui Chen, and Quanshi Zhang. Can we faithfully represent masked states to compute shapley values on a dnn? *International Conference on Learning Representations*, 2023b.
- Qihan Ren, Huiqi Deng, Yunuo Chen, Siyu Lou, and Quanshi Zhang. Bayesian neural networks
 avoid encoding perturbation-sensitive and complex concepts. *International Conference on Machine Learning*, 2023c.
- Qihan Ren, Jiayang Gao, Wen Shen, and Quanshi Zhang. Where we have arrived in proving the
 emergence of sparse interaction primitives in dnns. In *International Conference on Learning Representations*, 2024.
- Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based local-ization. In *Proceedings of the IEEE international conference on computer vision*, 2017.
- Tianhao Shen, Renren Jin, Yufei Huang, Chuang Liu, Weilong Dong, Zishan Guo, Xinwei Wu,
 Yan Liu, and Deyi Xiong. Large language model alignment: A survey. *arXiv preprint arXiv:2309.15025*, 2023.
- Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In International conference on machine learning. PMLR, 2017.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- ⁶⁷¹ Cunxiang Wang, Sirui Cheng, Qipeng Guo, Yuanhao Yue, Bowen Ding, Zhikun Xu, Yidong Wang, Xiangkun Hu, Zheng Zhang, and Yue Zhang. Evaluating open-qa evaluation. *Advances in Neural Information Processing Systems*, 36, 2024.
- Yufei Wang, Wanjun Zhong, Liangyou Li, Fei Mi, Xingshan Zeng, Wenyong Huang, Lifeng Shang,
 Xin Jiang, and Qun Liu. Aligning large language models with human: A survey. *arXiv preprint arXiv:2307.12966*, 2023.
- Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. Jailbroken: How does llm safety training
 fail? In Advances in Neural Information Processing Systems, volume 36, 2023.
- Chaojun Xiao, Haoxi Zhong, Zhipeng Guo, Cunchao Tu, Zhiyuan Liu, Maosong Sun, Yansong Feng, Xianpei Han, Zhen Hu, Heng Wang, et al. Cail2018: A large-scale legal dataset for judgment prediction. *arXiv preprint arXiv:1807.02478*, 2018.
- Quanshi Zhang, Xin Wang, Jie Ren, Xu Cheng, Shuyun Lin, Yisen Wang, and Xiangming Zhu.
 Proving common mechanisms shared by twelve methods of boosting adversarial transferability.
 arXiv preprint arXiv:2207.11694, 2022.
- Huilin Zhou, Hao Zhang, Huiqi Deng, Dongrui Liu, Wen Shen, Shih-Han Chan, and Quanshi Zhang.
 Explaining generalization power of a dnn using interactive concepts. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 17105–17113, 2024.
- 691 692

658

- 693 694
- 695
- 696
- 697
- 698 699
- 700
- 701

702 A RELATED WORK

704 Factuality and hallucination problems. Factuality in LLMs refers to whether the language gen-705 eralization results of LLMs align with the verificable facts. This includes the ability of LLMs to 706 avoid producing misleading or incorrect information (*i.e.*, factual hallucination), and to effectively 707 generate factually accurate results. For instance, several studies have evaluated the correctness of 708 LLM-generated answers to specific questions (Lin et al., 2021; OpenAI, 2023; Wang et al., 2024). 709 Other works have standardized fact consistency tasks into binary labels, evaluating whether there were factual conflicts within the input text (Honovich et al., 2022). Min et al. (2023) further de-710 composed language generation results into "atomic" facts, and calculated the proportion of these 711 facts that aligned with a given knowledge source. Additionally, Manakul et al. (2023) introduced 712 a sampling-based method to verify whether LLMs generated factually consistent results, based on 713 the assumption that if an LLM had knowledge of a concept, then the sampled generation results 714 contained consistent factual information. 715

Hallucination in LLMs typically refers to generated content that is nonsensical or unfaithful to the 716 provided source input (Filippova, 2020; Maynez et al., 2020; Huang et al., 2023). Hallucinations are 717 generally categorized into two primary types, namely intrinsic and extrinsic hallucinations (Maynez 718 et al., 2020; Huang et al., 2021; Dziri et al., 2021; Ji et al., 2023b). Intrinsic hallucinations oc-719 cur when the generated results contradict the source content, while extrinsic hallucinations arise 720 when the generated results cannot be verified from the provided source. For instance, Bang et al. 721 (2023) found extrinsic hallucinations in ChatGPT's responses, including both untruthful and fac-722 tual hallucinations, whereas intrinsic hallucinations were rarely observed. OpenAI's latest model, 723 GPT-4 (OpenAI, 2023), has further reduced the model's tendency to hallucinate compared to prior 724 models such as ChatGPT.

725 Value alignment. Value alignment in LLMs aims to ensure LLMs behave in accordance with hu-726 man intentions and values (Leike et al., 2018; Wang et al., 2023; Ji et al., 2023a). Recent research 727 has focused on improving the ability of LLMs to comprehend instructions, thereby aligning their be-728 havior with human expectations. For instance, OpenAI proposed Supervised Fine-Tuning (SFT) for 729 LLMs, which involved using human-annotated instruction data. LLMs such as InstructGPT (Ouyang 730 et al., 2022) and ChatGPT, both of which employed this technique, have demonstrated significant 731 improvements in understanding human instructions. Ouyang et al. (2022); OpenAI (2023); Touvron et al. (2023) have incorporated the Reinforcement Learning from Human Feedback (RLHF) method 732 to further fine-tune LLMs, enhancing their alignment with human preferences (OpenAI, 2023). 733

734 Using interactions to faithfully explain DNNs. Ren et al. (2023a) first proposed to quantify in-735 teractions between input variables encoded by the DNN, to explain the knowledge in the DNN. Li 736 & Zhang (2023) discovered the discriminative power of interactions between input variables. Ren 737 et al. (2024) further proved that DNNs usually only encoded a small number of interactions. Futher-738 more, Deng et al. (2024a) found that different attribution scores estimated by fourteen attribution methods, including the Grad-CAM (Selvaraju et al., 2017), Integrated Gradients (Sundararajan et al., 739 2017), and Shapley value (Lundberg & Lee, 2017)) methods, could all be represented as a combi-740 nation of interactions. Besides, Zhang et al. (2022) used interactions to explain the mechanism of 741 different methods of boosting adversarial transferability. Ren et al. (2023b) used interactions to de-742 fine the optimal baseline value for computing Shapley values. Deng et al. (2022) found that for most 743 DNNs it was difficult to learn interactions with median number of input variables, and it was dis-744 covered that DNNs and Bayesian neural networks were unlikely to model complex interactions with 745 many input variables (Ren et al., 2023c; Liu et al., 2024). Zhou et al. (2024) used the generalization 746 power of different interactions to explain the generalization power of DNNs. 747

Unlike evaluations on language generation results, we propose a method that leverages interaction-based explanations to evaluate the correctness of decision-making logic encoded by a LLM. This approach enables us to evaluate the alignment between the decision-making logic of LLMs and human cognition.

751 752 753

754

B PROOF OF THEOREM

Theorem 1 (Universal matching property) Given an input sample x, the network output score $v(\mathbf{x}_T) \in \mathbb{R}$ on each masked sample $\{\mathbf{x}_T | T \subseteq N\}$ can be well matched by a surrogate logical model

⁷⁵⁶ ⁷⁵⁷ ⁷⁵⁸ $h(\mathbf{x}_T)$ on each masked sample $\{\mathbf{x}_T | T \subseteq N\}$. The surrogate logical model $h(\mathbf{x}_T)$ uses the sum of AND interactions and OR interactions to accurately fit the network output score $v(\mathbf{x}_T)$.

$$\forall T \subseteq N, v(\mathbf{x}_T) = h(\mathbf{x}_T).$$

763 764

$$h(\mathbf{x}_{T}) = v(\mathbf{x}_{\emptyset}) + \sum_{S \subseteq N, S \neq \emptyset} \mathbb{1}\left(\frac{\mathbf{x}_{T} \text{ triggers}}{\text{AND relation } S}\right) \cdot I_{\text{and}}(S|\mathbf{x}_{T}) + \mathbb{1}\left(\frac{\mathbf{x}_{T} \text{ triggers}}{\text{OR relation } S}\right) \cdot I_{\text{or}}(S|\mathbf{x}_{T})$$

$$= \underbrace{v(\mathbf{x}_{\emptyset}) + \sum_{S \subseteq T, S \neq \emptyset} I_{\text{and}}(S|\mathbf{x}_{T})}_{v_{\text{and}}(\mathbf{x}_{T})} + \underbrace{\sum_{S \subseteq N, S \cap T \neq \emptyset} I_{\text{or}}(S|\mathbf{x}_{T})}_{v_{\text{or}}(\mathbf{x}_{T})}$$
(8)

765 766

> 801 802

804 805

808

Let us set a surrogate logical model $h(\mathbf{x}_T) = v(\mathbf{x}_T), \forall T \subseteq N$, which utilizes the sum of AND interactions $I_{and}(S|\mathbf{x})$ and OR interactions $I_{or}(S|\mathbf{x})$ in Equation (2) to fit the network output score $v(\mathbf{x}_T), i.e., v(\mathbf{x}_T) = h(\mathbf{x}_T) = v_{and}(\mathbf{x}_T) + v_{or}(\mathbf{x}_T).$

To be specific, we use the sum of AND interactions $I_{and}(S|\mathbf{x})$ to compute the component for AND interactions $v_{and}(\mathbf{x}_T)$, *i.e.*, $v_{and}(\mathbf{x}_T) = \sum_{S \subseteq T} I_{and}(S|\mathbf{x}_T)$. Then, we use the sum of OR interactions $I_{or}(S|\mathbf{x})$ to compute the component for OR interactions $v_{or}(\mathbf{x}_T)$, *i.e.*, $v_{or}(\mathbf{x}_T) = \sum_{S \subseteq N, S \cap T \neq \emptyset} I_{or}(S|\mathbf{x}_T)$. Finally, we use the sum of AND-OR interactions to fit the network output score, *i.e.*, $v(\mathbf{x}_T) = h(\mathbf{x}_T) = v_{and}(\mathbf{x}_T) + v_{or}(\mathbf{x}_T)$.

(1) Universal matching property of AND interactions.

Ren et al. (2023a) have used the Haranyi dividend (Harsanyi, 1963) $I_{and}(S|\mathbf{x})$ to state the universal matching property of AND interactions. The output of a well-trained DNN on all 2^n masked samples $\{\mathbf{x}_T | T \subseteq N\}$ could be universally explained by the all interaction primitives in $T \subseteq N$, *i.e.*, $\forall T \subseteq N, v_{and}(\mathbf{x}_T) = \sum_{S \subseteq T} I_{and}(S|\mathbf{x}).$

781 Specifically, the AND interaction (as known as Harsanyi dividend) is defined as $I_{and}(S|\mathbf{x}) := \sum_{L \subseteq S} (-1)^{|S|-|L|} v_{and}(\mathbf{x}_L)$ in Equation (2). To compute the sum of AND interactions $\forall T \subseteq N, \sum_{S \subseteq T} I_{and}(S|\mathbf{x}) = \sum_{S \subseteq T} \sum_{L \subseteq S} (-1)^{|S|-|L|} v_{and}(\mathbf{x}_L)$, we first exchange the order of summation of the set $L \subseteq S \subseteq T$ and the set $S \supseteq L$. That is, we compute all linear combinations of all sets S containing L with respect to the model outputs $v_{and}(\mathbf{x}_L)$, given a set of input variables L, *i.e.*, $\sum_{S:L \subseteq S \subseteq T} (-1)^{|S|-|L|} v_{and}(\mathbf{x}_L)$. Then, we compute all summations over the set $L \subseteq T$.

In this way, we can compute them separately for different cases of $L \subseteq S \subseteq T$. In the following, we consider the cases (1) L = S = T, and (2) $L \subseteq S \subseteq T$, $L \neq T$, respectively.

(1) When L = S = T, the linear combination of all subsets S containing L with respect to the model output $v_{and}(\mathbf{x}_L)$ is $(-1)^{|T|-|T|}v_{and}(\mathbf{x}_L) = v_{and}(\mathbf{x}_L)$.

(2) When $L \subseteq S \subseteq T, L \neq T$, the linear combination of all subsets S containing L with respect to the model output $v_{and}(\mathbf{x}_L)$ is $\sum_{S:L\subseteq S\subseteq T}(-1)^{|S|-|L|}v_{and}(\mathbf{x}_L)$. For all sets $S:T\supseteq S\supseteq L$, let us consider the linear combinations of all sets S with number |S| for the model output $v_{and}(\mathbf{x}_L)$, respectively. Let m := |S| - |L|, $(0 \le m \le |T| - |L|)$, then there are a total of $C^m_{|T|-|L|}$ combinations of all sets S of order |S|. Thus, given L, accumulating the model outputs $v_{and}(\mathbf{x}_L)$ corresponding to all $S \supseteq L$, then $\sum_{S:L\subseteq S\subseteq T}(-1)^{|S|-|L|}v_{and}(\mathbf{x}_L) = v_{and}(\mathbf{x}_L) \cdot \underbrace{\sum_{m=0}^{|T|-|L|}C^m_{|T|-|L|}(-1)^m}_{=0} = 0$.

Please see the complete derivation of the following formula.

$$\sum_{S\subseteq T} I_{\text{and}}(S|\mathbf{x}_T) = \sum_{S\subseteq T} \sum_{L\subseteq S} (-1)^{|S|-|L|} v_{\text{and}}(\mathbf{x}_L)$$

$$= \sum_{L\subseteq T} \sum_{S:L\subseteq S\subseteq T} (-1)^{|S|-|L|} v_{\text{and}}(\mathbf{x}_L)$$

$$= \underbrace{v_{\text{and}}(\mathbf{x}_T)}_{L=T} + \sum_{L\subseteq T, L\neq T} v_{\text{and}}(\mathbf{x}_L) \cdot \underbrace{\sum_{m=0}^{|T|-|L|} C_{|T|-|L|}^m}_{=0}$$
(9)

$$= a_{1} \cdot (\mathbf{x})$$

Furthermore, we can understand the above equation in a physical sense. Given a masked sample x_T , if \mathbf{x}_T triggers an AND relationship S (the co-appearance of all input variables in S), then $S \subseteq T$. Thus, we accumulate the interaction effects $I_{and}(S|\mathbf{x})$ of any AND relationship S triggered by \mathbf{x}_T as follows,

- \mathbf{x}_{T} triggers

$$v(\mathbf{x}_{\emptyset}) + \sum_{S \subseteq N, S \neq \emptyset} \mathbb{1} \begin{pmatrix} A_T \text{ argsers} \\ \text{AND relation } S \end{pmatrix} \cdot I_{\text{and}}(S | \mathbf{x}_T) \\ = v(\mathbf{x}_{\emptyset}) + \sum_{S \subseteq T, S \neq \emptyset} I_{\text{and}}(S | \mathbf{x}_T) \\ = \sum_{S \subseteq T} I_{\text{and}}(S | \mathbf{x}_T) \qquad //I_{\text{and}}(\emptyset | \mathbf{x}_T) = v_{\text{and}}(\mathbf{x}_{\emptyset}) = v(\mathbf{x}_{\emptyset}) \\ = v_{\text{and}}(\mathbf{x}_T).$$

$$(10)$$

(2) Universal matching property of OR interactions.

According to the definition of OR interactions in Equation (2), we will derive that $\forall T \subseteq$ $N, v_{\text{or}}(\mathbf{x}_T) = \sum_{S \subseteq N, S \subseteq T \neq \emptyset} I_{\text{or}}(S|\mathbf{x}_T), \text{ s.t., } I_{\text{or}}(\emptyset|\mathbf{x}_T) = v_{\text{or}}(\mathbf{x}_{\emptyset}) = 0.$

Specifically, the OR interaction is defined as $I_{or}(S|\mathbf{x}) := -\sum_{L \subset S} (-1)^{|S| - |L|} v_{or}(\mathbf{x}_{N \setminus L})$ in Equa-tion (2). To compute the sum of OR interactions $\forall T \subseteq N, \sum_{S \subseteq N, S \cap T \neq \emptyset} I_{\text{or}}(S|\mathbf{x}_T) =$ $\sum_{S \subseteq N, S \cap T \neq \emptyset} \left| - \sum_{L \subseteq S} (-1)^{|S| - |L|} v_{\text{or}}(\mathbf{x}_{N \setminus L}) \right|$, we first exchange the order of summation of the set $L \subseteq S \subseteq N$ and the set $S \cap T \neq \emptyset$. That is, we compute all linear combinations of all sets S containing L with respect to the model outputs $v_{or}(\mathbf{x}_{N\setminus L})$, given a set of input variables L, *i.e.*, $\sum_{S \cap T \neq \emptyset} \sum_{N \supset S \supset L} (-1)^{|S| - |L|} v_{\text{or}}(\mathbf{x}_{N \setminus L})$. Then, we compute all summations over the set $L \subseteq N$.

In this way, we can compute them separately for different cases of $L \subseteq S \subseteq N, S \cap T \neq \emptyset$. In the following, we consider the cases (1) $L = N \setminus T$, (2) L = N, (3) $L \cap T \neq \emptyset, L \neq N$, and (4) $L \cap T = \emptyset, L \neq N \setminus T$, respectively.

(1) When $L = N \setminus T$, the linear combination of all subsets S containing L with respect to the model output $v_{\text{or}}(\mathbf{x}_{N\setminus L})$ is $\sum_{S\cap T\neq\emptyset, S\supseteq L}(-1)^{|S|-|L|}v_{\text{or}}(\mathbf{x}_{N\setminus L}) = \sum_{S\cap T\neq\emptyset, S\supseteq L}(-1)^{|S|-|L|}v_{\text{or}}(\mathbf{x}_T)$. For all sets $S\supseteq L, S\cap T\neq\emptyset$ (then $S\neq N\setminus T, S\neq L$), let us consider the linear combinations of all sets S with number |S| for the model output $v_{or}(\mathbf{x}_T)$, respectively. Let |S'| := |S| - |L|, $(1 \le |S'| \le |T|)$, then there are a total of $C_{|T|}^{|S'|}$ combinations of all sets S of order |S|. Thus, given L, accumulating the model outputs $v_{\text{or}}(\mathbf{x}_T)$ corresponding to all $S \supseteq L$, then $\sum_{S \cap T \neq \emptyset, S \supseteq L} (-1)^{|S| - |L|} v_{\text{or}}(\mathbf{x}_{N \setminus L}) = 0$ $v_{\rm or}(\mathbf{x}_T) \cdot \underbrace{\sum_{|S'|=1}^{|T|} C_{|T|}^{|S'|}(-1)^{|S'|}}_{|S'|=-v_{\rm or}(\mathbf{x}_T).$

(2) When L = N (then S = N), the linear combination of all subsets S containing L with respect to the model output $v_{\text{or}}(\mathbf{x}_{N\setminus L})$ is $\sum_{S\cap T\neq\emptyset,S\supseteq L}(-1)^{|S|-|L|}v_{\text{or}}(\mathbf{x}_{N\setminus L}) = (-1)^{|\widetilde{N}|-|N|}v_{\text{or}}(\mathbf{x}_{\emptyset}) = v_{\text{or}}(\mathbf{x}_{\emptyset}) = 0, (I_{\text{or}}(\emptyset|\mathbf{x}_{T}) = v_{\text{or}}(\mathbf{x}_{\emptyset}) = 0).$

(3) When $L \cap T \neq \emptyset, L \neq N$, the linear combination of all subsets S containing L with re-spect to the model output $v_{\text{or}}(\mathbf{x}_{N\setminus L})$ is $\sum_{S\cap T\neq\emptyset,S\supset L}(-1)^{|S|-|L|}v_{\text{or}}(\mathbf{x}_{N\setminus L})$. For all sets $S\supseteq$ $L, S \cap T \neq \emptyset$, let us consider the linear combinations of all sets S with number |S| for the model output $v_{\text{or}}(\mathbf{x}_T)$, respectively. Let us split |S| - |L| into |S'| and |S''|, *i.e.*,|S| - |L| = |S'| + |S''|, where $S' = \{i|i \in S, i \notin L, i \in N \setminus T\}$, $S'' = \{i|i \in S, i \notin L, i \in T\}$ (then $0 \leq |S''| \leq |T| - |T \cap L|$) and S' + S'' + L = S. In this way, there are a total of $C_{|T|-|T\cap L|}^{|S''|}$ combinations of all sets S'' of order |S''|. Thus, given L, accumulating the model outputs $v_{\text{or}}(\mathbf{x}_{N\setminus L})$ corresponding to all $S \supseteq L$, then $\sum_{S \cap T \neq \emptyset, S \supseteq L} (-1)^{|S| - |L|} v_{\text{or}}(\mathbf{x}_{N\setminus L}) = v_{\text{or}}(\mathbf{x}_{N\setminus L}) \cdot \sum_{S' \subseteq N\setminus T\setminus L} \underbrace{\sum_{|S''|=0}^{|T| - |T \cap L|} C_{|T| - |T \cap L|}^{|S''|} (-1)^{|S'| + |S''|}}_{|S''| = 0} = 0.$

(4) When $L \cap T = \emptyset, L \neq N \setminus T$, the linear combination of all subsets S containing L with respect to the model output $v_{\text{or}}(\mathbf{x}_{N \setminus L})$ is $\sum_{S:S \cap T \neq \emptyset, S \supset L} (-1)^{|S| - |L|} v_{\text{or}}(\mathbf{x}_{N \setminus L})$. Similarly, let us split |S| - |L| into |S'| and |S''|, *i.e.*, |S| - |L| = |S'| + |S''|, where $S' = \{i | i \in S, i \notin L, i \in N \setminus T\}$, $S'' = \{i | i \in S, i \notin T\}$ (then $0 \le |S''| \le |T|$) and S' + S'' + L = S. In this way, there are a total of $C_{|T|}^{|S''|}$ combinations of all sets S'' of order |S''|. Thus, given L, accumulating the model outputs $v_{\text{or}}(\mathbf{x}_{N\setminus L})$ corresponding to all $S \supseteq L$, then $\sum_{S \cap T \neq \emptyset, S \supset L} (-1)^{|S| - |L|} v_{\text{or}}(\mathbf{x}_{N\setminus L}) = 0$ $v_{\rm or}(\mathbf{x}_{N\setminus L}) \cdot \sum_{S' \subseteq N\setminus T\setminus L} \underbrace{\sum_{|S''|=0}^{|T|} C_{|T|}^{|S''|}(-1)^{|S'|+|S''|}}_{|S''|=0} = 0.$ Please see the complete derivation of the following formula. $\sum\nolimits_{S \subseteq N, S \cap T \neq \emptyset} I_{\mathrm{or}}(S | \mathbf{x}_T) = \sum\nolimits_{S \subseteq N, S \cap T \neq \emptyset} \left[-\sum\nolimits_{L \subseteq S} (-1)^{|S| - |L|} v_{\mathrm{or}}(\mathbf{x}_{N \setminus L}) \right]$ $= -\sum\nolimits_{L \subseteq N} \sum\nolimits_{S \cap T \neq \emptyset, N \supseteq S \supseteq L} (-1)^{|S| - |L|} v_{\mathrm{or}}(\mathbf{x}_{N \setminus L})$ $= - \left| \sum_{|S'|=1}^{|T|} C_{|T|}^{|S'|} (-1)^{|S'|} \right| \cdot \underbrace{v_{\text{or}}(\mathbf{x}_T)}_{\mathbf{v}_{\text{or}}(\mathbf{x}_T)} - \underbrace{v_{\text{or}}(\mathbf{x}_{\emptyset})}_{\mathbf{v}_{\text{or}}(\mathbf{x}_T)} - \underbrace{v_{\text{or}}(\mathbf{x}_{\emptyset})}_{\mathbf{v}_{\text{or}}(\mathbf{x}_T)} \right|$ $-\sum_{L \cap T \neq \emptyset, L \neq N} \left[\sum_{S' \subseteq N \setminus T \setminus L} \left(\sum_{|S''|=0}^{|T|-|T \cap L|} C_{|T|-|T \cap L|}^{|S''|}(-1)^{|S'|+|S''|} \right) \right] \cdot v_{\text{or}}(\mathbf{x}_{N \setminus L})$ $-\sum_{L\cap T=\emptyset, L\neq N\setminus T} \left[\sum_{S'\subseteq N\setminus T\setminus L} \left(\sum_{|S''|=0}^{|T|} C_{|T|}^{|S''|}(-1)^{|S'|+|S''|} \right) \right] \cdot v_{\mathrm{or}}(\mathbf{x}_{N\setminus L})$ $= -(-1) \cdot v_{\text{or}}(\mathbf{x}_T) - v_{\text{or}}(\mathbf{x}_{\emptyset}) - \sum_{L \cap T \neq \emptyset, L \neq N} \left[\sum_{S' \subseteq N \setminus T \setminus L} 0 \right] \cdot v_{\text{or}}(\mathbf{x}_{N \setminus L})$ $-\sum_{L\cap T=\emptyset, L\neq N\setminus T} \left|\sum_{S'\subset N\setminus T\setminus T} 0\right] \cdot v_{\mathrm{or}}(\mathbf{x}_{N\setminus L})$ $= v_{\rm or}(\mathbf{x}_T)$ (11)

Furthermore, we can understand the above equation in a physical sense. Given a masked sample \mathbf{x}_T , if \mathbf{x}_T triggers an OR relationship S (the presence of any input variable in S), then $S \cap T \neq \emptyset$, $S \subseteq N$. Thus, we accumulate the interaction effects $I_{or}(S|\mathbf{x})$ of any OR relationship S triggered by \mathbf{x}_T as follows,

$$\sum_{\substack{S \subseteq N, S \neq \emptyset}} \mathbb{1}\left(\frac{\mathbf{x}_T \text{ triggers}}{\text{OR relation } S}\right) \cdot I_{\text{or}}(S|\mathbf{x}_T)$$

$$= \sum_{\substack{S \subseteq N, S \cap T \neq \emptyset}} I_{\text{or}}(S|\mathbf{x}_T)$$

$$= v_{\text{or}}(\mathbf{x}_T).$$
(12)

(3) Universal matching property of AND-OR interactions.

With the universal matching property of AND interactions and the universal matching property of OR interactions, we can easily get $v(\mathbf{x}_T) = h(\mathbf{x}_T) = v_{and}(\mathbf{x}_T) + v_{or}(\mathbf{x}_T) = v(\mathbf{x}_{\emptyset}) + \sum_{S \subseteq T, S \neq \emptyset} I_{and}(S|\mathbf{x}_T) + \sum_{S \subseteq N, S \cap T \neq \emptyset} I_{or}(S|\mathbf{x}_T)$, thus, we obtain the universal matching property of AND-OR interactions.

C PROOF OF LEMMA

Lemma 1 (Sparsity property) The surrogate logical model $h(\mathbf{x}_T)$ on each randomly masked sample $\mathbf{x}_T, T \subseteq N$ mainly uses the sum of a small number of salient AND interactions and salient OR interactions to approximate the network output score $v(\mathbf{x}_T)$.

915
916
917

$$v(\mathbf{x}_T) = h(\mathbf{x}_T) \approx v(\mathbf{x}_{\emptyset}) + \sum_{S \in \Omega^{\text{and}}} \mathbb{1} \left(\frac{\mathbf{x}_T \text{ triggers}}{\text{AND relation } S} \right) \cdot I_{\text{and}}(S|\mathbf{x}_T) + \sum_{S \in \Omega^{\text{or}}} \mathbb{1} \left(\frac{\mathbf{x}_T \text{ triggers}}{\text{OR relation } S} \right) \cdot I_{\text{or}}(S|\mathbf{x}_T)$$
(13)

Ren et al. (2024) have proven that under some common conditions⁶, the confidence score $v_{and}(\mathbf{x}_T)$ of a well-trained DNN on all 2^n masked samples $\{\mathbf{x}_T | T \subseteq N\}$ could be universally approximated by a small number of AND interactions $T \in \Omega^{\text{and}}$ with salient interaction effects $I_{\text{and}}(T|\mathbf{x})$, *s.t.*, $|\Omega^{\text{and}}| \ll 2^n$, *i.e.*, $\forall T \subseteq N$, $v_{\text{and}}(\mathbf{x}_T) = \sum_{S \subseteq T} I_{\text{and}}(S|\mathbf{x}) \approx \sum_{S \subseteq T: S \in \Omega^{\text{and}}} I_{\text{and}}(S|\mathbf{x})$.

According to Equation (10), $v_{\text{and}}(\mathbf{x}_T) = \sum_{S \subseteq T} I_{\text{and}}(S|\mathbf{x}) = v(\mathbf{x}_{\emptyset}) + \sum_{S \subseteq N, S \neq \emptyset} \mathbb{1}\begin{pmatrix} \mathbf{x}_T \text{ triggers} \\ \text{AND relation } S \end{pmatrix} \cdot I_{\text{and}}(S|\mathbf{x}_T).$ Therefore, $v_{\text{and}}(\mathbf{x}_T) \approx v(\mathbf{x}_{\emptyset}) + \sum_{S \in \Omega^{\text{and}}} \mathbb{1}\begin{pmatrix} \mathbf{x}_T \text{ triggers} \\ \text{AND relation } S \end{pmatrix} \cdot I_{\text{and}}(S|\mathbf{x}_T).$

Besides, as proven in Appendix D, the OR interaction can be considered as a specific AND in-teraction. Thus, the confidence score $v_{\rm or}(\mathbf{x}_T)$ of a well-trained DNN on all 2^n masked samples $\{\mathbf{x}_T | T \subseteq N\}$ could be universally approximated by a small number of OR interactions $T \in \Omega^{\text{or}} \text{ with salient interaction effects } I_{\text{or}}(T|\mathbf{x}), \text{ s.t., } |\Omega^{\text{or}}| \ll 2^n. \text{ Similarly, } v_{\text{or}}(\mathbf{x}_T) = \sum_{S \subseteq N, S \neq \emptyset} \mathbb{1}\binom{\mathbf{x}_T \text{ triggers}}{\operatorname{OR relation } S} \cdot I_{\text{or}}(S|\mathbf{x}_T) \approx \sum_{S \in \Omega^{\text{or}}} \mathbb{1}\binom{\mathbf{x}_T \text{ triggers}}{\operatorname{OR relation } S} \cdot I_{\text{or}}(S|\mathbf{x}_T).$

In this way, the surrogate logical model $h(\mathbf{x}_T)$ on each randomly masked sample $\mathbf{x}_T, T \subseteq N$ mainly uses the sum of a small number of salient AND interactions and salient OR interactions to approximate the network output score $v(\mathbf{x}_T)$, *i.e.*, $v(\mathbf{x}_T) = h(\mathbf{x}_T) = v_{and}(\mathbf{x}_T) + v_{or}(\mathbf{x}_T) \approx v_{or}(\mathbf{x}_T)$ $v(\mathbf{x}_{\emptyset}) + \sum_{S \in \Omega^{\text{and}}} \mathbb{1}(\underset{\text{AND relation } S}{\overset{\mathbf{x}_T \text{ triggers}}{}}) \cdot I_{\text{and}}(S|\mathbf{x}_T) + \sum_{S \in \Omega^{\text{or}}} \mathbb{1}(\underset{\text{OR relation } S}{\overset{\mathbf{x}_T \text{ triggers}}{}}) \cdot I_{\text{or}}(S|\mathbf{x}_T).$

OR INTERACTIONS CAN BE CONSIDERED SPECIFIC AND INTERACTIONS D

The OR interaction $I_{\rm or}(S|\mathbf{x})$ can be considered as a specific AND interaction $I_{\rm and}(S|\mathbf{x})$, if we inverse the definition of the masked state and the unmasked state of an input variable.

Given a DNN $v : \mathbb{R}^n \to \mathbb{R}$ and an input sample $\mathbf{x} \in \mathbb{R}^n$, if we arbitrarily mask the input sample, we can get 2^n different masked samples $\mathbf{x}_S, \forall S \subseteq N$. Specifically, let us use baseline values $\mathbf{b} \in \mathbb{R}^n$ to represent the masked state of a masked sample x_S , *i.e.*,

$$(\mathbf{x}_S)_i = \begin{cases} x_i, & i \in S \\ b_i, & i \notin S \end{cases}$$
(14)

Conversely, if we inverse the definition of the masked state and the unmasked state of an input variable, *i.e.*, we consider b as the input sample, and consider the original value x as the masked state, then the masked sample \mathbf{b}_S can be defined as follows.

$$(\mathbf{b}_S)_i = \begin{cases} b_i, & i \in S\\ x_i, & i \notin S \end{cases}$$
(15)

According to the above definition of a masked sample in Equations (14) and (15), we can get $\mathbf{x}_{N\setminus S} = \mathbf{b}_S$. To simply the analysis, if we assume that $v_{\text{and}}(\mathbf{x}_T) = v_{\text{or}}(\mathbf{x}_T) = 0.5v(\mathbf{x}_T)$, then the OR interaction $I_{or}(S|\mathbf{x})$ in Equation (2) can be regarded as a specific AND interaction $I_{and}(S|\mathbf{b})$ as follows.

$$I_{\rm or}(S|\mathbf{x}) = -\sum_{T \subseteq S} (-1)^{|S| - |T|} v_{\rm or}(\mathbf{x}_{N \setminus T}),$$

$$= -\sum_{T \subseteq S} (-1)^{|S| - |T|} v_{\rm or}(\mathbf{b}_T),$$

$$= -\sum_{T \subseteq S} (-1)^{|S| - |T|} v_{\rm and}(\mathbf{b}_T),$$

$$= -I_{\rm and}(S|\mathbf{b}).$$

(16)

Ε GENERALIZATION POWER OF INTERACTIONS OVER DIFFERENT ORDERS

In this section, we will give the definition and quantification of the generalization power of interactions over different orders. The generalization power of an interaction is defined as the transferability

⁶There are three assumptions. (1) The high order derivatives of the DNN output with respect to the input variables are all zero. (2) The DNN works well on the masked samples, and yield higher confidence when the input sample is less masked. (3) The confidence of the DNN does not drop significantly on the masked samples.



frequently occurs in the training samples to test samples. Specifically, if an interaction pattern $S \subseteq N$ frequently occurs in the training set, but rarely appears in the test set, then the interaction pattern S exhibits low generalization power. Conversely, if an interaction pattern S consistently appears in both the training and test sets, it demonstrates high generalization power.

1015 Specifically, for a given classification task, Zhou et al. (2024) defined the generalization power of 1016 *m*-order interactions *w.r.t.* the category c as the Jaccard similarity between the interactions observed 1017 in the training samples and those in the test samples for each category c.

$$\sin(\hat{I}_{\text{train},c}^{(m)}, \hat{I}_{\text{test},c}^{(m)}) = \frac{\|\min(\hat{I}_{\text{train},c}^{(m)}, \hat{I}_{\text{test},c}^{(m)})\|_{1}}{\|\max(\hat{I}_{\text{train},c}^{(m)}, \hat{I}_{\text{test},c}^{(m)})\|_{1}}$$
(17)

1020 1021

where $\hat{I}_{\text{train},c}^{(m)} = [(\max(I_{\text{train},c}^{(m)}, 0))^{\mathsf{T}}, (-\min(I_{\text{train},c}^{(m)}, 0))^{\mathsf{T}}]^{\mathsf{T}} \in \mathbb{R}^{2d}$ is conducted from $I_{\text{train},c}^{(m)}$ to ensure that all elements are non-negative. Here, $I_{\text{train},c}^{(m)} = [I_{\text{train},c}^{(m)}(S_1), I_{\text{train},c}^{(m)}(S_2), \cdots, I_{\text{train},c}^{(m)}(S_d)]^{\mathsf{T}} \in \mathbb{R}^{d}$ represents the distribution of *m*-order interactions over the training samples for category c, where $d = C_n^m$ enumerates all possible *m*-order interactions. Specifically, $I_{\text{train},c}^{(m)}(S_i) =$

Table 1: Average ratio (%) of reliable interaction effects of each order on the SaulLM-7B-Instruct model.

 order	1	2	3	4	5	6	7	Q	0	10
order	1	2	5	4	5	0	1	0	9	10
$s_o^{\text{reliable},+}$	60.19	56.53	49.51	48.86	43.74	30.92	42.86	NAN	NAN	NAN
$s_o^{\mathrm{reliable},-}$	NAN	66.79	53.89	58.67	50.34	35.20	52.38	NAN	NAN	NAN

Table 2: Average ratio (%) of reliable interaction effects of each order on the BAI-Law-13B model.

order	1	2	3	4	5	6	7	8	9	10
$s_o^{\text{reliable},+}$	56.22	71.24	49.02	49.17	46.56	40.10	31.70	25.00	22.22	NAN
$s_o^{\mathrm{reliable},-}$	58.15	71.06	69.68	63.02	49.73	35.98	42.86	NAN	NAN	NAN

1040 $\mathbb{E}_{\mathbf{x}\in\mathcal{D}_{\text{train},c}}[I(S_i|\mathbf{x})]$ denotes the average interaction effect of the set S_i across different training samples within category c.

1042 1043 Therefore, for each category c, a high similarity $\sin(\hat{I}_{\text{train},c}^{(m)}, \hat{I}_{\text{test},c}^{(m)})$ indicates that most *m*-order interactions from the training samples generalize well to the test samples.

Using the average similarity over different categories, *i.e.*, similarity = $\mathbb{E}_c[sim(\hat{I}_{train,c}^{(m)}, \hat{I}_{test,c}^{(m)})]$, Zhou et al. (2024) have empirically found that the low-order interactions usually exhibit stronger generalization power than high-order interactions. Specifically, Figure 4 in (Zhou et al., 2024) shows that compared to high order interaction patterns, DNNs are more likely to extract similar low order interaction patterns from both training and test data.

1050 1051

1052 F ACCURACY OF THE LEGAL LLM

1053

1054 Colombo et al. (2024) reported the accuracy of the SaulLM-7B-Instruct model, which achieved 1055 state-of-the-art results among 7B models, within the legal domain. Specifically, they followed (Guha 1056 et al., 2023) to use balanced accuracy as the metric. Balanced accuracy shows its strength for han-1057 dling imbalanced classification tasks. They tested the balanced accuracy on two popular bench-1058 marks, i.e., the LegalBench-Instruct benchmark (Guha et al., 2023) and the Massive Multitask Language Understanding (MMLU) benchmark (Hendrycks et al., 2021). The LegalBench-Instruct benchmark is a supplemental iteration of LegalBench (Guha et al., 2023), designed to evaluate the legal proficiency of LLMs. To further evaluate the performance of LLMs in legal contexts, the au-1061 thors incorporated legal tasks from the MMLU benchmark, focusing specifically on the international 1062 law, professional law and jurisprudence. 1063

Colombo et al. (2024) compared the SaulLM-7B-Instruct model to other 7B and 13B open-source models, including Mistral-7B (Jiang et al., 2023) and the Llama2 family (Touvron et al., 2023). Table 4 shows that SaulLM-7B-Instruct achieved state-of-the-art performance on the LegalBench-Instruct benchmark, outperforming its competitors in the legal domain.

1068 1069

107 107

Table 3: Comparison of LLMs on the LegalBench-Instruct benchmark.

0										
1	LLMs	SaulLM-7B-Instruct	Mistral-7B-v1	Mistral-7B-v2	Llama2-13B-chat	Llama2-7B-chat				
2	accuracy	0.61	0.55	0.52	0.45	0.39				

1073

To further confirm the observations on the LegalBench-Instruct, (Colombo et al., 2024) conducted additional experiments on the legal tasks from the MMLU benchmark. The SaulLM-7B-Instruct model exhibited strong performance across all three tasks, including international law, professional law, and jurisprudence tasks.

1079 Besides, Institute (2023) has not yet reported the specific classification accuracy of the BAI-Law-13B model, leaving its performance on certain benchmarks unclear.

1034 1035 1036

1039



actions were annotated as relevant tokens. We also translated the legal case to English as the input of
 the SaulLM-7B-Instruct model. (b) Judgements predicted by the two legal LLMs, which were both
 correct according to laws of the two countries. (c,d) We quantified the reliable and unreliable inter action effects of different orders. The SaulLM-7B-Instruct model used 66.1% reliable interaction
 effects, while the BAI-Law-13B model encoded 87.2% reliable interaction effects.



Figure 8: More results of judgments influenced by unreliable sentimental tokens. (d) The SaulLM-7B-Instruct model used 35.3% reliable interaction effects, while the BAI-Law-13B model encoded 48.5% reliable interaction effects.

¹¹¹⁹ G MORE EXPERIMENT RESULTS AND DETAILS

1121 G.1 More results of judgments influenced by unreliable sentimental tokens

1123 We conducted more experiments to show the judgments influenced by unreliable sentimental to-1124 kens in Figure 7, Figure 8, and Figure 9, respectively. We observed that a considerable number 1125 of interactions contributing to the confidence score $v(\mathbf{x})$ were attributed to semantically irrelevant 1126 or unreliable sentimental tokens. In different legal cases, the ratio of reliable interaction effects to 1127 all salient interactions was within the range of 32.6% to 87.1%. It means that about $13\sim 68\%$ of 1128 interactions used semantically irrelevant tokens or unreliable sentimental tokens for the judgment.

1129

1118

1120

```
1130 G.2 MORE RESULTS OF JUDGMENTS AFFECTED BY INCORRECT ENTITY MATCHING
1131
```

1132 We conducted more experiments to show the judgments affected by incorrect entity matching in Fig-1133 ure 10, Figure 11, and Figure 12, respectively. We observed that a considerable ratio of the confidence score $v(\mathbf{x})$ was mistakenly attributed to interactions on criminal actions made by incorrect



Figure 10: More results of judgments affected by incorrect entity matching. (a) A number of irrelevant tokens were annotated in the legal case, including the time and actions that were not the direct reason for the judgment. Criminal actions of the defendant were annotated as relevant tokens. Criminal actions of the unrelated person were annotated as forbidden tokens. (b) Judgements predicted by the two legal LLMs, which were both correct according to laws of the two countries. (c,d) We measured the reliable and unreliable interaction effects of different orders. The SaulLM-7B-Instruct model used 67.8% reliable interaction effects, while the BAI-Law-13B model encoded 64.1% reliable interaction effects.



Figure 11: More results of judgments affected by incorrect entity matching. (d) The SaulLM-7B-Instruct model used 63.7% reliable interaction effects, while the BAI-Law-13B model encoded 31.9% reliable interaction effects.

1184

entities. In different legal cases, the ratio of reliable interaction effects to all salient interactions was
within the range of 31.9% to 67.8%. It means that about 22~68% of interactions used semantically irrelevant tokens for the judgment, or was mistakenly attributed on criminal actions made by incorrect entities.



Figure 12: More results of judgments affected by incorrect entity matching. (d) The SaulLM-7B-Instruct model used 52.5% reliable interaction effects, while the BAI-Law-13B model encoded 42.2% reliable interaction effects.

1199

1203

G.3 MORE RESULTS OF JUDGMENTS BIASED BY DISCRIMINATION IN OCCUPATION

Experiment results of judgments biased by discrimination in occupation in Section 3. Figure 16 illustrates additional examples of how occupation influences the judgment of the legal case, which were tested on the SaulLM-7B-Instruct model. It shows that if we replaced "a judge" with law-related occupations, such as "a lawyer" and "a policeman," the judgment remained "robbery."
Besides, the occupation "a programmer" changed the judgment to "not mentioned." The interactions containing the occupation token (*i.e.*, "a judge", "a lawyer", "a policeman", "a programmer", and "a volunteer") were important factors that changed the ratio of reliable interactions from 81.4% to 84.0%. This suggested that the legal LLM sometimes had considerable occupation bias.

1212 Futhermore, Figure 17 shows the test of the BAI-Law-13B model on the legal case, in which Andy, 1213 the victim with varying occupations, was robbed of his belongings by two suspicious men. Similarly, 1214 we found that the BAI-Law-13B model encoded interactions with the occupation tokens "a judge," 1215 which boosted the confidence of the judgment "robbery." More interestingly, if we substituted the 1216 occupation tokens "a judge" to "a policeman," the interaction of the occupation "a policeman," decreased from 0.29 to 0.11. The interactions containing the occupation token were important factors 1217 1218 that changed the ratio of reliable interactions from 78.9% to 87.1%. This suggested that the legal 1219 LLM sometimes had considerable occupation bias.

1220 More results of judgments biased by discrimination in occupation. We conducted more ex-1221 periments to show the judgments biased by discrimination in occupation in Figure 13, Figure 14, 1222 and Figure 15, respectively. We found that the legal LLM usually used interactions on the occupa-1223 tion information to compute the confidence score $v(\mathbf{x})$. In different legal cases, the ratio of reliable interaction effects to all salient interactions was within the range of 30.1% to 63.7%. In particular, 1224 in Figure 13, changing the occupation from "lawyer" to "programmer" results in a decrease of the 1225 reliable interactions from 63.7% to 57.3%. The difference of interactions containing the occupation 1226 token changes the model output from "Larceny" to "Theft." 1227

1228 1229

G.4 EXPERIMENT DETAILS OF MASKED SAMPLES

This section discusses how to obtain the masked sample $\mathbf{x}_T, T \subseteq N$. Given the confidence score of a DNN $v(\mathbf{x})$ and an input sample $\mathbf{x} = [x_1, x_2, \dots, x_n]^{\mathsf{T}}$ with n input variables, if we arbitrarily mask the input sample \mathbf{x} , we can get 2^n different masked samples $\mathbf{x}_T, \forall T \subseteq N$. Specifically, for each input variable $i \in N \setminus T$, we replace it with the baseline value b_i to represent its masked state. Let us use baseline values $\mathbf{b} = [b_1, b_2, \dots, b_n]^{\mathsf{T}}$ to represent the masked state of a masked sample $\mathbf{x}_T, i.e.$,

$$(\mathbf{x}_T)_i = \begin{cases} x_i, & i \in T \\ b_i, & i \notin T \end{cases}$$
(18)

1236 1237

For sentences in a language generation task, the masking of input variables is performed at the embedding level. Following the approach of (Ren et al., 2024; Shen et al., 2023), we masked inputs at the embedding level by transforming sentence tokens into their corresponding embeddings. Given an input sentence $\mathbf{x} = [x_1, x_2, \dots, x_n]^{\mathsf{T}}$ with *n* input tokens, the *i*-th token x_i is mapped to its



Figure 13: More results of judgments biased by discrimination in occupation. (a) A number of irrelevant tokens were annotated in the legal case, including the occupation, time and actions that are not the direct reason for the judgment. Criminal actions of the defendant were annotated as relevant tokens. (b) The SaulLM-7B-Instruct model predicted the judgment based on the legal case with different occupations, respectively. (c,d) We measured the reliable and unreliable interaction effects of different orders. When the occupation was set to "*lawyer*," the LLM used 63.7% reliable interaction effects. In comparison, when the occupation was set to "*programmer*," the LLM encoded 57.3% reliable interaction effects.



Figure 14: More results of judgments biased by discrimination in occupation. (b) The SaulLM-7B-Instruct model predicted the judgment based on the legal case with different occupations, respectively. (d) When the occupation was set to "*telephone service*," the LLM used 30.1% reliable interaction effects. In comparison, when the occupation was set to "*volunteer*," the LLM encoded 32.7% reliable interaction effects.



Figure 15: More results of judgments biased by discrimination in occupation. (b) The BAI-Law-13B model predicted the judgment based on the legal case with different occupations, respectively. (d) When the occupation was set to "*former thief*," the LLM used 41.3% reliable interaction effects. In comparison, when the occupation was set to "*miner*," the LLM encoded 40.1% reliable interaction effects.

1294

embedding $e_i \in \mathbb{R}^d$, where d is the dimension of the embedding layer. To obtain the masked sample \mathbf{x}_T , if $i \in N \setminus T$, the embedding is replaced with the (constant) baseline value $b_i \in \mathbb{R}^d$, *i.e.*, $e_i = b_i$.



Figure 16: Visualization of judgments biased by discrimination in occupation. (a) A number of irrelevant tokens were annotated in the legal case, including the occupation, time and actions that are not the direct reason for the judgment. Criminal actions of the defendant were annotated as relevant tokens. (b) The SaulLM-7B-Instruct model predicted the judgment based on the legal case with different occupations, respectively. (c,d) We measured the reliable and unreliable interaction effects of different orders. When the occupation was set to "*a lawyer*," the LLM used 82.6% reliable interaction effects. In comparison, when the occupation was set to "*a policeman*," the LLM encoded 84.2% reliable interaction effects.

1325

Otherwise, the embedding remains unchanged, *i.e.*, $e_i = e_i$. Following (Ren et al., 2023b), we trained the (constant) baseline value $b_i \in \mathbb{R}^d$ to extract the sparsest interactions.

1326 G.5 EXPERIMENT DETAILS FOR USING THE SAME DATASET FOR COMPARISON

This section presents the experiment details of using the CAIL2018 dataset (Xiao et al., 2018) to ensure a fair comparison between two legal LLMs. For the BAI-Law-13B model, a Chinese legal LLM, we directly analyzed the Chinese legal cases from the CAIL2018 dataset. In contrast, for the SaulLM-7B-Instruct model, an English legal LLM, we translated the Chinese legal cases into English and performed the analysis on the translated cases, to enable fair comparisons. To simplify the explanation and avoid ambiguity, we only explained the decision-making logic on legal cases, which were correctly judged by the LLM.

Starting with a complete fact descriptions of the legal case from the CAIL2018 dataset, we first condensed the case by removing descriptive details irrelevant to the judgment, retaining only the most informative tokens, such as the time, location, people, and events. To prompt the model to deliver its judgment, we added a structured prompt designed to extract a concise answer. The format is as follows:

- "Question: [Fact descriptions of the case]. What crime did [the defendant] commit? Briefly answer
 the specific charge in one word. Answer: The specific charge is"
- Here, *[Fact descriptions of the case]* is replaced with the details of the specific legal case, and *[the defendant]* is substituted with the name of the defendant.
- To identify potential representation flaws behind the seemingly correct language generation results of legal LLMs, we introduced special tokens that were irrelevant to the judgments. For cases to assess if judgments were influenced by unreliable sentimental tokens, we added such tokens to describe actions in the legal case. We then observed whether a substantial portion of the interactions contributing to the confidence score $v(\mathbf{x})$ were associated with semantically irrelevant or unreliable sentimental tokens. Similarly, in cases where we aimed to detect potential bias based on occupation, we included irrelevant occupation-related tokens for the defendants or victims, and analyzed whether



1376 Figure 17: Visualization of judgments biased by discrimination in occupation. (a) A number of ir-1377 relevant tokens were annotated in the legal case, including the occupation, time and actions that are 1378 not the direct reason for the judgment. Criminal actions of the defendant were annotated as relevant 1379 tokens. (b) The BAI-Law-13B model predicted the judgment based on the legal case with different occupations, respectively. (c,d) We measured the reliable and unreliable interaction effects of dif-1380 ferent orders. When the occupation was set to "a judge," the LLM used 78.9% reliable interaction 1381 effects. In comparison, when the occupation was set to "a policeman," the LLM encoded 87.1% 1382 reliable interaction effects. 1383

the legal LLM leveraged these occupation-related tokens to compute the confidence score $v(\mathbf{x})$ in Equation (1).

1387 Finally, we show the selection of input variables for extracting interactions. As discussed in Sec-1388 tion 2.1, given an input sample x with n input variables, we can extracted at most 2^{n+1} AND-OR 1389 interactions to compute the confidence score $v(\mathbf{x})$. Consequently, the computational cost for extract-1390 ing interactions increases exponentially with the number of input variables. To alleviate this issue, 1391 we followed (Ren et al., 2024; Shen et al., 2023) to select a set of tokens as input variables, while 1392 keeping the remaining tokens as a constant background in Appendix G.4, to compute interactions 1393 among the selected variables. Specifically, we selected 10 informative input variables (tokens or phrases) for each legal case. These input variables were manually selected based on their informa-1394 tiveness for judgements. It was ensured that the removal of all input variables would substantially 1395 change the legal judgment result. 1396

- 1397
- 1398
- 1399
- 1400
- 1401 1402
- 1403