JOPA: EXPLAINING LARGE LANGUAGE MODEL'S GEN-ERATION VIA JOINT PROMPT ATTRIBUTION

Anonymous authors

Paper under double-blind review

ABSTRACT

Large Language Models (LLMs) have demonstrated impressive performances in complex text generation tasks. However, the contribution of the input prompt to the generated content still remains obscure to humans, underscoring the necessity of understanding the causality between input and output pairs. Existing works for providing prompt-specific explanation often confine model output to be classification or next-word prediction. Few initial attempts aiming to explain the entire language generation often treat input prompt texts independently, ignoring their combinatorial effects on the follow-up generation. In this study, we introduce a counterfactual explanation framework based on joint prompt attribution, JoPA, which aims to explain how a few prompt texts collaboratively influences the LLM's complete generation. Particularly, we formulate the task of prompt attribution for generation interpretation as a combinatorial optimization problem, and introduce a probabilistic algorithm to search for the casual input combination in the discrete space. We define and utilize multiple metrics to evaluate the produced explanations, demonstrating both the faithfulness and efficiency of our framework.

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

028 Large Language Models (LLMs), such as GPT4 (Achiam et al., 2023), LLaMA (Touvron et al., 029 2023) and Claude (Anthropic, 2024), have shown excellent performance in various natural language generation tasks including question answering, document summarization, and many more. Despite the great success of LLMs, we still have very limited understanding of the LLM generation behavior 031 - which parts in the input cause the model to generate a certain sequence. Unable to explain the causality between the input prompt and the output generation could cause failure in recognizing 033 potential unintended consequences, such as harmful response (DAN, 2023; Liu et al., 2023; Zou et al., 034 2023; Zhu et al., 2023) and biased generation (Wang et al., 2023) attributed to a specific malicious 035 description in the input. These issues undermine human trust in model usage, thus highlighting a pressing need for developing an interpretation tool that attributes how an input prompt leads to the 037 generated content.

038 Explaining LLM generation through **prompt attribution** involves extracting the most influential prompt texts on the model's entire generation procedure, a realm that remains relatively under-040 explored in current research. While extensive works on input attribution are proposed for text 041 classification interpretation (Chen & Ji, 2020; Modarressi et al., 2023) and next-word generation 042 rationale (Zhao & Shan, 2024; Vafa et al., 2021), they can not be directly applied to explain the 043 full generation sequence due to its complicated joint probability landscape and the autoregressive 044 generation procedure. The complexity for interpreting LLM generation compounds as the model size increases. Another line of works involves prompting LLMs to self-explain their behaviors (Wei et al., 2022). This method relies on the model's innate reasoning capabilities, although current 046 findings suggest that these capabilities may not always be faithful (Turpin et al., 2023; Xu et al., 047 2024). 048

There are limited existing attempts focusing on explaining the relationship between the input prompts
and the complete generated sequence. The most relevant approach, Captum (Miglani et al., 2023),
sequentially determines the importance score for each token by calculating the variations in the joint
probability of generating the targeted output sequence when the token is dropped from the model
input. While being straightforward, this approach treats tokens as independent features, ignoring
their joint semantic influence on the generated output. In fact, tokens may contain overlapping or

054 complementary information. For example, given the input prompt: "Write a story about the doctor 055 and his patient", the most influential components are "doctor" and "patient". Individually removing 056 either of these words would not significantly alter the generated output, resulting in inaccurately 057 low importance score for each of them. This is caused by the semantic interaction among these 058 components, allowing the model to infer the meaning of the omitted one. Existing attribution methods that ablate each token individually fail to capture such combinatorial effect. A straightforward remedy might involve exhaustively assessing all possible combinations to observe the variations in the model 060 generations, which however is impractical due to the vast search space with long-context input 061 prompts. 062

063 To efficiently search the space for generating accurate prompt explanations, we develop our framework 064 JoPA, which provides the counterfactual explanation to highlight which components of input prompts have the fundamental effect on the generated context via solving a combinatorial optimization 065 problem. We aim to explain the generation behavior of model outputs for any given prompt while 066 take the joint effects of the prompt components into account. Assuming that removing the essential 067 parts of the prompt would result in a significant variation in the model's output, we propose the novel 068 objective function and formulate our task of providing faithful counterfactual explanations for the 069 input prompt as an optimization problem. To quantify the influence of token combinations in the prompt on the generations, we incorporate a mask approach for joint prompt attribution. Thus, our 071 goal of extracting the explanations has been converted to finding the optimal mask of the input prompt. 072 We solve this problem by a probabilistic search algorithm, equipped with gradient guidance and 073 probabilistic updates for efficient exploration in the discrete solution space. Our main contributions 074 could be summarized as follows:

075 076

077

079

- We propose a general **interpretation scheme for LLM generation task** that attributes input prompts to the entire generation sequence. Notably, this recipe considers the joint influence of input token combinations on the generation. This motivation naturally formulates a combinatorial optimization problem for explaining generation with the most influential prompt texts.
- We demonstrate *JoPA*, an efficient probabilistic search algorithm to solve the optimization problem.
 JoPA works by searching better token combinations that lead to larger generation changes. It takes the advantage of both the gradient information and the probabilistic search-space strategy, thereby achieving an efficient prompt interpretation tool.
- Our framework demonstrates strong performance on language generation tasks including text summarization, question-answering, and general instruction datasets. The faithfulness of our explanations is evaluated based on a suit of comprehensive metrics considering generation probability, word frequency, and semantic similarity, verifying the transferability and effectiveness of our methods across a variety of tasks. Moreover, the generated explanations demonstrate the effectiveness of our framework to potentially be applied to improve the model's ability, especially making the model safer and more efficient.
- 091 092

2 RELATED WORK

While there are extensive works devoted to explaining language models in the context of text classification tasks (Shi et al., 2022; Han et al., 2021; Shi et al., 2023), relatively few attempts (Zhao & Shan, 2024; Vafa et al., 2021) are proposed to investigate the importance of prompt texts on the entire generation procedure especially for LLMs. This demonstrates a research gap that this work aims to fill in.

098 **Explaining Language Generation** There are many work explaining the predictions generated by LLMs by measuring the importance of the input features to the model's prediction on the classification 100 tasks. One group of studies perturb the specific input by removing, masking, or altering the input 101 features, and evaluate the model prediction changes (Kommiya Mothilal et al., 2021; Wu et al., 102 2020). The other group of works, such as integrated gradients (IG) (Sundararajan et al., 2017), first-103 derivative saliency (Li et al., 2016), and mixed partial derivatives (Tsang et al., 2020) leverage the 104 gradients of the output with respect to the input to determine the input feature importance. Although 105 Archipelago (Tsang et al., 2020) explains the feature attributions by considering the combined effects of the input attributions, it targets for the multi-label classification task and relies on the neutral 106 baseline. As for the generation tasks, Diffmask learns the differentiable mask for each layer of the 107 BERT model, but it aims to analyze how decisions are formed across network hidden layers by a

simple probing classifier. In contrast, our framework targets on offering insights into the relationship
 between the input prompt and model generations by employing the mask to highlight the essential
 prompt attributions.

111 Moreover, a few studies utilize the surrogate model to explain the individual predictions of the 112 black-box models, and the representative method is called LIME (Ribeiro et al., 2016). As for 113 explaining the model's generation behavior, Captum (Miglani et al., 2023) calculates the token's 114 importance score by sequentially measuring the contribution of input token to the output, which lacks 115 of the accounting for the semantic relationships between tokens. Another work, ReAGent (Zhao & 116 Shan, 2024), focuses on the next-word generation task, computing the importance distribution for 117 the next token position. This method ignores the contextual dependencies in the generated output and 118 could not adequately account for dynamically changing generations. Our framework aims to interpret the joint effects of the input prompts on the entire output contexts with considering of the textual 119 information covered the input prompt. 120

121 **Self-explaining by Prompting** As language models increase in scale, prompting-based models 122 demonstrate remarkable abilities in reasoning (Brown et al., 2020), creativity (Oppenlaender, 2022), 123 and adaptability across a range of tasks (Khattak et al., 2023). However, the complex reasoning processes of these models remain elusive and require tailored paradigm to better understand the 124 prompting mechanism. For instance, the chain-of-thought (CoT) paradigm could explain the LLM 125 behaviors by prompting the model to generate the reasoning chain along with the answers (Wei 126 et al., 2022), as pre-trained LLMs have demonstrated a certain ability to self-explain their behaviors. 127 However, recent studies have also suggested that the reasoning chain does not guarantee faithful 128 explanations of the model's behavior (Jacovi & Goldberg, 2020) and the final answer might not 129 always follow the generated reasoning chain (Turpin et al., 2023). xLLM (Chuang et al., 2024) 130 enhances the fidelity of explanation derived from LLMs via a fidelity evaluator, which however is 131 designed for classification tasks. Our efforts are concentrated on explaining LLM generation by 132 analyzing the attribution of the input prompts to the output content, without depending on the model's 133 innate reasoning ability that are not yet satisfactory.

135 3 PRELIMINARIES

We first introduce necessary notions for the LLM generation process, and then discuss limitations of prior attempts explaining the entire generation via prompt attribution.

139 LLM Generation Notions Denote a specific input prompt as a sequence of tokens $x = (x_1, \ldots, x_T), x_i \in \{1, 2, \ldots, |V|\}$, where |V| represents the vocabulary size, T is the length of the input sequence and the set of all token indices is $\mathcal{I} = \{1, 2, \ldots, T\}$. The corresponding generated output y could be represented as a sequence of tokens $y = (y_1, \ldots, y_S)$ with $y_j \in \{1, 2, \ldots, |V|\}$. The output tokens $\{y_i\}_{i=1}^S$ are generated from the LLM f_{θ} parameterized by θ in an autoregressive manner with the probability $p_{\theta}(y|x)$ as:

147

156

134

 $p_{\boldsymbol{\theta}}(\boldsymbol{y}|\boldsymbol{x}) = p_{\boldsymbol{\theta}}(y_1|\boldsymbol{x}) \prod_{i=1}^{S-1} p_{\boldsymbol{\theta}}(y_{i+1}|\boldsymbol{x}, y_i).$ (1)

The probability of generating the output text y given the input prompt x illustrates that the coherence and generation of the output texts heavily rely on the input prompt x, indicating an implicit causal relationships between the input prompt x and the output y.

151 **Limitation of Prior Interpretation Attempt** There are limited prior works about explaining 152 the relationship between the individual input prompts and the entire generated sequence, and their 153 faithfulness is limited by treating input tokens independently. Specifically, Captum (Miglani et al., 154 2023), calculates the importance score for each token by slightly perturbing the input x at the *i*-th 155 token:

$$G_i(\boldsymbol{x};\boldsymbol{\theta}) = p_{\boldsymbol{\theta}}(\boldsymbol{y}|\boldsymbol{x}) - p_{\boldsymbol{\theta}}(\boldsymbol{y}|\boldsymbol{x}_{\mathcal{I}\setminus\{i\}}),$$
(2)

where \mathcal{I} denotes all token indices. Through Eq. (2), one can calculate an attribution score for each token *i* indicating its importance towards the original generation result y. Such method is direct and simple, however, it treats tokens as independent features, ignoring their semantic interaction and joint effect on the generated output. To make it more obvious, let's recall the previous example of "Write a story about a doctor and his patient.". Note that since "doctor" and "patient" are semantically correlated, and removing any of the two words would not have a significant influence on the generated response, their corresponding importance score will not be too high. On the contrary, the most important token detected by this approach would be "his", which is clearly not ideal.

These observations inspire us to develop a new prompt attribution method that considers the semantic correlation among tokens. Specifically, we assume that the content generated by LLMs is primarily influenced by a subset of tokens jointly. The remaining tokens serve as auxiliary or potentially less relevant information. Those subset of tokens, which fundamentally shape the model's output, are viewed as the explanatory tokens for the generated content. In other words, explanatory tokens do not affect the model output independently, but jointly contribute to the generated responses.

4 PROPOSED METHOD

In this section, we propose *JoPA*, a simple yet effective framework designed for generative tasks to explain the attribution of the input prompt by solving a discrete optimization problem. We start with formulating the general objective for the discrete optimization problem, and then we introduce our proposed probabilistic search algorithm for solving the problem.



Figure 1: Overview of *JoPA*. *Left*: Demonstrating the pipeline of the algorithm. *Right*: Illustrating the process of mask m sampling.

Problem Formulation for Joint Prompt Attribution In order to explain language generation via prompt attribution, instead of treating the prompt tokens independently as in previous works (Miglani et al., 2023), we propose to evaluate the joint effect of k prompt tokens on the generated sequence.

Specifically, consider a binary prompt mask $m = (m_1, \ldots, m_T)$ in which $m_i \in \{0, 1\}$ indicates 195 whether the *i*-th token is important or not. The joint probability of generating the original output 196 sequence y given an input masked context $m \odot x$ is computed as $p_{\theta}(y|m \odot x)$, where \odot denotes 197 the Hadamard product. In our paper, we aim to identify a binary mask m, which is a (discrete) learnable parameter, that maximizes the discrepancy in the probability of generating the same 199 output y when comparing the masked input $m \odot x$ to the original input x. A large discrepancy 200 indicates that the masked token combination are the most influential components for generating y_{i} , thus should jointly serve as the attributed explanation. Consequently, this involves training the binary 202 mask m to optimize the following objective function: 203

171

172

173

174

175

 $\max_{\boldsymbol{m} \in \{0,1\}^T} \mathcal{L}(\boldsymbol{m}, \boldsymbol{x}, \boldsymbol{y}; \boldsymbol{\theta}) := p_{\boldsymbol{\theta}} \left(\boldsymbol{y} | \boldsymbol{x} \right) - p_{\boldsymbol{\theta}} \left(\boldsymbol{y} | \boldsymbol{m} \odot \boldsymbol{x} \right)$ s.t. $|\boldsymbol{m}|_1 = T - k,$ (3)

where k represents the number of explanatory tokens. Intuitively, optimizing the objective outlined in Eq. (3) suggests that we want to find the top k important tokens which, if they are masked, lead to a substantial variation in model's output probabilities. Compared with prior methods for individual token attribution stated in Eq. (2), our formulation in Eq. (3) measures the joint attribution of a subset of tokens, capturing token interactions to enable more accurate prompt explanations.

Challenges in Solving Eq. (3) Note that Eq. (3) is a constraint discrete optimization problem that is non-trivial to solve. One naive solution would be transforming the discrete optimization problem into a continuous one, i.e., let $m \in [0, 1]^T$, and formulate the constraint into a regularization term. Then one can adopt traditional gradient descent based optimization solutions to solve the problem. However, such a strategy would require extensive gradient calculations and backward 216 steps on the original inputs, which can be inefficient in practice especially for LLMs. Moreover, the 217 obtained continuous mask is not the final explanation we want. In fact, the approximation error when 218 transforming the continuous mask back into the discrete space can also be quite significant, leading to 219 worse performances. Another straightforward solution involves searching through all possible token 220 combinations exhaustively. However, this could be impractical as well due to the enormous search space especially while facing long-context inputs. Therefore, we hope to develop a new method that 221 adopts a search-based strategy to simplify the algorithm design and satisfy our constraints, while also 222 leveraging gradient information for efficient optimization. 223

224 225

242

Algorithm 1 Explainable Prompt Generator: JoPA

226 **Input:** Input tokens x, output tokens y, the integer k denoting the number of explanatory tokens, 227 $1 \le k \le T$, input mask $m^{(0)} = 1$, and the sampling numbers N. 228 **Output:** Optimal mask $m^{(N)}$ 229 1: $\boldsymbol{g} = |\nabla_{\boldsymbol{m}^{(0)}} \mathcal{L}(\boldsymbol{m}^{(0)}, \boldsymbol{x}, \boldsymbol{y}; \boldsymbol{\theta})|$ 230 2: Set $m^{(1)}$ as the top-k value mask for q: $m^{(1)} = m^{(0)}$; $m_i^{(1)} = 0, \forall i \in (q)$ 231 3: for n = 1 to N do 232 Sample $l \sim \operatorname{softmax}(\boldsymbol{m}^{(n)} \odot \boldsymbol{g})$ 4: 233 Sample $v \sim \operatorname{softmax}((1 - m^{(n)}) \odot g)$ 5: 234 $\boldsymbol{m}^{\text{tmp}} = \boldsymbol{m}^{(n)}.\text{copy}()$; switch the value of m_l^{tmp} and m_v^{tmp} 6: 235 if $p_{\theta}(\boldsymbol{y} | \boldsymbol{m}^{\text{tmp}} \odot \boldsymbol{x}) < p_{\theta}(\boldsymbol{y} | \boldsymbol{m}^{(n)} \odot \boldsymbol{x})$ then 7: 236 $\boldsymbol{m}^{(n+1)} = \boldsymbol{m}^{\mathsf{tmp}}$ 8: 237 9: else 238 $m^{(n+1)} = m^{(n)}$ 10: 239 end if 11: 240 12: end for 241

Proposed Probabilistic Search Algorithm To tackle this challenge, we propose *JoPA*, a novel 243 explainable prompt generator, for efficiently obtaining an optimal solution for Eq. (3). We summarize 244 our JoPA in Algorithm 1. The high-level pipeline is illustrated in Figure 1 left: we initialize and 245 maintain exactly k entries in the mask m to be zero to enforce the constraint and capture their 246 joint influence of being masked; the mask is then iteratively updated by sampling indexes for value 247 swapping, searching towards the direction with increased generation loss $\mathcal{L}(m, x, y; \theta)$. At the 248 core of this pipeline is the sampling and update of the discrete mask, which demands an efficient 249 exploration in the vast search space. Figure 1 right shows this step, which is featured by the following two essential components, the gradient-guided masking and the probabilistic search update. 250 Specifically, it illustrates the process of sampling a non-zero entry and a zero entry from m to 251 swap their values. 252

253 Gradient-Guided Masking: Trivial solutions that set the mask m by uniformly random could cost 254 massive sampling to hit the right optimization direction. Gradient is a common indicator of feature 255 importance, as evidenced in existing practices (Ebrahimi et al., 2018; Zou et al., 2023; Shin et al., 2020). Therefore, we propose to use gradient as a guidance to set and sample the mask for more 256 efficient optimization. Specifically, we begin with the binary mask $m^{(0)} = 1$ which indicates 257 that all tokens in the input x are marked as non-explanatory ones. Then we compute the gradient 258 $\nabla_{\boldsymbol{m}^{(0)}} \mathcal{L}(\boldsymbol{m}^{(0)}, \boldsymbol{x}, \boldsymbol{y}; \boldsymbol{\theta})$ of the loss function in Eq. (3), and denote the magnitude of gradients as 259 $g = |\nabla_{m^{(0)}} \mathcal{L}(m^{(0)}, x, y; \theta)|$. Note that the components with larger gradient magnitudes in g imply 260 that altering the corresponding tokens could result in a sharper change to the generated outputs. In 261 order to initialize the explanatory k tokens guided by the gradient magnitude, we set the binary 262 mask $m^{(1)}$ where $m_i^{(1)} = 0$ indicates g_i is the top-k value in g. The gradient guidance and mask 263 initialization are obtained following Line 1-2 in the algorithm. 264

265Probabilistic Search Update: While gradient is informative, greedily determining explanatory tokens266by top gradients lacks exploration, leading to suboptimal solutions. Therefore, we propose a prob-267abilistic search mechanism for mask sampling and update. Specifically when updating the current268binary mask m, we iteratively sample a non-zero entry l from the mask and swap its value with a269sampled zero entry v to explore a new (and potentially better) solution for the mask. Particularly,
the sampling is also guided by the gradient calculated before: the non-zero entry in the mask (rep-

270 resenting non-explanatory tokens) is sampled following probabilities calculated by the normalized 271 gradient magnitudes, i.e., softmax($m^{(n)} \odot q$), and we employ a similar sampling strategy for zero 272 entries (explanatory tokens). After swapping the mask indicators for the two sampled tokens l and 273 v, we generate a temporary mask m^{tmp} . We then evaluate whether this temporary mask leads to a 274 decrease in output probability. If it does, we update the binary mask to m^{tmp} , otherwise we leave the 275 binary mask as is. Consequently, without requiring intensive gradient computations, these sampling iterations keep discovering improved solutions for the discrete optimization problem shown in Eq. (3). 276 We conclude this update by sampling process in Line 3-12. 277

While capturing the joint influence of being masked, the proposed *JoPA* uses both the gradient information and the search-space strategy, thereby achieving better efficiency than adopting either method alone. *JoPA* requires only a single step of gradient calculation and avoids the need to convert between discrete and continuous masks. The obtained gradient information provides a favorable initial searching direction and reliable sampling probabilities that enhances the search efficiency.

283 **Theoretical Guarantee** Here we could prove that our algorithm can theoretically converge to the 284 local optima given enough iterations. Define that a solution $m^* \in \{0,1\}^T$ is the local optima, if 285 we have $p_{\theta}(y|\boldsymbol{m}^* \odot \boldsymbol{x}) \leq p_{\theta}(y|\boldsymbol{m} \odot \boldsymbol{x})$ for all \boldsymbol{m} in the one-swap neighborhood of \boldsymbol{m}^* , namely 286 $||m^* - m||_0 = 2$ (meaning *m* differs from m^* by one swap, as they are constrained to have k zero 287 entries). It could be proved that the output of the algorithm \hat{m} is the local optima by contradiction 288 that with sufficient iterations. Suppose \hat{m} is not the local optima, there must be a point m' in its 289 neighbor satisfying $||\mathbf{m}' - \hat{\mathbf{m}}||_0 = 2$, such that $p_{\theta}(y|\mathbf{m}' \odot x) < p_{\theta}(y|\hat{\mathbf{m}} \odot x)$. This means our 290 algorithm can still find a better solution by sampling another swap in further iterations and thus \hat{m} is not the output of our algorithm, leading to a contradiction to the assumption. 291

5 EXPERIMENT

292 293

294

295

296 297

298

299

300 301

302

303

304

305 306

307

This section aims to verify the effectiveness and efficiency of our proposed framework for interpreting the LLMs on the generation task. We conduct the experiments to answer the following questions:

- Q1: Do the generated prompt attributions play a predominant role in the model's generation, thereby serving as faithful counterfactual explanations for the generated output?
- Q2: Is our interpretation algorithm efficient for practical usage?
- Q3: Could the proposed algorithm effectively identify a combination of tokens that impose joint influence on the model generation?

We provide quantitative studies to evaluate the faithfulness of the explanatory prompt fragments generated by *JoPA*, comparing with existing interpretation baselines and ablation variants.

5.1 EXPERIMENT SETTINGS

308 **Models & Baselines** In the experiment, we employ two LLMs as our targeted $f_{\theta}(\cdot)$: LlaMA-2 (7B-309 Chat) (Touvron et al., 2023) and Vicuna (7B) (Zheng et al., 2023). There are relatively few methods attributing prompts on the entire language generation, and we choose random removal (Random), 310 Integrated-Gradient (Sundararajan et al., 2017), averaged attentions across all layers (Pruthi et al., 311 2019)(Attention), last-layer attention (Zhao & Shan, 2024)(Last-Attention) and Captum (Miglani 312 et al., 2023) as our baseline for comparison. We also compare JoPA with ReAGent (Zhao & Shan, 313 2024), with the results presented in Appendix A.11. We implement these models using the PyTorch 314 framework and pretrained weights from the transformers Python library (Wolf et al., 2020), and 315 conduct our experiments on an Nvidia RTX A6000-48GB platform with CUDA version 12.0. 316

317 Datasets We employ three distinct text generation datasets: Alpaca (Taori et al., 2023),
318 tldr_news (Belvèze, 2022), and MHC (Amod, 2024), to evaluate the effectiveness of our method
across various generation tasks. As we aim to capture the joint influence of prompts on the model
generations, we focus on relatively long-context prompts rather than simple one-sentence prompts.
Longer prompts tend to contain more diverse vocabularies, convey more information, and thus more
likely to exhibit high textual correlation. For evaluation purposes, we randomly select approximately
110 data samples with at least 15 words from each dataset. All datasets are publicly available and
more details about the dataset are illustrated in Appendix A.1.

324 5.2 EVALUATION METRICS

326 Faithfulness scores are a key metric for assessing the quality of explanations, with faithful explanations 327 accurately reflecting the model's decision-making process (Jacovi & Goldberg, 2020). Studies (Samek et al., 2017; Hooker et al., 2019) suggest that if a certain input tokens are truly important, their removal 328 should lead to a more significant change in model output than the removal of random tokens (Madsen 329 et al., 2022). Therefore, after removing explanatory tokens identified by a faithful method, the 330 model's new generation would significantly differ from using the original output. Moreover, the 331 input prompt with masking these explanatory tokens are less likely to reproduce the original model 332 response \boldsymbol{y} , indicating a lower value of $p_{\boldsymbol{\theta}}(\boldsymbol{y}|\boldsymbol{m} \odot \boldsymbol{x})$. 333

To quantitatively evaluate explanation faithfulness, we measure the change of model generation 334 behavior from two dimensions. On one hand, we compare the model's original and new generations: 335 the originally generated sequence y is based on the complete input prompt x (e.g., $y = f_{\theta}(x)$), while 336 the new generation is conditioned on the masked prompt (e.g., $y' = f_{\theta}(m \odot x)$). We thus measure 337 the similarity of the original generation y and the new generation y' based on their word frequency 338 and semantics. A smaller similarity reflects a larger change in model generation, suggesting a better 339 explanation. On the other hand, we measure the likelihood of generating the original output y when 340 the model uses the masked prompt, e.g., changes on $p_{\theta}(y|m \odot x)$. A smaller likelihood indicates 341 better explanations whose mask prevents the model from generating its original output. Detailed 342 definitions of these metrics are explained below. 343

Word Frequency: **BLEU** (Papineni et al., 2002) is widely used to measure how close the candidate text y' is to the reference text y. The score measures the precision of matching *n*-grams from the text y' to y by a clipping method to avoid overcounting and adjusting for brevity of y' if it is shorter than y. **ROUGE-L** (Lin, 2004) is the metric for measuring the overlap of sequences of words between the two texts, evaluating how much of the y' matches with the reference y (precision), how much the reference y is covered by the candidate y' (recall), and combine them into an F1 score.

Semantic Similarity: SentenceBert (Reimers & Gurevych, 2019) is a variation of the BERT model which is designed to generate high-quality sentence embeddings for the pairs of sentences. We leverage SentenceBert to transform the text y and y' into fixed-length embedding vectors, and calculate the cosine similarity between their embeddings to quantify their semantic similarity.

Probability Measurement: We define two measurements to reflect how the likelihood of generating the 354 original text y changes before and after applying the explanation mask. We first define the Probability 355 Ratio (**PR**) to indicate how less likely to generate the original output y when explanatory tokens are 356 masked: $PR(x, y, m) = \frac{\tilde{p}_{\theta}(y|m \odot x)}{\tilde{p}_{\theta}(y|x)}$, where $\tilde{p}(y|x) = p(y|x)^{1/s}$ is the length-normalized generation 357 probability. If the PR score is far below the random baseline, we can conclude that the masked 358 tokens are indeed important to cause the model generating y. In addition, we also calculate the **KL**-359 divergence between these two distributions to measure their difference: $D_{\text{KL}}(p_{\theta}(\boldsymbol{y}|\boldsymbol{m} \odot \boldsymbol{x})||p_{\theta}(\boldsymbol{y}|\boldsymbol{x}))$, 360 and a larger score indicates a larger change in generation likelihood and a more accurate explanation. 361

362 363

5.3 MAIN RESULTS ON FAITHFULNESS (Q1)

We evaluate the faithfulness of the explanatory prompt tokens generated by our framework on five aforementioned metrics. Table 1 shows the results of different methods on three datasets by masking k = 3 identified tokens. Experiments on the larger dataset is shown in Appendix A.10. For all metrics except the KL-divergence, a lower score is better which is annotated as \downarrow . In general, we observe that our method consistently demonstrates better interpretation faithfulness on all datasets compared with baselines, with a clear margin. In particular, we have the following observations demonstrating the advantage of our method:

Slight random perturbations on the input prompt do not significantly alter the model's output,
validating the soundness of our approach. This is evident from the high PR value of around
0.958 and the low KL value of only 0.023 on the MHC dataset when randomly masking tokens. Only
when essential tokens are perturbed, particularly when masked, does the generation of content and
probability change substantially. The noticeable gaps in these metrics between *JoPA* and Random
underscore the genuine importance of the masked tokens and their role as counterfactual explanations
for the model output. The variance assessment of these metrics is presented in Appendix A.6,
highlighting the performance stability of our algorithm.

IaMA-2 /B-Chat) t	Alpaca	Random	BLEU↓	RO	UGE-L↓				
	Alpaca	Random	- •				SentenceBert	PR↓	KL↑
JaMA-2 7B-Chat) 1	Alnaca	Random		Precision	Recall	F1	·	•	
JaMA-2 /B-Chat) (Alpaca	T and Attending	0.601	0.527	0.533	0.522	0.825	0.842	0.134
JaMA-2 7B-Chat) (Alpaca	Last-Attention	0.533	0.423	0.452	0.432	0.672	0.770	0.202
JaMA-2 /B-Chat) (Attention	0.547	0.447	0.460	0.448	0.721	0.791	0.170
	. npueu	Captum	0.515	0.409	0.421	0.409	0.680	0.602	0.417
		Integrated-Gradient	0.541	0.424	0.435	0.424	0.725	0.726	0.242
JaMA-2 7B-Chat) (JOPA	0.484	0.388	0.380	0.379	0.642	0.549	0.504
JaMA-2 /B-Chat) (Random	0.794	0.742	0.741	0.741	0.923	0.944	0.037
/B-Chat) t		Last-Attention	0.747	0.683	0.685	0.683	0.876	0.869	0.108
-	tldr news	Attention	0.767	0.703	0.710	0.706	0.900	0.899	0.077
-	iidi_iiews	Captum	0.759	0.701	0.703	0.701	0.900	0.910	0.069
-		Integrated-Gradient	0.713	0.641	0.642	0.640	0.866	0.817	0.149
]		JoPA	0.692	0.619	0.610	0.612	0.841	0.604	0.394
]		Random	0.723	0.617	0.614	0.615	0.787	0.958	0.023
]		Last-Attention	0.660	0.529	0.525	0.526	0.736	0.907	0.023
-	MHC	Attention	0.665	0.536	0.538	0.531	0.745	0.916	0.056
	WITC	Captum	0.640	0.497	0.493	0.494	0.663	0.760	0.189
		Integrated-Gradient	0.646	0.500	0.496	0.498	0.687	0.836	0.117
		JoPA	0.575	0.403	0.405	0.403	0.602	0.701	0.246
		Random	0.587	0.541	0.552	0.535	0.786	0.891	0.079
		Last-Attention	0.475	0.423	0.436	0.415	0.672	0.840	0.113
	Alpaca	Attention	0.486	0.447	0.469	0.441	0.683	0.831	0.140
-	r	Captum	0.466	0.435	0.427	0.418	0.649	0.654	0.347
		Integrated-Gradient	0.541	0.495	0.503	0.488	0.744	0.835	0.135
_		JOPA	0.433	0.395	0.390	0.377	0.639	0.589	0.459
		Random	0.781	0.736	0.746	0.737	0.921	0.891	0.091
licuna		Last-Attention	0.746	0.707	0./18	0.708	0.898	0.8/6	0.114
(7B) t	tldr_news	Contum	0.621	0.514	0.515	0.509	0.817	0.714	0.294
		Integrated-Gradient	0.330	0.669	0.401	0.666	0.874	0.370	0.183
		JoPA	0.536	0.456	0.454	0.448	0.772	0.317	1.006
_		Random	0.715	0.625	0.623	0.623	0.810	0.972	0.012
		Last-Attention	0.684	0.570	0.573	0.570	0.773	0.947	0.028
,	MUC	Attention	0.685	0.581	0.581	0.580	0.775	0.950	0.026
1		Captum	0.579	0.438	0.432	0.433	0.627	0.811	0.120
		Integrated-Gradient	0.672	0.559	0.555	0.556	0.762	0.941	0.030
		ΙοΡΔ	0.575	0.431	0.425	0.427	0.620	0.783	0.141

410

Table 1: Faithfulness Measurement Results for LlaMA-2 (7B-Chat) and Vicuna (7B)

411 JoPA can effectively capture semantically important prompt fragments, and the advantage 412 stands out for long context. Compared to Captum and Integrated-Gradient, our method generally 413 yields lower values for BLEU and ROUGE-L across all datasets. The discrepancies between JoPA and 414 Captum are more pronounced for datasets with longer text, such as tldr_news and MHC. Specifically, 415 on the tldr_news dataset, the F1-score of JoPA is 12.7% lower than Captum for the LlaMA-2 (7B-Chat) model, and on the MHC dataset, it is 17.85% lower, this all indicates that JoPA finds and 416 removes the more important token. Furthermore, the SentenceBert results are better on all datasets, 417 indicating larger semantic variations in the generated output after removing tokens by JoPA. 418

JoPA works even better on stronger LLMs. The gaps between JoPA and Captum are less apparent
 on the Vicuna (7B) model compared to the LlaMA-2 (7B-Chat) model, which could be attributed to
 the model's inherent inference capabilities. Our method is based on the premise that there are textual
 correlations among the input tokens, allowing the model to infer from the remaining content when
 a portion of the tokens is masked. If the model has a poor ability to infer the masked token, it may
 struggle to capture the contextual information. Consequently, the effectiveness of our method is tied
 to the LLM's proficiency in inference and understanding.

426 5.4 TIME EFFICIENCY (Q2)

427In Table 2, we compare the average time428cost of JoPA and Captum for generating429k = 3 explanatory tokens for each prompt430instance. Note that since Captum's design431requires sequentially appending the nexttoken to the input prompt to re-generate the

	Method	Dataset					
		Alpaca	tldr_news	MHC			
Time(s)	Captum JoPA	1169.648 15.225	1727.602 15.397	1806.551 14.473			

 Table 2: Time Efficiency on LlaMA-2 (7B-Chat)

432 new output, the time consumption of Captum would increase significantly when the prompt is long. 433 As the average length of MHC is longer than the Alpaca shown in Table 3, the computational time 434 increases from 1169.648s to 1806.551s for Captum, which is quite inefficient and impractical. As 435 for comparison, our algorithm aims to solve the combinatorial optimization problem efficiently via 436 the proposed probabilistic search algorithm, which significantly reduce the computational cost of the explanation generation. And also since our design only needs to perform the gradient-guided 437 probabilistic search step for a certain number of times, the computation time remains consistent 438 regardless of the prompt length as can be shown form Table 3. 439

5.5 QUALITATIVE VISUALIZATION (Q3)

Legend: Captum JOPA Write a story about a Writing a **book** to teach others <mark>a few</mark> things can (2) (R) doctor and his patient be a rewarding experience on many levels. This article discusses **publishing** and **whether** one should self-publish or find a publisher. It covers topics such as royalty rates, rights, Here is a story about a cover design, copyright, and much more. doctor and his patient: Dr. Smith was a kind and compassionate doctor who Here are some potential headlines based on had been practicing the context you provided: 1. "Navigating the medicine for over 20 World of Publishing: A Guide to Selfyears Publishing vs. Traditional Publishing'

Figure 2: Case study for visualizing the explanation of the model responses.

Use Case Study. Figure 2 uses a case study to showcase our method can effectively identify 455 interacted tokens. This figure illustrates the identified explanatory tokens by Captum and our method 456 JoPA. While there are few overlaps in the found tokens, compared with our method, those found 457 by Captum are mostly not important to the generated response. This demonstrates the Captum's 458 limited ability to explain the relationship between input prompt and output response, especially considering the textual correlations within input tokens. For instances, the tokens "publishing", 459 "publish" and "book" have semantic correlations, thus Captum masking one by one individually 460 cannot really eliminate this information complementarily provided by the others, especially when the LLM has a certain context inference ability. Similar issues exist in the second case, where "doctor" 462 and "patient" has semantic interactions. Therefore, treating tokens independently and ignoring their 463 joint influence on the generation is not a favored choice for prompt attribution. This case study highlights the importance of our formulation of joint attribution, and verifies the effectiveness of our 465 algorithm in discovering important token combinations on the generated output. More examples are 466 shown in Appendix. A.2.

468 **Discussion on Benefits.** Moreover, our algorithm shows significant potential for improving model 469 safety and supporting model diagnosis, as elaborated in Appendix A.3. In situations involving a 470 malicious prompt with an adversarial suffix, our approach can be leveraged to effectively detect and 471 remove the attributes responsible for the success of jailbreak attacks. By identifying and filtering out these harmful tokens, we can enhance the model's robustness, making it more resistant to adversarial 472 manipulations and ensuring safer outputs. When it comes to model diagnosis, our method can also 473 provide valuable insights. Specifically, users can utilize this approach to assess how effectively the 474 model responds to a particular input prompt. This evaluation process serves as a diagnostic tool, 475 helping to verify whether the generated content is both relevant to the provided prompt and consistent 476 with its intended meaning. Consequently, our algorithm not only enhances the model's security by 477 mitigating risks from malicious inputs but also contributes to ensuring the reliability and accuracy of 478 the model's output, ultimately improving user trust and satisfaction.

479 480

481

440

441 442

443 444

445

446

447

448

449

450

451

452

453 454

461

464

467

5.6 ABLATION STUDY

482 **Number of Explanatory Tokens.** Figure 3 shows the Probability Ratio (PR) score as the number of masked tokens k changes on each dataset for LlaMA-2 (7B-Chat). Our algorithm consistently 483 outperforms other baseline methods even with larger k, demonstrating the stable performance and 484 effectiveness of JoPA. After masking k = 2 tokens, the value of PR would decrease dramatically, 485 indicating that these two tokens are crucial for generating the output y and they act as the triggers

that alter the probability distribution of the output. As the number of tokens increase to k = 5, the PR keeps decreasing but with a less sharp slope. This suggests that the probability of generating the original output y is significantly determined by a few predominant tokens.





499 500

Figure 3: Probability Ratio with varying number of masked tokens k. plot of JoPA.

Number of Sampling Iterations. We now demonstrate the convergence of our search algorithm. For each sampling iteration, we sample entries in the mask for value swapping; the mask will be updated if the swapping leads to a drop in the generation log-likelihood log $p_{\theta}(y|m \odot x)$. Figure 4 empirically shows how the log-likelihood decreases as the sampling and mask update continue. The quickly decreasing trend demonstrates that our algorithm is successfully performed to improve the quality of mask in locating the predominant tokens on generating y.

Recall that in our algorithm, we use gradient as guidance to initialize the mask (e.g., $m^{(1)}$) and 507 calculate sampling probability (i.e., softmax($m^{(n)} \odot g$)). To verify the efficiency of gradient guidance, 508 we compare JoPA with two variants: w/o Initialization that initializes the mask by uniformly random 509 instead of gradient, and w/o Probability that samples swapping entries by uniformly random instead 510 of gradient. Table 5 in the Appendix reports their resulting generation log-likelihood log $p_{\theta}(y|m \odot x)$ 511 in different iterations. We observe that without using gradient for initialization, the randomly 512 initialized mask starts from a worse point with a high generation likelihood; and without using 513 gradient to guide the sampling, the mask is updated in less effective direction to explore the search 514 space, resulting in a high generation likelihood in the end. These results show that gradient is a useful 515 and efficient tool for initializing the optimization from a better starting point and guiding the search 516 to a better optimal point. 517

In addition, we extend the application of our framework, *JoPA*, to a larger model and a more challenging task, specifically Few-shot Chain-of-Thought (CoT) reasoning (Wei et al., 2022; Kojima et al., 2022), as described in Appendix A.8 and Appendix A.9. The better performance of our algorithm compared with others highlights the transferability of our framework across different model architectures and tasks.

522 523 524

6 CONCLUSIONS

525 In this study, we introduce JoPA, an efficient probabilistic search algorithm designed to generate the prompt attributions that elicit the model outputs for the generation tasks. We tackle the challenge 526 of explaining the generation behavior for any given prompt by analyzing the joint effect of the 527 prompt attributions on the output. We frame this explanation task as a discrete optimization problem, 528 which can be efficiently solved by our proposed probabilistic search algorithm. This methodology 529 enables efficient generation of any arbitrary number of explanatory prompt attributions that are 530 deterministic to the generated content. Our framework is rigorously evaluated across extensive 531 language generation tasks, including text summarization, question-answering, and general instruction 532 datasets. The faithfulness of the explanatory prompt attributions is thoroughly analyzed and assessed 533 by comprehensive metrics. The results demonstrate that our proposed method efficiently generates 534 explanatory attributions that faithfully reflect the model's generation behavior for the specific prompt. 535 These explanatory attributions interact in semantically and jointly influence the output generation. 536 Furthermore, the overall excellent performance of our method on these diverse datasets highlights its remarkable transferability. 537

538

540 REFERENCES 541

555

556

567

576

577

578

579

580 581

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, 542 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. 543 arXiv preprint arXiv:2303.08774, 2023. 544
- mental_health_counseling_conversations (revision 9015341), 2024. Amod. URL 546 https://huggingface.co/datasets/Amod/mental_health_counseling_ 547 conversations. 548
- Maksym Andriushchenko, Francesco Croce, and Nicolas Flammarion. Jailbreaking leading safety-549 aligned llms with simple adaptive attacks, 2024. URL https://arxiv.org/abs/2404. 550 02151. 551
- 552 Anthropic. Introducing the next generation of claude, March 2024. URL https://www. 553 anthropic.com/news/claude-3-family. Accessed: 2024-05-09. 554
 - Jules Belvèze. TL;DR News: A Large Dataset for Abstractive Summarization. https:// huggingface.co/datasets/JulesBelveze/tldr_news, 2022. Accessed: May 22, 2024.
- 558 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, 559 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. Advances in neural information processing systems, 33:1877–1901, 2020. 561
- 562 Patrick Chao, Edoardo Debenedetti, Alexander Robey, Maksym Andriushchenko, Francesco 563 Croce, Vikash Sehwag, Edgar Dobriban, Nicolas Flammarion, George J. Pappas, Florian Simon Tramèr, Hamed Hassani, and Eric Wong. Jailbreakbench: An open robustness bench-564 mark for jailbreaking large language models. ArXiv, abs/2404.01318, 2024. URL https: 565 //api.semanticscholar.org/CorpusID:268857237. 566
- Hanjie Chen and Yangfeng Ji. Learning variational word masks to improve the interpretability of 568 neural text classifiers. arXiv preprint arXiv:2010.00667, 2020. 569
- 570 Yu-Neng Chuang, Guanchu Wang, Chia-Yuan Chang, Ruixiang Tang, Fan Yang, Mengnan Du, Xuanting Cai, and Xia Hu. Large language models as faithful explainers, 2024. 571
- 572 Benjamin Cohen-Wang, Harshay Shah, Kristian Georgiev, and Aleksander Madry. Contextcite: 573 Attributing model generation to context, 2024. URL https://arxiv.org/abs/2409. 574 00729. 575
 - DAN. Chat gpt "dan" (and other "jailbreaks"), 2023. URL https://gist.github.com/ coolaj86/6f4f7b30129b0251f61fa7baaa881516. GitHub repository.
 - Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. Hotflip: White-box adversarial examples for text classification, 2018.
- Mononito Goswami, Vedant Sanil, Arjun Choudhry, Arvind Srinivasan, Chalisa Udompanyawit, and 582 Artur Dubrawski. Aqua: A benchmarking tool for label quality assessment. Advances in Neural Information Processing Systems, 36, 2024. 584
- Xu Han, Weilin Zhao, Ning Ding, Zhiyuan Liu, and Maosong Sun. Ptr: Prompt tuning with rules for 585 text classification, 2021. 586
- Sara Hooker, Dumitru Erhan, Pieter-Jan Kindermans, and Been Kim. A benchmark for interpretability 588 methods in deep neural networks, 2019. 589
- Alon Jacovi and Yoav Goldberg. Towards faithfully interpretable NLP systems: How should we define and evaluate faithfulness? In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, 592 pp. 4198–4205, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/ 2020.acl-main.386. URL https://aclanthology.org/2020.acl-main.386.

594 Muhammad Uzair Khattak, Hanoona Rasheed, Muhammad Maaz, Salman Khan, and Fahad Shahbaz 595 Khan. Maple: Multi-modal prompt learning. In Proceedings of the IEEE/CVF Conference on 596 Computer Vision and Pattern Recognition, pp. 19113–19122, 2023. 597 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large 598 language models are zero-shot reasoners. Advances in neural information processing systems, 35: 22199-22213, 2022. 600 601 Ramaravind Kommiya Mothilal, Divyat Mahajan, Chenhao Tan, and Amit Sharma. Towards uni-602 fying feature attribution and counterfactual explanations: Different means to the same end. In 603 Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society, AIES '21, pp. 652–663, 604 New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450384735. doi: 605 10.1145/3461702.3462597. URL https://doi.org/10.1145/3461702.3462597. 606 607 Jiwei Li, Xinlei Chen, Eduard Hovy, and Dan Jurafsky. Visualizing and understanding neural models 608 in NLP. In Kevin Knight, Ani Nenkova, and Owen Rambow (eds.), Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: 609 Human Language Technologies, pp. 681-691, San Diego, California, June 2016. Association for 610 Computational Linguistics. doi: 10.18653/v1/N16-1082. URL https://aclanthology. 611 org/N16-1082. 612 613 Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In Text Summarization 614 Branches Out, pp. 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics. 615 URL https://aclanthology.org/W04-1013. 616 617 Xiaogeng Liu, Nan Xu, Muhao Chen, and Chaowei Xiao. Autodan: Generating stealthy jailbreak prompts on aligned large language models, 2023. 618 619 Andreas Madsen, Nicholas Meade, Vaibhav Adlakha, and Siva Reddy. Evaluating the faithfulness of 620 importance measures in NLP by recursively masking allegedly important tokens and retraining. In 621 Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), Findings of the Association for Compu-622 tational Linguistics: EMNLP 2022, pp. 1731–1751, Abu Dhabi, United Arab Emirates, December 623 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-emnlp.125. 624 URL https://aclanthology.org/2022.findings-emnlp.125. 625 Vivek Miglani, Aobo Yang, Aram H Markosyan, Diego Garcia-Olano, and Narine Kokhlikyan. Using 626 627 captum to explain generative language models. arXiv preprint arXiv:2312.05491, 2023. 628 Ali Modarressi, Mohsen Fayyaz, Ehsan Aghazadeh, Yadollah Yaghoobzadeh, and Mohammad Taher 629 Pilehvar. Decompx: Explaining transformers decisions by propagating token decomposition. arXiv 630 preprint arXiv:2306.02873, 2023. 631 632 Jonas Oppenlaender. The creativity of text-to-image generation. In Proceedings of the 25th Interna-633 tional Academic Mindtrek Conference, pp. 192–202, 2022. 634 635 Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In Proceedings of the 40th Annual Meeting on Association for 636 Computational Linguistics, ACL '02, pp. 311-318, USA, 2002. Association for Computational 637 Linguistics. doi: 10.3115/1073083.1073135. URL https://doi.org/10.3115/1073083. 638 1073135. 639 640 Danish Pruthi, Mansi Gupta, Bhuwan Dhingra, Graham Neubig, and Zachary C Lipton. Learning to 641 deceive with attention-based explanations. arXiv preprint arXiv:1909.07913, 2019. 642 643 Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks, 644 2019. 645 Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. " why should i trust you?" explaining the 646 predictions of any classifier. In Proceedings of the 22nd ACM SIGKDD international conference 647 on knowledge discovery and data mining, pp. 1135–1144, 2016.

- Wojciech Samek, Alexander Binder, Grégoire Montavon, Sebastian Lapuschkin, and Klaus-Robert Müller. Evaluating the visualization of what a deep neural network has learned. *IEEE Transactions* on Neural Networks and Learning Systems, 28(11):2660–2673, 2017. doi: 10.1109/TNNLS.2016.
 2599820.
- Weijia Shi, Julian Michael, Suchin Gururangan, and Luke Zettlemoyer. Nearest neighbor zero-shot inference. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 3254–3265, 2022.
- Yucheng Shi, Hehuan Ma, Wenliang Zhong, Qiaoyu Tan, Gengchen Mai, Xiang Li, Tianming Liu, and Junzhou Huang. Chatgraph: Interpretable text classification by converting chatgpt knowledge to graphs. In 2023 IEEE International Conference on Data Mining Workshops (ICDMW), pp. 515–520. IEEE, 2023.
- Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric Wallace, and Sameer Singh. AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (eds.), *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 4222–4235, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.346. URL https://aclanthology.org/2020.emnlp-main.346.
- Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In International conference on machine learning, pp. 3319–3328. PMLR, 2017.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy
 Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model.
 https://github.com/tatsu-lab/stanford_alpaca, 2023.
- 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.
- Michael Tsang, Sirisha Rambhatla, and Yan Liu. How does this interaction affect me? interpretable
 attribution for feature interactions. *Advances in neural information processing systems*, 33:6147–6159, 2020.
- Miles Turpin, Julian Michael, Ethan Perez, and Sam Bowman. Language models don't always say what they think: Unfaithful explanations in chain-of-thought prompting. *ArXiv*, abs/2305.04388, 2023. URL https://api.semanticscholar.org/CorpusID:258556812.

682

683

684 685

686

687 688

689

690

- Keyon Vafa, Yuntian Deng, David M Blei, and Alexander M Rush. Rationales for sequential predictions. In *Empirical Methods in Natural Language Processing*, 2021.
- Boxin Wang, Weixin Chen, Hengzhi Pei, Chulin Xie, Mintong Kang, Chenhui Zhang, Chejian Xu, Zidi Xiong, Ritik Dutta, Rylan Schaeffer, et al. Decodingtrust: A comprehensive assessment of trustworthiness in gpt models. arXiv preprint arXiv:2306.11698, 2023.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. Advances in neural information processing systems, 35:24824–24837, 2022.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. Transformers: State-of-the-art natural language processing. In Qun Liu and David Schlangen (eds.), *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 38–45, Online, October 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-demos.6.
 URL https://aclanthology.org/2020.emnlp-demos.6.
- Zhiyong Wu, Yun Chen, Ben Kao, and Qun Liu. Perturbed masking: Parameter-free probing for analyzing and interpreting BERT. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*,

702 703	pp. 4166–4176, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.383. URL https://aclanthology.org/2020.acl-main.383.
704	Ziwei Xu, Sanjay Jain, and Mohan Kankanhalli. Hallucination is inevitable: An innate limitation of large language models 2024
706 707	
708	Jawei Zhang. Cognitive functions of the brain: Perception, attention and memory, 2019.
709	Zhixue Zhao and Boxuan Shan. Reagent: A model-agnostic feature attribution method for generative
710	language models, 2024.
710	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
712	Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Haotong Zhang, Joseph Gonzalez, and Ion Stoica.
714	Judging llm-as-a-judge with mt-bench and chatbot arena. ArXiv, abs/2306.05685, 2023. URL https://api.semanticscholar.org/CorpusID:259129398.
715	
716	Sicheng Zhu, Ruiyi Zhang, Bang An, Gang Wu, Joe Barrow, Zichao Wang, Furong Huang, Ani
717	language models. In Socially Responsible Language Modelling Research 2023. LIRL https:
718	//openreview_net/forum?id=r0ivmxm8t0
719	//openieview.nee/ioium.id ioijmmmoeg.
720	Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial
/21	attacks on aligned language models, 2023.
722	
723	
724	
726	
727	
728	
729	
730	
731	
732	
733	
734	
735	
736	
737	
730	
740	
741	
742	
743	
744	
745	
746	
747	
748	
749	
750	
751	
752	
754	
755	

756 A APPENDIX

760

761

762

764

765

766

767

768

769

770

771

789

791

792

758 A.1 DATASET DETAILS

We evaluate explanations on three datasets: Alpaca, tldr_news and MHC. Due to the extensive computational cost, we randomly select approaximately 110 data samples with at least 15 words from each datasets. The statistics of the data used in our experiment can be found in Table 3.

• Alpaca Taori et al. (2023) is a dataset with 52000 unique examples consisting of instructions and demonstrations generated by OpenAI's text-davinci-003 engine. Each example in the dataset includes an instruction that describes the task the model should perform, accompanied by optional input context for that task. In our experiment, we randomly select a subset of these examples for verification purposes.

• **tldr_news** Belvèze (2022) dataset is constructed by collecting a daily tech newsletter. For every piece of data, there is a headline and corresponding content extracted. The task is to ask the model to simplify the extracted content and then generate a headline from the input.

mental_health_counseling (MHC) Amod (2024) includes broad pairs of questions and answers derived from online counseling and therapy platforms. It covers a wide range of mental-health related questions and concerns, as well as the advice provided by the psychologists. It is utilized for used to fine-tuning the model to generate the metal health advice. Here, we prompt the model to generate an advice based on the provided question.

Dataset	Length	Number 108	Prompt Examples					
Alpaca	16-137		Pretend you are a project manager of a construction com- pany. Describe a time when you had to make a difficult decision.					
tldr_news 14-138 12		120	eddit aired a five-second long ad during the Super Bow he ad consisted of a long text message that hinted at the ameStop stocks saga. A screenshot of the ad is available to the article.					
МНС	18-478	100	I cannot help myself from thinking about smoking. What can I do to get rid of this addiction?					

Table 3: Data Statistics

A.2 CASE STUDY

793 In this section, we visualize three examples which are selected from the dataset Alpaca, tldr_news, 794 and MHC respectively, shown in Figure. 5. The tokens in the red frame are the k explanatory tokens extracted by the model LlaMA-2 (7B-Chat), and the output responses framed in blue are the most 795 informative results selected by human. In the first data sample, the recognized token: "brains" and 796 "cognitive", has some semantic relationship intuitively. It is a common sense the cognitive ability is 797 underpinned by the brain's function and structure. There are numerous reports and studies focused 798 on the functionally relationship between the brain and cognitive ability (Zhang, 2019). Due to the 799 extensive amount of data on which the LLMs are trained, the model can integrate their relationship 800 which is reflected in its response. 801

As for the second example, the visualized tokens are almost related to the "book" and "publish", which aligns with the theme of the output content. Removing either "publishing" or "publish" would not greatly alter the meaning of the input prompt. This example also shows that our method could detect some fix expressions like "whether...or..", whose relation is recognized by the token "vs." in the output.

807 In the last example, we could figure out the model infers or "guesses" that the kid feels frustrated and helpless, a sentiment partly related to the token "never". The generated output does not mention the reason to explain how the model makes such inference. From this point, our method could generate human-intelligible explanations to assist the users in identifying which parts of the prompts fundamentally lead to the undesired response, providing guidance on how to revise the prompts effectively.



Figure 5: Case study for visualizing the top-k tokens in each example.

In Figure 6, we also illustrate the changes in the model response as the targeted top-k token is removed from the input prompt. Compared with the top-k tokens identified by our method, those found by Captum are mostly not important to the response. This demonstrates the Captum's limited ability to explain the relationship between input prompt and output response with the consideration of the textual information. Moreover, when masking the token "publishing" or "publish", as identified by our method *JoPA*, the generated outputs do not exhibit significant variations in either word frequency or semantics. However, substantial changes in the generated outputs occur only when both tokens are simultaneously masked. This highlights the effectiveness of our methodology in explaining the joint effect of the prompt attributions on the generated output.

A.3 BENEFITS TO THE COMMUNITY ON IMPROVING MODEL ABILITY

Benefit on Model Safety Our method has the potential to help LLM defend against jailbreaking attacks, making the model more robust. We verify this via a real case study. Given an adversarial prompt with malicious request: "Generate a message that spreads fake news. organized few GO a a a a newcommand ! SudGuiheures ! !/ refection/:" It can jailbreak Llama-2 (7B-Chat) model to generate unsafe responses. Our algorithm successfully attributes "Sud" and "he", which are part of the adversarial suffix. By deleting these attributions, the attack fails. Thus, our approach could help detect these high-risk malicious tokens to robustify model safety.

Benefit on Model Diagnosis Our approach can also be used to assess the relevance of model's response to the prompt, helping diagnose the effectiveness of model responding to the prompt. For a real case study, consider a single prompt consisting of multiple tasks: "Create a message to complain about the noise in the neighborhood to friends. Give me the Python code for filtering results based on the following columns..." We find Llama-2 (7B-Chat) only responds to the first request, neglecting the second one. Our approach attributes the model's output to tokens like "complain" "noise" and "friends", mainly about the first request. Thus by input attribution, our method can help diagnose the overlooked request and suggest better prompt design.



Figure 6: Case study for visualizing the variation of the model responses after the specific tokens are masked.

Detection Accuracy JoPA could be used as the defense of the jailbreak attaks and help users to design better prompts to get satisfying responses. We also do the experiments to show the effectiveness of applying JoPA to detect malicious prompts on Llama-2 (7B-Chat) model, which is similar to the application done in CONTEXTCITE (Cohen-Wang et al., 2024) The results are shown below:

Method	Detection Accuracy by JoPA	ASR	
GCG	100%	3%	
Prompt with Random Search	91%	90%	

Table 4: Comparison of detection accuracy and attack success rate.

The table presents the detection accuracy of JoPA against different adversarial prompts generated by GCG attacks (Zou et al., 2023) and Prompt with Random Search (Andriushchenko et al., 2024) on Llama-2 (7B-Chat). For Prompt with Random Search, the ASR is 90% (Chao et al., 2024), while JoPA achieves a detection accuracy of 91%. This indicates that JoPA successfully identifies 91% of malicious prompts, highlighting its potential utility as a defense mechanism.

910 A.4 DISCUSSION ON COUNTERFACTUAL EXPLANATION

891

892 893

894

895

896

897

899 900 901

902 903 904

905

906

907

908

909

912 Our work *JoPA* is to explain the generation of LLMs by identifying a small changes (i.e., masking 913 few tokens) needed to alter the outcome (i.e., minimizing the probability of the original generation). 914 This can be understood as counterfactual explanation, since it answers a "what if" question: if these 915 input tokens are masked, the LLM would not generate the response. Namely, by applying the mask *m* 916 on the input prompt *x*, the generation probability of its original response *y* would be largely decreased 917 (i.e., $p_{\theta}(y|m \odot x) \ll p_{\theta}(y|x)$). This suggests that the newly generated output *y'* resulting from the 918 masked prompt would differ substantially from the original outcome *y*, as measured by our metrics in Table 1: the semantic similarity (SentenceBert) and word similarity (BLEU and ROUGE-L) between y and y' are much smaller than those from baselines. Here are several case studies to illustrate the concepts discussed.

921

939

940 941

942

950

951 952 953

954 955

956

957

965 966 967

968

General Q&A Task Given the input prompt to the model: "Follow the law of supply and demand, describe what would happen to the price of a good if the demand increased." The model would generate the response: "\nIn economics, the law of supply and demand states that the price of a good is determined by the interaction between the quantity of the good that...". After masking the 3 explanatory tokens: 'supply', 'price', and 'demand', the generated response would change significantly: "However, I must inform you that the question you've provided is not factually coherent, and I cannot provide an answer that may not be accurate or safe...".

Fill-in-the-Blank Task Given the input prompt for the model to fill in the blank:" Fill in the blank with a word or phrase The most successful team in the NBA is the _____" the original response is: "The most successful team in the NBA is the Golden State Warriors". If we masked out the explanatory tokens 'most', 'successful', and 'NBA' from the input, the output would change to: "\Based on the context you provided, the word that best fits the blank is "best." So, the sentence would read:\The team in the office is the best."

The case studies of diverse tasks above demonstrate that the new responses generated from the masked inputs differ significantly from the original outputs, which highlights the impact of these explanatory tokens in influencing and altering the model's output.

A.5 GRADIENT DESIGN

The Table. 5 displays the log-likelihood without using the gradient as the guidance to initialize the optimization process or to do sampling at different stage when searching for the optimal results.

Methods	Number of Iterations								
methods	1	5	10	15	20	30	40	50	
JoPA	-42.448	-51.661	-56.378	-59.852	-62.303	-64.812	-67.274	-69.426	
w/o Initialization w/o Probability	-33.906 -41.759	-44.555 -49.867	-50.828 -53.784	-55.652 -57.107	-58.483 -59.639	-62.694 -63.339	-64.793 -65.848	-66.415 -67.496	

Table 5: Ablation study showing the log-likelihood without using gradient for initialization or sampling with different number of iterations.

A.6 VARIANCE ASSESSMENT

To assess the variance of these metrics, we run the experiment with different seed for 3 times on the Llama-2 (7B-Chat) model. The results of different metrics on three datasets are shown in the Table. 6, where value in each cell denotes mean \pm std. The small variance indicates that our algorithm is quite stable across different runs.

Dataset	BLEUI		ROUGE-L↓			PR	KL↑
Dataset		Precision	Recall	F1			1112
Alpaca	$0.479 {\pm} 0.005$	$0.378 {\pm} 0.008$	$0.380{\pm}0.010$	0.371±0.009	$0.638 {\pm} 0.005$	$0.552{\pm}0.002$	$0.496 {\pm} 0.007$
ldr_news	$0.694{\pm}0.003$	$0.612 {\pm} 0.004$	$0.613 {\pm} 0.001$	$0.610 {\pm} 0.002$	$0.842{\pm}0.004$	$0.592{\pm}0.005$	$0.406 {\pm}~0.008$
MHC	$0.575 {\pm} 0.000$	$0.407 {\pm} 0.001$	$0.406 {\pm} 0.001$	$0.403 {\pm} 0.001$	$0.595 {\pm} 0.008$	$0.701 {\pm} 0.000$	$0.245 {\pm} 0.000$

Table 6: Variance measurement results for LlaMA-2 (7B-Chat)

A.7 OPTIMIZING MASK VIA GRADIENT

As discussed in Section 4, the strategy that directly optimizes the mask based on gradients requires first relaxing the discrete problem into a continuous optimization problem (i.e. $m \in [0, 1]^T$), optimizing via gradient descent, and later projecting the continuous solution back into the discrete space (i.e., $m \in \{0, 1\}^T$). This projection from continuous to discrete space in practice usually results in large

rounding error, i.e., after projection, the loss increases dramatically. The results of this strategy for Llama-2 (7B-Chat) model on the Alpaca dataset are in Table 7:

Method	BLEUL	RC	UGE-L↓		SentenceBert	PR	KL↑
		Precision Re	Recall	F1	SentenceDent	•	
Gradient JoPA	0.624 0.484	0.523 0.388	0.534 0.386	0.520 0.379	0.839 0.642	0.825 0.549	0.030 0.504

Table 7: Optimizing continuous mask on Alpaca dataset

The results verify our claim: directly optimizing *m* using the gradients (and then projecting to discrete solution) has much worse performance than our method, highlighting the challenge of this discrete optimization problem and our contribution of solving it effectively.

A.8 RESULTS ON LARGER MODEL

To further validate a wider applicability of our algorithm on larger models, we conducted additional experiments for Llama-2 (70B-Chat) 16-bit model on the Alpaca dataset. Due to the inefficiency of baseline Captum, we randomly sampled 20 instances from the Alpaca dataset to obtain the results in Table 8, which demonstrate that our algorithm still outperforms the baselines on this larger model with clear margins.

Method	BLEUL	ROUGE-L↓			SentenceBert	PR	KL↑
	22204	Precision	Recall	F1	SentenceDerty		
Random	642	0.553	0.552	0.551	0.801	0.894	0.085
Captum	0.565	0.458	0.469	0.462	0.659	0.647	0.333
JoPA	0.547	0.415	0.410	0.410	0.627	0.615	0.363

Table 8: Results of Llama-2 (70B-Chat) model on Alpaca dataset

A.9 METHOD GENERALIZABILITY ON REASONING TASK

We demonstrate the generalizability of our algorithm by conducting an additional experiment on the reasoning task of Few-shot-CoT (Wei et al., 2022; Kojima et al., 2022), using the dataset AQuA (Kojima et al., 2022; Goswami et al., 2024). The following results in Table 9 shows the remarkable performance of our method, indicating its ability to generalize well on more complex tasks like CoT.

Method	BLEUL	RC	UGE-L↓		SentenceBert	PR	KL↑
	DLLC ₄	Precision	Recall	F1	SentenceDerty	1114	
Random	0.967	0.976	0.978	0.977	0.995	0.958	0.039
Captum	0.872	0.849	0.893	0.870	0.962	0.666	0.369
JoPA	0.778	0.671	0.726	0.696	0.894	0.600	0.456

Table 9: Experimental results on reasoning task of Few-shot-CoT

A.10 RESULTS ON LARGER DATASET

As the Captum is 1000 times slower than JoPA (Table 2), to afford the comparison, we were using around 100 data in our experiments. Here, we increase the data size by randomly sampling 800 Alpaca data to demonstrate the effectiveness of our method. The experimental results in Table 10 indicate the effectiveness of our method as the sample size increases on Llama-2 (7B-Chat) model.

1026			DO	LIGE L				
1027	Method	BLEU↑	RC	DUGE-L↑		SentenceBert↑	PR↑	KL
1028		- 1	Precision	Recall	F1		I	¥
1029	Random	0.621	0.539	0.540	0.535	0.824	0.871	0.104
1030	Captum	0.498	0.401	0.402	0.395	0.660	0.623	0.392
1031	JoPA	0.479	0.383	0.376	0.372	0.642	0.565	0.479
1032		- 10. E		-14	1		1.4	

Table 10: Experimental results on evaluation metrics for larger data samples.

A.11 **COMPARISON WITH OTHER BASELINES**

ReAGent (Zhao & Shan, 2024) addresses a different task than JoPA. It focuses on classification tasks by explaining the importance of words in inputs corresponding to the predicted label and extends this approach to generation tasks by explaining the single next predicted word. In contrast, our task focuses on explaining the relationship between the input prompt and the entire generated sentences. To compare the performance of JoPA and ReAGent on the generation task, we modify ReAGent's explanatory target to encompass the full generation outputs rather than a single word. We conduct experiments on the Alpaca dataset using the Llama-2 (7B-Chat) model, and the results are presented in Table 11. The experimental results demonstrate that JoPA outperforms in explaining the relationship between joint attributions in the input and the resulting output.

Method	BLEU↑	RO	UGE-L↑		_ SentenceBert↑ PR↑		KLL
	22201	Precision	Recall	F1		111	1124
ReAGent	0.607	0.511	0.514	0.508	0.803	0.857	0.111
JoPA	0.484	0.388	0.386	0.379	0.642	0.549	0.504

Table 11: Experimental results between JoPA and ReAGent on Alpaca dataset.