# Usable XAI: 10 Strategies Towards Exploiting Explainability in the LLM Era

XUANSHENG WU\*, University of Georgia, USA

HAIYAN ZHAO\*, New Jersey Institute of Technology, USA

YAOCHEN ZHU\*, University of Virginia, USA

YUCHENG SHI\*, University of Georgia, USA

FAN YANG, Wake Forest University, USA

LIJIE HU, King Abdullah University of Science and Technology, Saudi Arabia

TIANMING LIU, University of Georgia, USA

XIAOMING ZHAI, University of Georgia, USA

WENLIN YAO, Amazon, USA

JUNDONG LI, University of Virginia, USA

MENGNAN DU, New Jersey Institute of Technology, USA

NINGHAO LIU, University of Georgia, USA

Explainable AI (XAI) refers to techniques that provide human-understandable insights into the workings of AI models. Recently, the focus of XAI is being extended toward explaining Large Language Models (LLMs). This extension calls for a significant transformation in the XAI methodologies for two reasons. First, many existing XAI methods cannot be directly applied to LLMs due to their complexity and advanced capabilities. Second, as LLMs are increasingly deployed in diverse applications, the role of XAI shifts from merely opening the "black box" to actively enhancing the productivity and applicability of LLMs in real-world settings. Meanwhile, the conversation and generation abilities of LLMs can reciprocally enhance XAI. Therefore, in this paper, we introduce Usable XAI in the context of LLMs by analyzing (1) how XAI can explain and improve LLM-based AI systems and (2) how XAI techniques can be improved by using LLMs. We introduce 10 strategies, introducing the key techniques for each and discussing their associated challenges. We also provide case studies to demonstrate how to obtain and leverage explanations. The code used in this paper can be found at: https://github.com/JacksonWuxs/UsableXAI_LLM.

CCS Concepts: • **Computing methodologies** → **Natural language generation**; **Machine learning**.

Additional Key Words and Phrases: Large Language Models, Explainable AI (XAI), Usability

---

\*The authors contributed equally to this research.

Authors' addresses: Xuansheng Wu, xuansheng.wu@uga.edu, University of Georgia, Athens, Georgia, USA; Haiyan Zhao, hz54@njit.edu, New Jersey Institute of Technology, Newark, New Jersey, USA; Yaochen Zhu, uqp4qh@virginia.edu, University of Virginia, Charlottesville, Virginia, USA; Yucheng Shi, yucheng.shi@uga.edu, University of Georgia, Athens, Georgia, USA; Fan Yang, yangfan@wfu.edu, Wake Forest University, Winston-Salem, North Carolina, USA; Lijie Hu, lijie.hu@kaust.edu.sa, King Abdullah University of Science and Technology, Saudi Arabia; Tianming Liu, tliu@uga.edu, University of Georgia, Athens, Georgia, USA; Xiaoming Zhai, xiaoming.zhai@uga.edu, University of Georgia, Athens, Georgia, USA; Wenlin Yao, ywenlin@amazon.com, Amazon, Seattle, Washington, USA; Jundong Li, jundong@virginia.edu, University of Virginia, Charlottesville, Virginia, USA; Mengnan Du, mengnan.du@njit.edu, New Jersey Institute of Technology, Athens, Georgia, USA; Ninghao Liu, ninghao.liu@uga.edu, University of Georgia, Athens, Georgia, USA.

---

## 1 INTRODUCTION

Explainability holds great promise for understanding machine learning models and providing directions for improvement. In practice, users have high expectations for model explainability:

> *1. Through explanation, can we know if a model works properly?*

> *2. Does explainability tell us how to develop better models?*

First, explanations are expected to illuminate whether a model operates as humans expect. For example, does the model leverage reliable evidence and domain knowledge in its decision making? Does the model contain bias and discrimination? Does the model show any vulnerabilities to potential attacks? Will the model output harmful information? Second, in recognition of model imperfections, we aspire for explainability to inform the development of better models. For example, how to adjust the behaviors of a model if we find that it uses unreliable or unreasonable features in making predictions? Can we improve the performance of a model by aligning its behavior with human preferences?

Therefore, the question arises: **Have these expectations been met?** Since the rise of deep learning, the body of literature on Explainable AI (XAI) has expanded rapidly to improve model transparency [63, 64, 206, 233, 260], encompassing a wide array of methods customized for different data modalities, including visual [332], textual [54], graph [326], and time-series data [342]. In addition, some offer reviews of general principles and initiate discussions on evaluating the faithfulness of explanations [62, 308]. **Despite the progress, the last mile of XAI – making use of explanations – has not received the attention it merits.** In many cases, we seem to be satisfied with just acquiring explanations and their associated visualizations, sometimes followed by qualitative discussions of the model's strengths and failure cases. However, how to quantify model properties (e.g., fairness, security, rationality) and improve models by utilizing explanations remains a difficult task.

While the opacity issues have not yet been fully resolved for traditional deep models (e.g., multi-layer perceptrons, convolutional and recurrent neural networks), the recent advancements of Large Language Models (LLMs) [1, 31, 47, 263] appear to have exacerbated the challenge we are facing. Firstly, LLMs typically possess a significantly larger model size and a greater number of parameters. This increased model complexity intensifies the difficulty of explaining their inner workings. Second, unlike traditional ML models that primarily focus on low-level pattern recognition tasks such as classification and parsing, LLMs can handle more complex tasks [340] such as generation, reasoning, and question answering. Understanding the exclusive abilities of LLMs presents novel challenges for XAI techniques. Considering the transformative impact of LLMs across various applications, ensuring the explainability and ethical use of LLMs has become an imminent and pressing need. Meanwhile, the emerging capabilities of LLMs also present new opportunities for XAI research. Their human-like communication and commonsense reasoning skills offer prospects for achieving explainability in ways that could potentially augment or replace human involvement.

**Defining "Usable XAI".** In light of the above considerations, in the context of LLMs, we define Usable XAI which includes two aspects. *(1) Leveraging XAI to improve LLMs and AI Systems.* Beyond just producing explanations of LLMs, we explore whether these explanations can pinpoint issues for model debugging or improve the overall performance of LLMs or AI models, such as accuracy, controllability, and trustworthiness. *(2) Enhancing the usability of XAI through LLMs.* The human-like communication ability of LLMs can enhance model explanations in terms of user-friendliness, by converting the numerical values into understandable language. Also, the commonsense knowledge stored in LLMs can significantly boost the practicality of existing XAI frameworks, by playing the role of humans and alleviating the need for real human involvement in AI workflows.
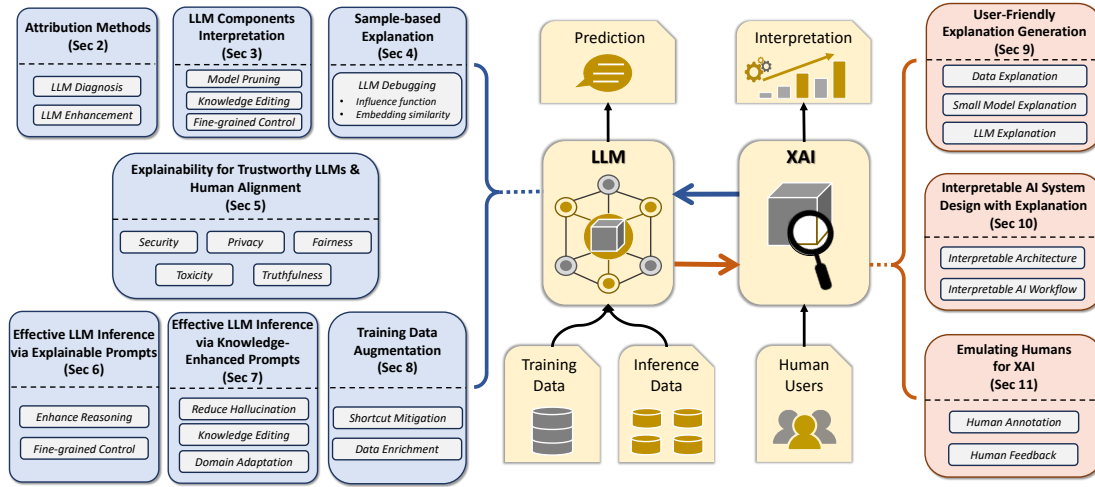
Fig. 1. The contributions and outline of this paper. We define Usable XAI in the context of LLMs with seven strategies of enhancing LLMs with XAI, and three strategies of enhancing XAI with LLMs.

**Contribution of this paper.** In this paper, we investigate 10 strategies towards usable XAI techniques in the context of LLMs, as shown in Figure 1. For each strategy, we also explore the open challenges that require further investigation in future work.

- **Leveraging XAI to improve LLMs and AI Systems** We introduce how interpretation can be utilized to enhance LLMs and AI models. We follow the LLM development pipeline in reverse order—spanning post-hoc analysis, prompt improvement, and data enhancement. First, we investigate how post-hoc explanations could be utilized to diagnose and enhance LLMs. We study three types of post-hoc explanation methods, targeting LLM predictions (**Section 2**), LLM components (**Section 3**), and training samples (**Section 4**), respectively, followed by how explanations could be leveraged to scrutinize and boost LLM trustworthiness (**Section 5**), including security, fairness, toxicity, and truthfulness. Second, we discuss two strategies for designing explainable prompts to guide LLM reasoning: Chain-of-Thought prompts (**Section 6**) and knowledge-enhanced prompts (**Section 7**). Third, we introduce the use of LLM explanations to augment training data toward improving AI models (**Section 8**).

- **Enhancing XAI Usability through LLMs** In this part, we investigate strategies for leveraging the advanced capabilities of LLMs to address the challenges in traditional XAI domains, thus enhancing the usability of XAI in practice. First, we examine ways to enhance the user-friendliness of explanations through the generative capabilities of LLMs (**Section 9**). Second, we introduce how to automate the design of interpretable AI workflows by leveraging the planning abilities of LLMs (**Section 10**). Third, we introduce how to facilitate the evaluation of XAI methods by utilizing the unique property of LLMs in emulating human cognition processes (**Section 11**).

**Differences between this paper and existing surveys.** Many surveys have been conducted to examine Explainable AI [63, 64, 206, 260]. This paper differs from existing work as we focus on explanation methods for large language models. Meanwhile, different from the existing survey [76, 338] that reviews explanation methods for LLMs, our paper places an emphasis on the XAI usability in LLM studies. To the best of our knowledge, the most related paper to our survey is [185], which also discusses several aspects where explanations can improve LLM performance. Nevertheless, this light-weight investigation lacks a thorough examination of XAI methods (e.g., sample-based explanation, interpretable workflows,

explainable prompts) and how LLMs can benefit existing XAI frameworks (e.g., data augmentation, improving user-friendliness, XAI evaluation). Finally, our paper contributes further by providing detailed case studies and open-sourced codes, fostering future research in applying explanations effectively within the LLM context.

## 2  LLM DIAGNOSIS AND ENHANCEMENT VIA ATTRIBUTION METHODS

This section introduces attribution methods as post-hoc explanations for LLMs, and how we can discover model defects with attribution scores. We start by reviewing existing attribution techniques and discussing which methods are still suitable for explaining LLMs. After that, we explore cases that apply attribution methods to assess LLM-generated output quality. Finally, we discuss future work of designing novel post-hoc explanations for LLMs.

### 2.1  Attribution Methods for LLMs

The attribution-based explanation quantifies the importance of each input feature that contributes to the prediction. Given a language model $f$ with a prediction $\hat{y} = f(x)$ according to the $N$-words input prompt $x$, the explainer $g$ assesses the influence of input words in $x$ as: $a = g(x, \hat{y}, f)$, where $a$ stores the importance scores of words. In *text classification*, $\hat{y}$ denotes a specific class label. In *text generation*, $\hat{y}$ represents a varying length of generated text. The primary distinction between them is that: classification is limited to a specific set of predictions, while generation encompasses an endless array of possibilities. In the following, we review related works based on the tasks to which they are applicable.

*2.1.1  Attribution Methods for Classification.* Common attribution methods [64, 206] for deep models include gradient-based methods, perturbation-based methods, surrogate methods, and decomposition methods.

- **Perturbation-based explanation** assesses the importance of input features by perturbing them and monitoring changes in prediction confidence, i.e., $a_n = p(\hat{y}|x) - p(\hat{y}|x_{/n})$, where $x_{/n}$ refers to the input sequence with the $n$-th feature being perturbed or deleted [156, 158, 302]. This approach has limitations, particularly in its assumption that features are independent, which is not always the case with textual data. Additionally, it is computationally intensive for explaining LLMs, requiring $N$ inferences for an input of $N$ words.

- **Gradient-based explanation** offers a computationally efficient approach for estimating model sensitivity to input features based on gradients $\frac{\partial p(\hat{y}|x)}{\partial \mathbf{x}_n}$, where $\mathbf{x}_n$ refers to the embedding of word $x_n$ [70, 141, 156, 202, 245]. This approach only requires a single inference and one backpropagation pass, achieving an efficient way to estimate the input word sensitivities. However, it is still expensive to compute a backward pass on such giant LLMs.

- **Surrogate-based explanation** understands the target model $f$ by constructing a simpler model $g$ trained on $\mathcal{D}(x, \hat{y}) = \{(\widetilde{x}_k, \widetilde{y}_k)\}_{k=1}^{K}$, where $\widetilde{x}_k$ is usually obtained by perturbing $x$, and $\widetilde{y}_k = f(\widetilde{x}_k)$ [145, 183, 229]. The surrogate model $g$, ranging from basic linear models to sophisticated decision trees, serves as a proxy to approximate the decision boundary of $f$ for the target instance $(x, \hat{y})$. Nevertheless, a significant limitation of this approach is its repeated interactions with the target model, which leads to huge computing requirements. In addition, it is impractical to find a simple and explainable model $g$ that can estimate the advanced LLMs well.

- **Decomposition-based explanation** assigns linearly additive relevance scores to inputs, effectively breaking down the model's prediction [204, 205, 265, 266, 303]. However, a primary challenge in implementing such an approach is the need for tailored decomposition strategies to accommodate different model architectures. This challenge poses a limitation on the universal applicability of decomposition methods for general-purpose interpretation.

*2.1.2  Attribution Methods for Generation.* With advancements in generative AI, interpreting why an output is generated has become increasingly important. The explanation of LLM generation can be defined as attributing the overall

Table 1. Time complexity analysis on different attribution methods for the generative task.

| Method | Forward | Backward | Notes |
|---|---|---|---|
| Mask Perturbation | $O(N)$ | 0 | - |
| Gradient×Input | $O(1)$ | $O(M)$ | - |
| Integrated Gradients | $O(N_{step})$ | $O(N_{step} \cdot M)$ | $N_{step}$ is the number of steps for integrating gradients. |
| LIME | $O(N_{aug})$ | 0 | $N_{aug}$ is the number of augmented samples. |
| SHAP | $O(2^N)$ | 0 | - |

confidence $p(\hat{y}|x)$ to the input $x$, where $\hat{y}$ denotes the generated response $\hat{y} = [\hat{y}_1, ..., \hat{y}_M]$ with $M$ words. An approach is to treat the text generation process as a sequence of word-level classification tasks. This perspective allows for the application of **existing classification-based explanation techniques** to assess the influence of each input word $x_n$ in relation to each output word $\hat{y}_m$. In particular, Wu et al. [298] proposes a gradient-based method for input-output attribution. Specifically, the importance score $a_{n,m}$ defined between input token $x_n$ to output token $\hat{y}_m$ can be estimated as:

$$a_{n,m} = p(\hat{y}_m|x, \hat{y}_{1:(m-1)}) - p(\hat{y}_m|x_{/n}, \hat{y}_{1:(m-1)}) \approx \frac{\partial f(\hat{y}_m|x, \hat{y}_{1:(m-1)})}{\partial \mathbf{E}[x_n]} \cdot \mathbf{E}[x_n], \tag{1}$$

where $\hat{y}_{1:(m-1)}$ denotes the first $m-1$ tokens of response, $x_{/n}$ means removing the $n$-th input token from $x$, and $\mathbf{E}[x_n]$ is the embedding of token $x_n$. The total contribution of each input word $x_n$ can be estimated by averaging the attributions $a_{n,m}$ for $m = 1, ..., M$ [241]. Figure 2 (left) plots an example, where each index in the Y-axis refers to an input token, while that in the X-axis is an output token. A higher attribution score is brighter. Given that not all output tokens warrant interpretation, Qi et al. [220] outlines a two-step procedure where informative tokens are first identified, and then attribution is made to these tokens using gradient-based or other attribution techniques [317]. It also develops a contrastive explanation method, which aims to explain why the model predicted one target token $\hat{y}_m$ over an alternative token $\hat{y}'_m$. Furthermore, a Python library has been developed to incorporate existing attribution methods for sequential generation of language models [236].

**Time complexity** is a critical consideration when applying attribution methods to language generation models like LLMs. In particular, the perturbation-based and gradient-based approaches may be limited by their huge demand for computing resources, while the surrogate-based and decomposition-based strategies are restricted by their limited generalizability across different model architectures. To emphasize the computing costs of applying attribution-based approaches for generative LLMs, we present the time complexity of several representative methods in Table 1, by assuming the model takes $N$ words as input and predicts $M$ words as output. It shows that existing attribution methods require a large number of forward or backward operations, emphasizing the exploration of efficient methods.

Another category of attribution approaches leverages the **generative capabilities of LLMs** to provide citations to one or more text passages that support the output they generate. The attribution can be obtained by explicitly prompting LLMs to quote the curated sources of information [292]. The cited contents can originate either from **external knowledge sources** or from the **LLMs' own memory**. For instance, Sun et al. [253] introduces RECITE, which first retrieves one or more relevant passages from the LLMs' memory via sampling and then generates the final response based on these passages. Conversely, many systems assume the presence of an external retrieval corpus. In these cases, the LLM output includes both an answer string and a reference to a short segment of text from the corpus that substantiates the answer [13, 28, 81]. Nevertheless, LLMs may hallucinate and produce fabricated information. To address this, various benchmarks have been developed to assess the faithfulness of LLM-generated attributions, i.e., whether the generated statements are fully supported by the cited references. For details on evaluation methods, we refer readers to Section 11.
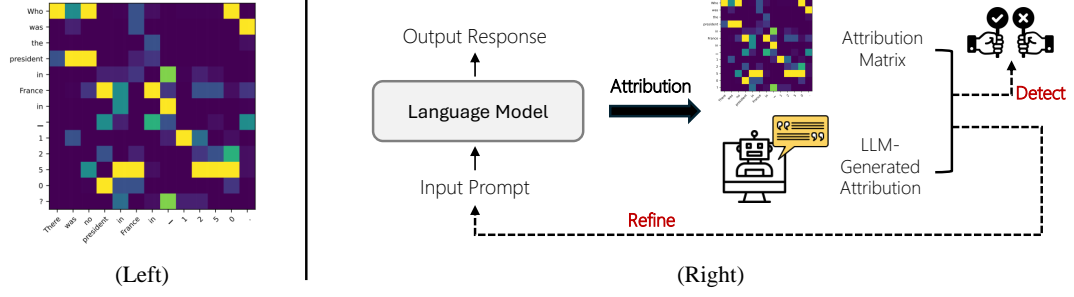
Fig. 2. Left: An example of attribution saliency maps. Right: The pipeline of making use of attribution results to improve LLMs.

## 2.2 Usability of Attribution Methods for LLMs

Explanations provide a lens to evaluate whether a model's predictions are grounded in reasonable rationale and to identify opportunities for improving prompt instructions. Based on this, we discuss two scenarios where LLMs can benefit from attributions: model diagnosis and model enhancement.

- **Model Diagnosis.** By comparing the LLM's output with attribution results, it is possible to detect ungrounded or unsupported information. Formally, let $x$ denote the input, $\hat{y} = f(x)$ be the target model's output, and $a$ be the attribution. The goal is to construct a detector $g_{\text{detect}}(a, x, \hat{y}) \in \{-1, +1\}$ where $+1$ indicates that the prediction can be trusted, while $-1$ means the opposite. Here $g_{\text{detect}}$ can be implemented as an LLM. For instance, Gao et al. [79] proposes a framework called RARR, which leverages attribution to identify unsupported or misleading content generated by LLMs. In particular, the attribution step is automated by using a query generator to produce comprehensive questions about different aspects of the text that need verification. In addition, Agrawal et al. [4] proposes "consistency checks", which asks the LLM to perform attribution multiple times for the same output. Inconsistencies in the attribution results serve as indicators of potential issues with the factuality of the LLM's output.

- **Model Enhancement.** Attribution results can serve as additional information to refine prompts, thereby improving LLM performance. For instance, Sun et al. [253] proposes "recitation" as a new two-step approach where an LLM $f$ first explicitly recalls relevant factual knowledge $a$ from its own memory before generating an answer for input $x$. The LLM inference is augmented as $\hat{y} = f(x, a)$ instead of the standard $f(x)$. This recitation step allows the model to better access and verify its stored knowledge, similar to how a student might first recite relevant facts they've learned before answering an exam question, leading to more accurate responses. In addition, Krishna et al. [146] proposes a framework called AMPLIFY, which uses attribution results to improve LLM performance by crafting effective prompts. It selects samples in the validation set that were misclassified by the LLM, and uses another model (e.g., GPT-2 or BERT) to compute their attribution scores with respect to the ground truth labels. These samples and their corresponding explanations are used to construct prompts for future test samples.

## 2.3 Challenges

We discuss two major challenges regarding attribution methods for explaining LLMs, including **efficiency** and **faithfulness**. First, according to our analysis in Table 1, many traditional attribution methods are computationally expensive when adapted to LLMs. Furthermore, these methods were primarily developed for tasks in the context of *classification or regression*, while the more urgent need lies in designing attribution methods tailored for LLM *generation*. The diverse and dynamic nature of LLM generation demands the development of novel explanation paradigms that go beyond existing frameworks. Second, while some approaches utilize the generative capabilities of LLMs to produce free-text attributions,

Fig. 3. Taxonomy of explanation methods for LLM internals, including explaining *model components* and *hidden representations*.

the faithfulness of these attributions remains underexplored. Ensuring that the generated explanations truly reflect the model's reasoning process is a challenge to be tackled for building trust and reliability in LLM explanations.

## 3 LLM DIAGNOSIS AND ENHANCEMENT VIA INTERPRETING MODEL INTERNALS

This section discusses XAI methods that interpret LLM internals. We delve into the insights these methodologies offer, which can be instrumental in refining and enhancing the design of language models. Specifically, we review research that focuses on interpreting **model components** and **latent representations**, respectively, as shown in Figure 3.

### 3.1 Explanation of Model Components

LLMs adopt transformers as their basic architectures, which typically comprise two types of major components: Self-attention Mechanism (ATT) and Feed-forward Network (FFN). Researchers have focused on analyzing individual components as well as their interactions to gain insights into how LLMs make their predictions. We start by reviewing the connection between these components, focusing on how information is transformed through the residual stream.

We then discuss methods to analyze the contribution of single or multiple components, respectively. We conclude by summarizing how these methods can be applied to improve model performance.

*3.1.1* **Linear Relation between Components**. We consider a decoder-only LLM $f$ with $L$ internal layers, each of which is either an ATT or FNN layer, denoted as $f^l$. There is an input embedding matrix $\mathbf{W}_e$ that encodes each word of input text $x$ into a $D$-dimensional vector, i.e., $\mathbf{x}^0 = \mathbf{W}_e[x] \in \mathbb{R}^{|x| \times D}$. The internal representation at each layer is computed by $\mathbf{x}^l = f^l(\mathbf{x}^{l-1}) + \mathbf{x}^{l-1}$. Finally, there is an output embedding matrix $\mathbf{W}_u$ that decodes the hidden representations for next word prediction, i.e., $f(x) = \mathbf{x}^L \cdot \mathbf{W}_u$. By decomposing the final prediction [71], we have

$$f(x) = \underbrace{\mathbf{x}^0 \mathbf{W}_u}_{\text{Residual Stream}} + \sum_{l=1}^{L} \underbrace{f^l\left(\mathbf{x}^{l-1}\right) \mathbf{W}_u}_{\text{Layer Update}}. \tag{2}$$

Equation (2) demonstrates two critical insights in understanding LLM mechanism: **(1) Model prediction is a sum of layer outputs [200].** That is to say, the final prediction is an additive combination of the initial input embedding and the contribution from each layer. This linear transformation allows us to directly trace how each component shapes the overall prediction. Please note that the residual mechanism does not necessarily lead to a linear structure of model hidden spaces. For example, ResNet [109] does not show such a linear structure because there are non-linear operations applied to its residual stream. **(2) The residual stream serves as a communication channel across layers [71].** In particular, each layer first extracts information from the residual stream, and then injects its processed output back into it. This shared residual stream enables multiple layers to work collaboratively, indirectly contributing to specific tasks by continuously updating the shared residual stream.

*3.1.2* **Importance Analysis on Single Component**. Inspired by the first insight of Equation (2), early works [75, 87, 274] quantify the **contextual impact** of a single layer $f^l$ in predicting a certain word $w$ under the context $x$ by computing $f^l(\mathbf{x}^{l-1}) \mathbf{W}_u[:, w]$, where a greater value indicates a stronger impact. Following this trend, some researchers [71, 196] propose to more precisely quantify the contextual impact by conducting interventions on layer outputs, i.e., $\text{diff}(\mathbf{h} \mathbf{W}_u[:, w], \widetilde{\mathbf{h}} \mathbf{W}_u[:, w])$, where $\text{diff}(\cdot, \cdot)$ measures the difference of two inputs, $\mathbf{h} = f^l(\mathbf{x}^{l-1})$, and $\widetilde{\mathbf{h}}$ is a intervention of $\mathbf{h}$. In practice, the choices of $\widetilde{\mathbf{h}}$ can vary, such as a constant vector [203, 275], a noise vector [196], and a hidden representation of other context [104]. On the other hand, researchers [55, 88, 298] extend this framework to analyze the **global impact** of a single layer $f^l$ in predicting word $w$ by computing $\mathbf{W}_o^l \cdot \mathbf{W}_u[:, w]$, where $\mathbf{W}_o^l$ is the output weight matrix of layer $f^l$ that maps the processed results back to the residual stream.

*3.1.3* **Circuit Analysis on Multiple Components**. Motivated by many emerging abilities of LLMs, it is interesting to explore how **different layers work collaboratively** to perform such complex tasks, called circuit analysis [35]. Specifically, the goal of circuit analysis is to identify a subset of layers, whose combination performs a particular function, such as "indirect objects identification" [275] and "number comparison" [104]. Recall that the residual stream enables the communication across layers, researchers [84] formalize the forward pass of an LLM as a directed acyclic graph, where each node is one particular layer, and an edge is a representation flow from one layer to another via the residual stream. Under this setting, a circuit is a sub-graph with a distinct function. Identifying a specific circuit consists of three steps [48]: (1) selecting an interested task and constructing a dataset, (2) building the model's graph, and (3) measuring the importance of sub-graphs with interventions (Sec. 3.1.2). This circuit identification process is time-consuming as

it needs a large number of forward passes for interventions. To overcome this challenge, researchers [105, 193, 236] propose a gradient-based method to approximate the intervention results with one more backward pass for each sample.

*3.1.4* **Usability of Component Analysis**. The success of component analysis enables researchers to isolate the entire LLM into different parts for specific tasks. Therefore, the direct application of this technique is either modifying (enhancing or weakening) or removing certain parts that are essential or redundant to downstream tasks, resulting in two common applications: Knowledge Editing and Model Pruning.

- **Knowledge Editing.** Component analysis can be used to locate the exact position storing a certain piece of knowledge [52, 281], and then update it accordingly to control the minimum effect on other knowledge. Specifically, we aim to update an outdated knowledge $t = (s, r, o)$ to a new one $t' = (s, r, o')$. Given the prefix $[s, r]$ of $t$, ROME [196] first identifies a specific layer $f^{l*} \in f$ whose activation maximally contributes to predicting the original object $o$. This is done by maximizing the difference $f(o|[s, r]) - f(o|[s, r], do(\mathbf{h}^l = \mathbf{h}^l + \epsilon))$, where $\mathbf{h}^l$ denotes the output activation for input $[s, r]$ at layer $f^l$, and $\epsilon$ denotes a random noise to simulate a perturbation. ROME then modifies only the parameters of this identified layer $f^{l*}$ to update the prediction to the new object $o'$, without affecting unrelated knowledge stored elsewhere. The key limitation of ROME is that it can update only one piece of knowledge at a time, leading to an unacceptable time complexity to update a batch of knowledge. To overcome this challenge, follow-up works [97, 107, 167, 197] scale up this process by skipping the knowledge locating step in ROME, and instead, they focus on learning the representation of updated knowledge at each layer. This line of work has been further applied to edit user personal information for privacy [297], and edit out harmful concepts for safety [276].

- **Model Pruning.** Component analysis allows researchers to identify the specific roles of each model component for a given task, and thus, redundant or irrelevant components can be pruned to speed up the inference. Dalvi et al. [53] directly measure the functionality of $f^l \in f$ by monitoring the expected difference $\mathbb{E}_{x,y \in \mathcal{D}}[p(y|x; f) - p(y|x, do(f^l(l) = \mathbf{0}); f)]$, where $\mathcal{D}$ is a collected dataset for the downstream task. Recent studies [91, 195] directly measure the redundancy by quantifying the similarities of hidden representations from different layers. Empirical studies [53, 72] found that even over half of the model parameters are redundant for a certain task, and dropping them off can significantly reduce the inference time without sacrificing performance.

## 3.2 Explanation of Latent Representations

Many researchers have studied latent representations of LLMs (i.e., $\mathbf{x}^l$) to interpret model behaviors. It is critical to recognize that the representations are not naturally interpretable because of their ***polysemantic*** nature [12, 30, 237], indicating that each dimension of the latent space represents multiple semantic meanings, called **concepts**. Thus, we **cannot** clearly understand a latent representation by reviewing the values in each dimension. There are various lines of work to tackle this challenge, and we categorize them according to whether they require additional training.

*3.2.1* **Non-Training based Methods**.

Both the *input and output layer* of transformer-based LLM $f$ are directly interpretable as they operate on human-readable text. Hidden representations across different layers share the same residual stream (see Section 3.1.1), allowing information to flow seamlessly between layers. Consequently, non-training based interpretation methods leverage this insight by transferring the hidden representation $\mathbf{x}^l$ to either the first or the last layer to facilitate interpretation.

- **Projecting $\mathbf{x}^l$ to Output Embeddings.** Given a hidden representation $\mathbf{x}^l$ at layer $l$, *Logit Lens* [209] understands its semantics by projecting it to the output word embedding matrix $\mathbf{W}_u$. Specifically, we interpret $\mathbf{x}^l$ by collecting the $K$

words whose output embeddings could maximally activate $\mathbf{x}^l$, i.e.,

$$\mathcal{I} = \underset{\mathcal{V}' \subset \mathcal{V}, |\mathcal{V}'|=K}{\arg\max} \sum_{w \in \mathcal{V}'} \mathbf{x}^l \cdot \mathbf{W}_u[:, w], \tag{3}$$

where $\mathcal{V}$ is a pre-defined vocabulary set. The Logit Lens can be viewed as an instance of skipping all later layers after $\mathbf{x}^l$ [132]. However, researchers [23, 60] find that there is a distribution shift between the intermediate layer $l$ and the last layer $L$, and thus, they propose to train a *translator* to learn the distribution shift to improve its interpretability for representations from shallow layers.

- **Feeding $\mathbf{x}^l$ as Input Embeddings.** In contrast with moving hidden representation $\mathbf{x}^l$ to the output layers, *Self-IE* [40] proposes feeding $\mathbf{x}^l$ as an input word embedding in another round of model inference, where we prompt LLM to explain the meaning of the fed representation. For example, we may instruct LLM to interpret $\mathbf{x}^l$ with the prompt "Please summarize [X] in one sentence.", where "[X]" is a placeholder, whose embedding will be replaced with our interested hidden representation $\mathbf{x}^l$. Since LLMs are trained to follow human instructions, they are expected to generate a sentence to describe the information encoded in $\mathbf{x}^l$.

### 3.2.2 Supervised Training based Methods.

The supervised training-based interpretation methods break the polysemantic nature of LLM latent space by leveraging some *annotated* datasets. They focus on exploring two main research questions: (1) whether the LLM latent space encodes a specific concept; and if so, (2) whether this concept impacts LLM's predictions and how significant.

- **Identifying Existence of Certain Concepts.** *Probing* [96] is the most traditional technique in analyzing the knowledge encoded within hidden representations of LLMs. In specific, the Probing technique requires an annotated dataset to train a probe $p$ with parameter $\theta$ to predict the occurrence of interested knowledge, i.e., $p : \mathbb{X} \rightarrow [0, 1]$, where $\mathbb{X}$ is the latent space of LLMs at a certain layer. Here, the dataset consists of positive samples (with label "1") that clearly demonstrate our interest concept and negative samples that clearly exclude such concept. For example, early works [22] use this technique to identify some morphology concepts (e.g., part-of-speech tags and morph tags) in the LLM's latent space by constructing a dataset where each example has human annotations on these tags. If the probe $p$ achieves a significantly high accuracy on a certain dataset, then we could conclude that the latent representation encodes that concept. This technique has been used to study diverse concepts, ranging from low-level linguistic patterns [22, 38, 117, 117, 174] to high-level semantic meanings [32, 190, 344, 352]. The improvements in the probing technique [116, 218, 267] mainly focus on designing baselines to confirm that the high accuracy achieved by probe $p$ is led by the knowledge provided by the latent space $\mathbf{x}$.

- **Measuring Importance of Certain Concepts.** To quantify how much a specific concept impacts the model's predictions, researchers often use *TCAV* [138]. We first train a linear probe $p$ with parameter $\theta \in \mathbb{R}^D$ to distinguish latent representations containing concept $c$ or not, where $\theta$ is named as the concept activation vector (CAV). TCAV measures how sensitive the model's output is to perturbations along this concept direction (e.g., via directional derivatives of the model's logits). Intuitively, if the model's predictions change substantially when the representation is shifted along the learned concept direction, we can conclude that the concept has a high *importance* for the model's decision. Mathematically, the perturbation is realized by computing the gradient of $f(x)$ with respect to $\mathbf{x}^l$ and projecting it onto $\theta$ as $\Delta = \nabla_{\mathbf{x}^l} f(x) \cdot \theta$, so that a larger $|\Delta|$ indicates that small shifts in the latent space along $\theta$ lead to substantial changes in $f(x)$, underscoring the concept's influence.

### 3.2.3 Unsupervised Training based Methods.

Breaking the polysemantic constraint for interpretation can be formalized as an unsupervised learning problem on the basis vectors of the LLM latent space [73], where each basis vector refers to a clear and concise **semantic concept**, indicating a *monosemantic* nature. Mathematically, given a total of $N$ hidden representations $\mathbf{X} \in \mathbb{R}^{N \times D}$, we aim to learn two matrices $\mathbf{A} \in \mathbb{R}^{N \times C}$ and $\mathbf{C} \in \mathbb{R}^{C \times D}$ by minimizing the loss

$$\underset{\mathbf{A},\mathbf{C}}{\arg\min} \ \|\mathbf{X} - \mathbf{AC}\|^2, \tag{4}$$

where $\mathbf{C}$ serve as the learned concepts, and $\mathbf{A}$ is the coefficient of the instances on each concept. There are some additional conditions to ensure that the learned concept vectors satisfy our expectations on monosemantic.

- **Orthogonal Assumption on Concepts.** Early works [201, 298] applies the *Singular Vector Decomposition* to learn a set of concept vectors $\mathbf{C}$ that are perfectly orthogonal, i.e.,

$$\mathbf{CC}^\top = \mathbf{I}, \tag{5}$$

  where $\mathbf{I}$ is an identical matrix. However, satisfying this hard request is computing costly in practice, leading to limited scalability in interpreting hidden representations from giant LLMs with both large $N$ and $D$.

- **Sparse Assumption on Concepts.** *Sparse Autoencoder* [210] (SAE) is another alternative that relaxes the orthogonal condition to be a smoother one in which the sample coefficients matrix is asked to be sparse, i.e.,

$$\mathbf{A} = \text{Top-K}(\mathbf{X} \cdot \mathbf{C}^\top), \tag{6}$$

  where Top-K activation enforces all values to be zero unless the $K$ largest ones from the inputs. Specifically, researchers [30, 50] train a SAE $g(\mathbf{x}) = \text{Top-K}(\mathbf{x} \cdot \mathbf{W}^\top) \cdot \mathbf{W}$ to reconstruct hidden representations $\mathbf{x}$ by minimizing $\|\mathbf{x} - g(\mathbf{x})\|^2$, where $\mathbf{W} \in \mathbf{R}^{C \times D}$ and Top-K activation enforces all values to be zero unless the $K$ largest ones from the input. This approach has shown effectiveness in decomposing the latent representations from giant LLMs with a hundred billion parameters, such as GPT-4 [80] and Claude 3 [258].

Once these concept vectors are obtained, their representing meanings can be further described with natural language by selecting $M$ input texts that could maximally activate the related concept vectors [30], i.e.,

$$\mathcal{I}_c = \arg \max_{\mathcal{X}' \subset \mathcal{X}, |\mathcal{X}'|=M} \sum_{x \in \mathcal{X}'} f^l(x)\mathbf{C}[c]. \tag{7}$$

However, the latest research [98, 300] has shown that such an input-based explanation cannot well interpret their impacts on controlling LLM responses, and thus, they propose to interpret these learned concept vectors by building the connection to the output texts. We refer audiences to check [246] for a more detailed review of designing and interpreting SAEs for LLMs.

### 3.2.4 *Usability of Interpreting Latent Representations*.

Latent space interpretations improve LLMs further at each stage of their entire life-cycles. In particular, we consider three stages of model development, including the "pre-processing" stage for data preparation, the "training" stage that trains LLM on prepared data, and the "inference" stage that uses a trained model to operate user requests.

- **Data Pre-processing.** Existing studies have empirically found that the same model but trained on different datasets with different diversity [83], complexity [305], and quality [162] levels may lead different downstream task performance. The latent space interpretation techniques enable researchers to understand the geometric structure of a dataset and, further, to select a subset that preserves a certain geometric structure. For example, to obtain a subset of the dataset preserving as diverse topics as possible, SAE-GreedSelect [310] greedily select examples to maximize the

total number of activated SAE features. Results show that a selected subset with only 5% of total samples can train a model with comparable performance to a full-trained one, significantly reducing the training cost.

- **Training-Time Regularization.** Either the trained parameter $\theta$ of probe $p$ or each column vector $\mathbf{C}[c]$ learned by SAE or SVD indicates a feature in the LLM latent space. Therefore, these learned feature vectors can be used to perform feature engineering in the LLM latent space for model training. For example, Yin et al. [318] propose a feature-level constrained fine-tuning approach. This method effectively prevents the fine-tuned model from changing significantly from its pretrained counterpart without explicitly involving the pretrained model during fine-tuning, thereby substantially reducing computational costs. In contrast, Wu et al. [299] introduce a more semantically meaningful regularization approach for model training. They first determine whether each learned feature vector $\mathbf{C}[c]$ is causally relevant to the downstream task by analyzing its natural language explanation $\mathcal{I}_c$. Subsequently, they impose a penalty during fine-tuning to discourage reliance on features identified as task-irrelevant (i.e., shortcuts). This approach effectively enhances the generalizability of the fine-tuned model.

- **Inference-Time Control.** As the learned feature vectors by probing or SAEs break the polysemantic nature and achieve monosemantic, it is possible to directly steer LLM latent space toward a specific direction without sacrificing other abilities. Formally, let $\mathbf{z}$ denote the learned feature vector to an interested concept, we can *erase* it from the hidden representation of input $x$ by $\mathbf{x}' = \mathbf{x} - \mathrm{ReLU}(\mathbf{x} \cdot \mathbf{z}^\top) \cdot \mathbf{z}$, or constantly *amplify* it to a certain level $\alpha$ by $\mathbf{x}' = \mathbf{x} + \alpha \cdot \mathbf{z}$. The combination of these two techniques on different categories of feature vectors leverages model steering for different purposes. For example, Wu et al. [300] observe that certain feature vectors semantically correlate with specific safety-related concepts, such as "First Aid," based on summarizing their natural language explanations $\mathcal{I}_c$. By enforcing a constant activation of these safety-related feature vectors during model generation, they demonstrate that the model reliably switches from generating harmful responses to explicitly rejecting inappropriate or malicious user inputs. Generally, this technique has shown success in steering LLMs during inference to improve their helpfulness [112, 126, 230, 252], safety [10, 25, 230, 352], faithfulness [5, 39], fairness [249, 352], and currentness [240].

## 3.3 Challenges

*3.3.1 Explanation Efficiency.* The most significant challenge in analyzing the inner workings of LLMs is their demand for a large amount of computing resources. Circuit analysis (Section 3.1.3) is one of the examples that is bounded by this bottleneck. Specifically, performing circuit analysis on one sample still requires a large amount of forward inference of LLMs under different perturbed internals of that sample, leading to limited scalability to analyze over a large dataset. Therefore, existing works on **circuit analysis focus on studying particular behaviors of LLMs**, such as "Name Mover Head" and "Duplicate Token Head" for object identifications [275], "Single Letter Head" and "Correct Letter Head" for multiple-choice question answering [170], and "Capitalize Head" as well as "Antonym Head" for writing [261]. On the other hand, although SAEs (Section 3.2.3) have shown strong promise in understanding LLM latent representations, training such proxy models requires collecting LLM representations over a large corpus. For example, researchers [80] found that **training an SAE with only 1 million features requires over 10 billion training tokens**. Future works may explore computing-efficient methods to understand LLMs' diverse behaviors comprehensively.

*3.3.2 Explanation Faithfulness.* Although many proposed methods have shown strong promise with certain usability cases (Section 3.1.4 and Section 3.2.4), some researchers still hold concerns about whether they faithfully interpret the inner workings of LLMs since they may not practically outperform some trivial baselines. For example, although many

researchers [30, 300] can apply SAEs (Section 3.2.3) to steer LLM responses for specific purposes, some research [29, 69, 301] has observed that **SAE-based methods cannot significantly outperform prompting baselines**, posing a serious concern about the effectiveness of SAEs. Furthermore, other researchers [113, 215] suspect that SAEs cannot faithfully explain LLMs' internals, as they observe that **SAEs can even learn features from randomized LLMs**. These concerns suggest two critical questions for future works: (1) How do we ensure our learned features faithfully represent LLM internals? (2) In which cases are SAE-based methods more effective than prompting strategies?

## 4 LLM DEBUGGING WITH SAMPLE-BASED EXPLANATION

In this section, we discuss sample-based explanation strategies for LLMs, which aim to trace back the answers generated by LLMs to specific training samples (i.e., documents) or document fragments in the corpora. The utility of sample-based explanations for LLMs is multifaceted. First, tracing back the predictions of LLM to the training samples can provide evidence for the generation results, which facilitates model debugging in cases of errors and increases the trustworthiness of the model from users when the outcomes are accurate. Second, it can also help researchers understand how LLMs generalize from training samples. If the outputs of LLMs can be traced back to exact subsequences directly spliced from the training data, it might suggest that the LLM is simply memorizing the data.

### 4.1 Literature Review of Sample-based Explanation

*4.1.1 Influence Function-based Methods.* One strategy to quantify the influence of a training sample $z_i$ in the dataset $\mathcal{D}_{train}$ to a test sample $z$ is through the influence function [103, 142], which measures the change of the prediction loss $\mathcal{L}(z, \theta)$ for $z$, when the training sample $z_i$ is removed from the dataset:

$$\mathcal{I}(z_i, z) = -\nabla_\theta \mathcal{L}(z, \hat{\theta})^\top \mathbf{H}_{\hat{\theta}}^{-1} \nabla_\theta \mathcal{L}(z_i, \hat{\theta}), \tag{8}$$

where $\nabla_\theta \mathcal{L}(z, \hat{\theta})$ is the gradient of loss $\mathcal{L}$ on $z$, and $\mathbf{H}_{\hat{\theta}} \stackrel{\text{def}}{=} \frac{1}{N} \sum_{i=1}^{N} \nabla_\theta^2 \mathcal{L}(z_i, \hat{\theta})$ is the Hessian matrix. To improve efficiency, Koh and Liang [142] adopt an iterative approximation process to calculate the Hessian-Vector Product (HVP) in Eq. (8), where the memory complexity can be reduced to $O(P)$ and time complexity to $O(NPr)$ ($r$ is the number of iterations). To further reduce the complexity, [219] propose TracIn, which measures the influence by calculating the total reduction of loss whenever $z_i$ is included in the minibatch during model training as follows:

$$\mathcal{I}_{\text{TracIn}}(z_i, z) = \sum_{t:z_i \in \mathcal{B}_t} \mathcal{L}(z, \theta_t) - \mathcal{L}(z, \theta_{t+1}) \approx \frac{1}{|B_t|} \sum_{t:z_i \in B_t} \eta_t \nabla_\theta \mathcal{L}(z_i, \theta_t) \cdot \nabla_\theta \mathcal{L}(z, \theta_t), \tag{9}$$

where $\mathcal{B}_t$ is the $t$-th mini-batch, $\theta_t$ is the parameter at the $t$-th step, $\eta_t$ is the step size. TracIn only leverages gradients, which substantially improves the efficiency. In addition, Schioppa et al. [238] propose to use Alnordi iteration [11] to find the dominant eigenvalues and eigenvectors of $\mathbf{H}_{\hat{\theta}}$ on randomly sampled subsets $\mathcal{D}_{sub}$, with $|\mathcal{D}_{sub}| \ll |\mathcal{D}_{train}|$, where the diagonalized Hessian can be cheaply cached and inverted. Observing that finding the most influential training sample on $z$ needs to iterate overall $N$ training samples, [94] propose to use fast KNN to pre-filter a small subset of influence-worthy data points from $\mathcal{D}_{train}$ as candidates to explain small pretrained language models, whereas [102] propose to find a small subset $\mathcal{D}_{sub} \subset \mathcal{D}_{train}$ whose gradient is the most similar to the downstream task examples. [92] propose to use the Eigenvalue-corrected Kronecker-Factored Approximate Curvature (EK-FAC) approximation to scale influence functions to 52B LLMs. For adaptation, only influences mediated by the feed-forward networks (FFNs) are considered. Based on the assumption that weights from different FFN layers are independent, the EK-FAC approximated

influence is finally formulated as the sum of influences mediated by each layer as follows:

$$\mathcal{I}_{\text{EKFAC}}(z_i, z) = \sum_l \nabla_{\theta^{(l)}} \mathcal{L}(z, \hat{\theta})^\top (\hat{\mathbf{G}}_{\hat{\theta}^{(l)}} + \lambda^{(l)} \mathbf{I})^{-1} \nabla_{\theta^{(l)}} \mathcal{L}(z_i, \hat{\theta}), \tag{10}$$

where $\theta^{(l)}$ denotes the weights of the $l$-th MLP layer, and $\hat{\mathbf{G}}_{\hat{\theta}^{(l)}}$ is the EK-FAC approximated Gauss-Newton Hessian for $\theta^{(l)}$. Since the inversion of $L$ small $K_l \times K_l$ matrices (i.e., $O(L \times K_l^3)$) is substantially more efficient than the inversion of a large $LK_l \times LK_l$ matrix (i.e., $O((LK_l)^3)$), $\mathcal{I}_{\text{EKFAC}}$ can be adaptable to very large models.

Recently, influence functions have also been used to explain and improve the in-context learning ability of LLMs. Based on the finding [268] that during in-context learning, LLMs can be viewed as implicitly "learning" an internal kernelized least square surrogate objective on the few-shot examples included in the prompt, Zhou et al. [348] propose to use this hypothesized objective as the loss function to calculate the influence of each in-context sample on the LLM generation, and based on which select and re-rank the samples. They show that the in-context learning ability of LLMs can be improved after the sample selection and demonstration order re-ranking.

*4.1.2 Embedding-based Methods.* Another strategy for sample-based explanation leverages hidden representations within the transformer architecture, which is recognized for encoding high-level semantics from textual data, to calculate the semantic similarity between $z$ and $z_i$. The similarity can also be used to measure the influence of $z_i$ on $z$ as explanations [226]. Specifically, [6] propose to represent the training sample $z_i$ and test sample $z$ by concatenating the input and output as $z_i^{cat} = [x_i || y_i]$, $z^{cat} = [x || y]$. The concatenation is feasible for generation tasks where the output $y$ lies in the same token sequence space as the input prompt $x$. The similarity between $z_i$ and $z$ can then be calculated as:

$$\mathcal{I}_{\text{emb}}(z_i, z) = \frac{f_{\hat{\theta}}^{(l)}(z_i^{cat})^\top \cdot f_{\hat{\theta}}^{(l)}(z^{cat})}{\left\| f_{\hat{\theta}}^{(l)}(z_i^{cat})^\top \right\| \left\| f_{\hat{\theta}}^{(l)}(z^{cat}) \right\|}, \tag{11}$$

where $f_{\hat{\theta}}^{(l)}$ is the sub-network that outputs the $l$-th layer intermediate activation of $f_{\hat{\theta}}$. The Eq. (11) has a similar form as the vanilla influence function defined in Eq. (8) as well as its TracIn alternative defined in Eq. (9), which assigns a score $\mathcal{I}$ for the explainee $z$ for each training sample $z_i$ in the dataset $\mathcal{D}_{train}$ as the explanation confidence of the sample $z_i$.

## 4.2   Case Study: EK-FAC-based Influence Estimation

This case study implements the EK-FAC-based influence function [92] and verifies its scalability and effectiveness on LLMs with billions of parameters, including GPT2-1.5B [224], LLaMA2-7B [263], Mistral-7B [134], and LLaMA2-13B.

*4.2.1 Experimental Design.* We use the SciFact dataset [270] as the corpora, which contains the abstract of 5,183 papers from basic science and medicine. The LLMs are obtained by finetuning the pretrained LLMs for 20,000 iterations, where AdamW is used as the optimizer [180]. Then, we use 500 samples from the corpora to estimate the *(i)* uncentered covariance matrices of the activations and pre-activation pseudo-gradients $\mathbf{Q}_A^{(l)}$, $\mathbf{Q}_S^{(l)}$, and *(ii)* the variances of the projected pseudo-gradient $\mathbf{\Lambda}^{(l)}$ for each selected dense layer $l$, and cache them on the hard disk (details see Eqs. (16) and (20) in [92]). We select the c_fc layer for GPT2-1.5B, and gate_proj layer for LLaMA2-7B, Mistral-7B, and LLaMA2-13B. For evaluation, we randomly select 200 samples to construct the test set, which we name SciFact-Inf. Specifically, for the $j$-th sample $z_j = x_j$, we use the first three sentences in $x_j$, i.e., $\hat{x}_j$, to generate a completion $\hat{y}_j$ with the finetuned LLM. Ideally, the $j$-th training sample $z_j$ itself should be the most influential sample w.r.t. the generation of $\hat{y}_j$.

Table 2. Effectiveness of EK-FAC approximated influence function on the SciFact-Inf dataset. Time (Pre.) stands for the time for precomputing the $Q_A$, $Q_S$, and $\Lambda$. Time (Inf.) stands for time for calculating the influence of 100 training samples per test sample.

| Strategy | LLM | Recall@5 | Recall@10 | Time (Pre.) | Time (Inf.) |
|---|---|---|---|---|---|
| | GPT2-1.5B | 0.6368 | 0.7363 | 0h 27min | 0min 28sec |
| | Mistral-7B | 0.6418 | 0.6866 | 2h 05min | 1min 47sec |
| Inf. Func. | LLaMA2-7B | 0.8063 | 0.8308 | 1h 37min | 1min 34sec |
| | LLaMA2-13B | 0.7811 | 0.8940 | 3h 11min | 3min 08sec |

*4.2.2 Results and Analysis.* The results are summarized in Table 2. We observe that the EK-FAC approximated influence function achieves good accuracy in finding the training sample with greatest influence on the generation of LLMs. In addition, we find that the main computational bottleneck in calculating the EK-FAC-based influence is to estimate the covariances $Q_A^{(l)}$, $Q_S^{(l)}$ and variance $\Lambda^{(l)}$, which takes hours when 500 training samples are used for the estimation. However, after that, it is relatively cheap to calculate the influence training samples for each test sample.

## 4.3 Challenges

Overall, explaining the generation of LLMs by tracing back to the training samples is still an emerging area. Open questions need to be addressed to further advance the field.

*4.3.1 Strong Assumptions for Scalability.* The unprecedented number of parameters in modern LLMs causes severe scalability issues for sample-based explanation strategies. This is especially evident for the gradient-based methods, as the HVP in Eq. (8) induces both high computational and space complexity. To address the bottleneck, strong assumptions are usually required to make it feasible for large models. For example, TracIn [219] simplifies the second-order term in Eq. (8) via first-order approximation. [238] assume the Hessian to be low rank. [92] that assume that the weights from different layers of the LLMs are independent, as well as the tokens in different steps, such that EK-FAC can be appropriately applied to approximate the influence function. From the above analysis, we can find that while the method from [92] has the best scalability, it also has the strongest assumption, which may fail to hold in practice. Therefore, how to improve the scalability with weak assumptions needs to be investigated in the future.

*4.3.2 LLM-Oriented Sample-based Explanations.* Finally, we observed that both gradient-based and embedding-based methods are loosely connected to the LLM, as well as the backbone transformer networks. For example, algorithms like TracIn [219] are designed to scale up influence functions to large models, which are not specific for LLMs. Similarly, the embedding-based method proposed in [6] is applicable to most machine learning models with latent representations. [92] considers the specialty of LLMs by utilizing the knowledge neuron assumption of the backbone transformers [279] to simplify the influence function, where the weights considered are constrained to the MLP layers, which may not fully utilize the property of transformers. Therefore, how to further utilize the property of the LLM and the backbone transformer to design LLM-tailored sample-based influence/similarity (either to reduce the computational/space overhead or to improve the explanation quality) is highly promising for future work.

## 5 EXPLAINABILITY FOR TRUSTWORTHY LLMS AND HUMAN ALIGNMENT

In previous sections, we explore the use of explanation techniques for assessing and improving the performance of LLMs. In this section, we shift the focus towards examining LLM trustworthiness. As LLMs are increasingly integrated into high-stakes areas like healthcare, finance, and legal advice, it is crucial that their responses not only are accurate but also
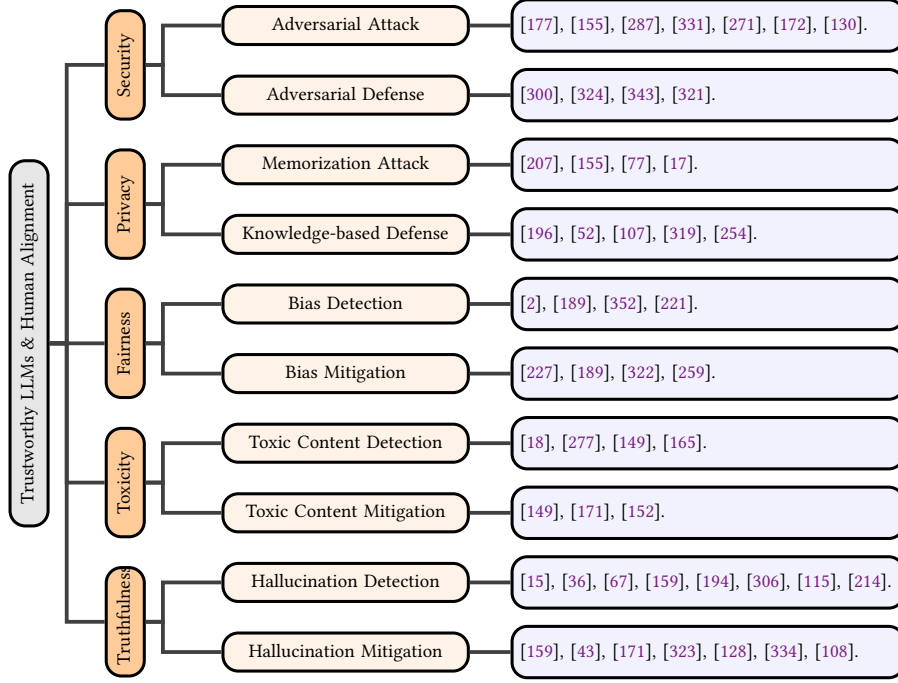
Fig. 4. Taxonomy of explanation techniques for trustworthiness aspects of LLMs.

*align* with human ethical standards and safety protocols [122, 168, 178]. Herein, we examine how explanation techniques, discussed in the previous sections, can be instrumental in assessing LLMs across key aspects of trustworthiness like security, privacy, fairness, toxicity, and honesty. It is worth noting that while explainability itself is an aspect of trustworthiness, it holds the promise of serving as a foundational tool for addressing other trustworthiness concerns.

### 5.1 Security

LLMs are known to be vulnerable to attacks and exploitation, such as spreading misinformation, and poisoning training data [57]. For enhanced safety, LLMs are designed to reject prompts that may result in harmful contents. However, jailbreak techniques can circumvent these restriction measures and manipulate LLMs into producing malicious contents. Malevolent users can craft special prompts that compel LLMs to prioritize instruction following over rejections [155, 177]. For example, through Prefix Injection, attackers can use out-of-distribution prompt prefixes that are less likely to be rejected [271, 287]. Another approach, called Refuse Suppression, involves directing or persuading models to ignore established safety protocols [287, 331], where the instruction following ability is then employed to perform the attack.

By engineering latent representations of LLMs, explanation methods provide a viable way to discover the potential vulnerabilities of LLMs by simulating advanced attacks [172]. For example, a recent work extracts "safety patterns" via explaining the latent space. These patterns are captured from the activation differences between malicious and benign queries. They reflect the internal protection mechanisms within LLMs. Circumventing these patterns leads to novel attacks, which helps exploring potential vulnerabilities of LLMs [165]. Besides, a deeper understanding of fine-tuning can shed light on the reliability of existing safety measures. In particular, Jain et al. [130] apply network pruning, attention maps, and probing classifiers to track the changes of model capabilities from pre-training to fine-tuning.

These tools are helpful in demonstrating that the capabilities gained during fine-tuning can be easily removed through fine-tuning on other unrelated tasks. This finding casts doubt on the robustness of current safety alignments in LLMs. Finally, by probing LLM representations, Zheng et al. [343] investigates the impact of safety prompts that motivates a novel safety prompt optimization method to safeguard LLMs.

## 5.2 Privacy

Recent studies have revealed that LLMs such as ChatGPT can leak extensive amounts of training data through divergence attacks [207]. These attacks utilize specially crafted prompts to lead the model away from its standard chatbot-style generation. The risk of private data exposure through such means poses a serious challenge to the development of ethically responsible models. This issue is compounded by strategies similar to jailbreak attacks, where misalignment is exploited to induce LLMs into operating in an unconventional "developer mode" via out-of-distribution prompts [155].

Enhancing LLM privacy involves two strategic approaches: (1) preventing LLMs from memorizing sensitive data, and (2) establishing safeguards against the release of sensitive information. The second strategy involves employing techniques from jailbreak defenses, treating prompts that solicit private information as potentially malicious. The first approach requires identifying specific knowledge within models, which is traditionally achieved with prompt engineering. However, prompt engineering faces limitations due to LLMs' sensitivity to the phrasing of prompts. Explanatory techniques can serve as an auxiliary tool to confirm whether LLMs have internalized certain knowledge. For instance, factual knowledge is localized via explaining the relation between factual knowledge and neuron activations [52, 107, 196]. In addition, [319] recently proposes the concept of "knowledge boundary" and develops a gradient-based method to explore whether LLMs master certain knowledge independent of the input prompt.

## 5.3 Fairness

The widespread applications of LLMs also bring concerns about exacerbating bias issues in society [78]. For example, in a gender stereotype case, "[He] is a doctor" is much more likely than "[She] is a doctor". In this subsection, we focus on fairness issues related to race, gender, and age [169]. Prompt engineering is a major method to quantify fairness issues within LLMs [2]. Interpretation complements this method by unraveling the mechanisms through which biases are embedded into LLMs. One direction is examining biased attention heads. Ma et al. [189] unveil that approximately 15% to 30% of attention heads linking to specific stereotypes through probing attention heads. Furthermore, recent work has placed LLM representations under scrutiny [352]. Typically, representations of crafted templates showing specific biases are examined to derive a vector that identifies a certain bias.

To achieve fair model outputs, a diverse range of mitigation techniques have been proposed. One stream of work proposes to debias LLMs at the embedding level. For example, one work alters biased embeddings with minimal alterations to make them orthogonal to neutral embeddings [227]. Additionally, some studies remove biases at the level of attention heads. Ma et al. [189] address this by pruning biased attention heads. Beyond modifying embeddings and pruning attention heads, another strategy targets a specific group of neurons known to propagate biases. It unlearns the biases by retraining weight vectors for these neurons [322]. Besides, bias mitigation can also be approached from a data-centric perspective using the most biased training samples, then modifying them to fine-tune the model [259].

## 5.4 Toxicity

Toxicity is another form of harmful content that LLMs may produce. This issue arises from training documents containing elements of toxicity that can hardly be fully eliminated. Toxicity can be identified by interpreting LLM components like

the feed-forward layers and attention heads. For instance, recent work reveals how toxicity is represented within LLMs by identifying multiple vectors and relevant dimensions promoting toxicity within the MLP layers [149]. Furthermore, the exploration of geometric structures in per-layer representations offers another way to detect toxicity. [18] extract input features from MLPs to describe the domain of prompts and classifying toxic remarks.

Aforementioned insights also shed light on mitigation strategies. Motivated by the finding that toxicity can be reduced by manipulating relevant vectors, Lee et al. [149] utilizes paired toxic and non-toxic samples to fine-tune models so that non-toxic content is promoted. By examining the changes in the parameter matrices during the fine-tuning process, it substantiates that even minor adjustments to these critical vectors can reduce toxicity. Built on the observation that LLMs' representations are updated by outputs from attention layers [71], another work attempts to reduce toxicity by identifying the "toxicity direction" and then adjusting representations in the opposite direction [152].

## 5.5 Truthfulness

LLMs tend to confidently produce false statements that fall into two main categories: 1) statements that contradict learned knowledge, also known as *dishonesty*; 2) statements that are factually incorrect and fabricated by models, that is *hallucination*. In the following, we look into explainability tools that aim to understand the above two behaviors.

*5.5.1  Dishonesty.* Dishonesty of LLMs describes models' ability to produce false statements compared to their learned information. Studies have been undertaken to understand how and why it happens. One notable work distinguishes dishonesty by training a classifier to predict the accuracy of statements, on top of activations from true/false statements [15]. The classifier reaches an accuracy range between 60% and 80% [15]. Furthermore, Campbell et al. [36] localize dishonesty behaviors at the level of attention heads. Despite adopting above probing-based method, it also employs activation patching to substitute lying activations with honesty ones. Both approaches have witnessed the importance of layer 23−29 in flipping dishonesty behaviors. Besides, another popular method tries to study the geometric structure of true/false statements[194]. A linear structure and the truth directions are derived to mitigate the dishonest behaviors.

*5.5.2  Hallucinations.* Hallucinations in LLMs can arise due to poor data quality and the lack of explicit knowledge [307]. However, whether LLMs are aware of their hallucination behaviors remains unknown. Recent work investigates this question by examining models' hidden representation space [67]. It defines an "awareness" score to quantify the uncertainty of LLMs' answers, finding that adversarially induced hallucination can increase models' awareness. Additionally, Li et al. [159] identifies differences between models' output and their inner activations as a potential source of hallucination. By training two orthogonal probes to mitigate hallucinations, they reveal that "truth" might exist in a subspace instead of a single direction [159]. Another work investigates the source of hallucination by analyzing patterns of source token contributions through perturbations [306]. They find that hallucinations may stem from the model's excessive dependence on a restricted set of source tokens, and the static distribution of source token contributions.

## 5.6 Challenges

We discuss challenges in employing explanations to improve models' trustworthiness and enhance alignment: 1) limitations of existing evaluation techniques, and 2) shortcomings of mitigation strategies based on explanations.

*5.6.1  Challenges of Existing Evaluation Methods.* Current detection methods primarily focus on layers, attention heads, and representations. However, we still lack robust understanding of how a certain trustworthiness concept (e.g., fairness, harmlessness) is comprehended by LLMs. Furthermore, we lack general strategies to identify this knowledge. Existing methods often rely on building an evaluation dataset, but it is difficult to guarantee the its comprehensiveness. Meanwhile,

the definition of what constitutes a trustworthiness concept often varies among different stakeholders and contexts, complicating the evaluation process. Moreover, existing localization approaches rely on probing classifiers or casual scrubbing, which might not be reliable, since the performance of probing classifiers depends on the comprehensiveness of pre-defined biases, while casual scrubbing usually introduces new variables that complicate the analysis.

*5.6.2 Challenges of Mitigation Strategies.* Currently, mitigation methods for LLMs are typically developed based on explanations. Existing explanations are implemented using techniques from mechanistic interpretability and representation engineering [339]. However, neither stream of methods can fully address the issues. For example, the identified directions from representation engineering and patched activations from mechanistic interpretability can only mitigate issues to a certain extent. Moreover, the changes to either representations or activations could also influence other aspects of models' capabilities, which we are yet unable to evaluate.

## 6 EFFECTIVE LLM INFERENCE VIA EXPLAINABLE PROMPTING

A key distinction between LLMs and traditional machine learning models lies in the LLMs' ability to accept flexibly manipulated input data, namely prompts, during model inference [176]. To mitigate the opacity issue in LLM predictions, we can enhance prompts with *explainable content*. In response, LLMs can reveal their decision-making processes during inference, which improves the explainability and even rationality of their behaviors.

### 6.1 Explainable Prompting Paradigms and Effects

To better understand and control how LLMs generate outputs, many researchers directly prompt LLMs to provide the rationales along with their predictions [24, 123, 136, 289, 315]. **Chain-of-Thought** (CoT) prompting [289] stands out by employing explicit reasoning to guide the generation process. Formally, we define the language model as $f_\theta$, and input prompt as $X = \{x_1, y_1, x_2, y_2, ..., x_n\}$, where $x_1, y_1, x_2, y_2, ..., x_{n-1}, y_{n-1}$ denote the example question-response pairs for in-context learning, and $x_n$ is the actual question. In a standard prompting scenario, we have the model output as $y_n = \arg\max_Y p_\theta(Y|x_1, y_1, x_2, y_2, ..., x_n)$. However, this approach does not provide insight into how the answer $y_n$ is generated. Therefore, CoT proposes including human-crafted explanations $e_i$ for the $i$-th in-context example, resulting in a modified input format $X = \{x_1, e_1, y_1, x_2, e_2, y_2, ..., x_n\}$. Given the input, the model will output not only prediction $y_n$ but along with an explanation $e_n$: $e_n, y_n = \arg\max_Y p_\theta(Y|x_1, e_1, y_1, x_2, e_2, y_2, ..., x_n)$. Specifically, the usefulness of CoT methods lies in several key aspects:

- **Reducing Errors in Reasoning:** By breaking down complex problems into a series of smaller tasks, CoT reduces errors in logic-oriented tasks, leading to a more precise resolution of intricate problems [222, 282, 289, 336].
- **Providing Adjustable Intermediate Steps:** CoT enables the outlining of traceable intermediate steps within the problem-solving process. This feature enables users to trace the model's thought process from inception to conclusion, and to adjust the prompts if undesirable model behaviors are observed [188, 279].
- **Facilitating Knowledge Distillation:** The step-by-step reasoning derived from larger LLMs can serve as a specialized fine-tuning dataset for smaller LLMs. It allows smaller models to learn complex problem solving by following explanations, effectively teaching them to tackle intricate questions with enhanced reasoning capabilities [191].

In addition, techniques beyond simple CoT have been developed to broaden the range of reasoning paths available to LLMs [24, 42, 314, 315, 328, 329]. For example, **Tree-of-Thoughts (ToT)** [315] advances beyond the traditional linear chain-of-thought reasoning, offering a more versatile structure that allows models to navigate through multiple reasoning paths. Similarly, **Graph of Thoughts (GoT)** [24] transforms the outputs into a graph, where information

pieces are nodes and their connections are edges, enabling a more intricate and connected form of reasoning paths. By organizing data into nodes (individual concepts or pieces of information) and edges (relationships between these concepts), GoT makes the logical connections within complex systems more understandable [314].

However, the above prompting-based methods may not always trigger LLMs to provide explanations. Thus, many efforts [46, 213, 311, 333] have been made to empower LLMs to acquire the inherent ability of providing their reasoning paths. DeepSeek-R1 [93] is one of the remarkable checkpoints, where they train LLMs to explain before making final predictions, and only give a positive reward if it provides a correct prediction. Even though there are no human-annotated explanations as ground truth to guide the LLM training, the researchers found that LLM can still learn to provide human-understandable explanations and achieve high accuracy in solving reasoning problems. To ensure the efficiency and faithfulness of LLM-generated explanations, recent works focus on encouraging LLMs to provide a shorter but necessary explanation [120, 184, 309].

## 6.2 Case Study: Is CoT Really Making LLM Inferences Explainable?

*6.2.1 Background and Experimental Settings.* Despite the apparent intuitiveness of the CoT prompt design, a critical question remains unanswered: *Does CoT really make LLM inferences explainable?* In other words, can the information provided through CoT faithfully reflect the underlying generation process of LLMs? We use multi-hop question-answering (QA) as the scenario to investigate this problem. In QA systems, answering multi-hop questions remains a significant challenge. Instead of leveraging a single information source, multi-hop questions require synthesizing information from multiple pieces or sources of data into a coherent and logical sequence [139, 255, 346].

**CoT Prompts for Multi-hop QA.** To address this, our case study applies the CoT technique with in-context examples. We provide an example as follows, where $[x]$ denotes the question. The "Thoughts" following each "Question" are step-by-step problem-solving statements for the multi-hop questions.

**CoT Faithfulness for Explanation:** To quantitatively measure the faithfulness of CoTs, we select fidelity as the corresponding metrics [269, 338]:

$$Fidelity = \frac{\sum_{i=1}^{N} \left( \mathbb{1}\left(\hat{\mathbf{y}}_i = \mathbf{y}_i\right) - \mathbb{1}\left(\hat{\mathbf{y}}_i^{mislead} = \mathbf{y}_i\right) \right)}{\sum_{i=1}^{N} \left( \mathbb{1}\left(\hat{\mathbf{y}}_i = \mathbf{y}_i\right) \right)} \times 100\%,$$

where $\mathbf{y}_i$ denotes the ground truth label, $\hat{\mathbf{y}}_i$ denotes the original model output with CoT, while $\hat{\mathbf{y}}_i^{mislead}$ denotes the model output with misleading information inserted. In the following, we give an example. Given the target question, the correct step-by-step thoughts should be: "*Ellie Kemper is a citizen of the United States of America. The president of the United States of America is Joe Biden.*" To mislead the model, we replace the thoughts with incorrect information (the underlined text) and ask the model to generate a new answer based on incorrect thoughts. If the model still generates the correct answer after the modification, we believe that the CoT information does not faithfully reflect the true process of the answer generation. On the other hand, if it generates an answer corresponding to the incorrect thoughts, then we claim the thoughts are faithful.

**Experimental Settings.** We evaluate the performance on the MQUAKE-CF dataset [346], which includes 1,000 cases for each $K$-hop questions, $K \in$

**Question:** What is the capital of the country where Plainfield Town Hall is located?
**Thoughts:** Plainfield Town Hall is located in the country of the United States of America. The capital of the United States is Washington, D.C.
**Answer:** Washington, D.C.
...

**Question:** Who has ownership of the developer of the Chevrolet Corvette (C4)?
**Thoughts:** The developer of Chevrolet Corvette (C4) is Chevrolet. Chevrolet is owned by General Motors.
**Answer:** General Motors

**Question:** [x]

Fig. 5. CoT Prompts for Multi-hop QA.

**Question:** What is the capital of the country where Plainfield Town Hall is located?
**Thoughts:** Plainfield Town Hall is located in the country of the United States of America. The capital of United States is Washington, D.C.
**Answer:** Washington, D.C.
...

**Question:** Who has ownership of the developer of the Chevrolet Corvette (C4)?
**Thoughts:** The developer of Chevrolet Corvette (C4) is Chevrolet. Chevrolet is owned by General Motors.
**Answer:** General Motors

**Question:** [Who is the head of state of the country where Ellie Kemper holds a citizenship?]
**Thoughts:** Ellie Kemper is a citizen of Croatia. The head of state in Croatia is Zoran Milanović.
**Answer:**

Fig. 6. CoT Faithfulness for Explanation.

{2, 3, 4}, which totally consists of 3,000 questions. Our evaluation applies various language models, including GPT-2 (1.5B) [224], GPT-J (6B) [272] with 6 billion parameters, LLaMA (7B) [262], Vicuna-v1.5 (7B) [47], LLaMA2-chat-hf (7B) [263], Falcon [8] with 7 billion parameters, Mistral-v0.1 (7B) [134], and Mistral-Instruct-v0.2 (7B) [134]. These models have demonstrated proficiency in both language generation and comprehension.

*6.2.2 Experiment Results.* **Faithfulness Evaluation of CoT.** Table 3 illustrates the impact of accurate versus misleading CoTs on the performance of LLMs. The Fidelity metric indicates how faithfully the model's output reflects the reasoning process described in the CoT. Ideally, a high Fidelity score suggests that the model's final response is directly based on the provided CoT,

Table 3. CoT Faithfulness Evaluation on MQUAKE-CF.

| Datasets | MQUAKE-CF | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Question Type | 2-hops | | | 3-hops | | | 4-hops | | |
| Edited Instances | C | M | F | C | M | F | C | M | F |
| GPT-2 (1.5B) | 15.9 | 5.2 | 67.3% | 8.9 | 2.9 | 67.4% | 8.4 | 1.3 | 84.5% |
| GPT-J (6B) | 51.9 | 7.3 | 85.9% | 30.5 | 1.8 | 94.1% | 49.8 | 2.0 | 96.0% |
| LLaMA (7B) | 65.1 | 9.9 | 84.8% | 39.3 | 6.1 | 84.5% | 62.9 | 6.0 | 90.5% |
| Vicuna-v1.5 (7B) | 56.3 | 21.7 | 61.5% | 29.7 | 12.7 | 57.3% | 53.1 | 16.1 | 69.7% |
| LLaMA2 (7B) | 58.7 | 17.0 | 71.0% | 30.3 | 8.3 | 72.5% | 49.1 | 12.0 | 75.6% |
| Falcon (7B) | 61.7 | 24.0 | 61.1% | 31.6 | 15.0 | 52.6% | 48.6 | 23.1 | 52.4% |
| Mistral-v0.1 (7B) | 69.3 | 24.0 | 65.4% | 42.3 | 13.0 | 69.3% | 63.2 | 18.4 | 70.8% |
| Mistral-v0.2 (7B) | 56.3 | 47.9 | 14.8% | 37.7 | 22.0 | 41.6% | 56.2 | 37.3 | 33.6% |

*C: CoT; M: Misleaded CoT; F: Fidelity Score.*

validating it as a faithful explanation of the model's reasoning pathway. However, as we will discuss below, a low Fidelity may not always imply a lack of faithfulness in the model's reasoning, which calls for developing more effective evaluation methods in future research.

GPT-J and LLaMA exhibit high fidelity scores across different question types, indicating a strong adherence to the given reasoning paths. Conversely, other models show relatively high mislead accuracy scores with lower fidelity scores. In the experiments, we observe that these models usually rely on their own generated thoughts instead of using incorrect information provided in the CoT. Mistral-v0.2, in particular, demonstrates the lowest fidelity scores and highest misleading accuracy scores, suggesting a potential self-defense ability against false information. The lower fidelity scores of later models may be attributed to their improved training processes on more diverse and high-quality datasets, enabling them to develop a better understanding of context and reasoning. As a result, they are more likely to generate their own correct reasoning paths.

## 6.3 Challenges

An explanation is considered as faithful if it causes the model to make the same decision as the original input [166]. In this context, the challenge faced by explainable prompting (e.g., CoT prompt) lies in two aspects: (1) directing language models to generate explanations that are genuinely representative of the models' internal decision-making processes, and (2) preventing language models from depending on potentially biased CoT templates.

Regarding the first challenge, our case study has revealed that relatively small language models may generate answers that do not align with the provided CoT rationales. Therefore, these rationales do not accurately represent the decision-making process within these models. Some efforts have been made to bolster the CoT capabilities of smaller language models by implementing instruction tuning with CoT rationales [119, 140]. Nevertheless, it remains a challenging problem of how to ensure the generated explanations (i.e., "what the model says") are faithful to the internal mechanism (i.e., "what the model thinks") of language models. Regarding the second challenge, recent research shows that explanations in the CoT can be heavily influenced by the introduction of biasing prompt templates into

model input [264]. If incorrect or biased information is encoded in such templates, the generated explanations could be misleading. Recently, [282] propose a novel decoding strategy to implement CoT with prompting, which could mitigate this issue. However, how to effectively help language models get rid of the template reliance still remains underexplored.

## 7 EFFECTIVE LLM INFERENCE VIA KNOWLEDGE-AUGMENTED PROMPTING

Enhancing models with external knowledge can significantly improve the control and interpretability of decision-making processes. While LLMs acquire extensive knowledge through pre-training on web-scale data, this knowledge is embedded implicitly within the model parameters, making it challenging to explain or control how this knowledge is utilized during inference. To address these limitations, this section discusses Retrieval-Augmented Generation (RAG) for the explicit integration of external knowledge into LLMs, aiming to yield more interpretable predictions.

### 7.1 Retrieval-Augmented Generation with Knowledge

Generally, RAG operates in two steps: (1) *Retrieval*: It locates and fetches pertinent information from an external source based on the user's query; (2) *Generation*: It incorporates this information into the model's generated response. Given an input query $x$ and the desired output $y$, the objective function of RAG can be formulated as [100]: $\max_{\phi,\theta} \log p(y|x) = \max_{\phi,\theta} \log \sum_{z \in \mathcal{K}} p_\phi(y|x,z) \cdot p_\theta(z|x)$, where $z$ stands for the external knowledge retrieved from a knowledge base $\mathcal{K}$. Thus, the target distribution is jointly modeled by a knowledge retriever $p_\theta(z|x)$ and an answer reasoning module $p_\phi(y|x,z)$. The knowledge $z$ serves as a latent variable. An RAG model is trained to optimize the parameters, so that it learns to retrieve relevant knowledge $z$ and to produce correct answers $y$ based on $z$ and $x$. As LLMs possess stronger text comprehension and reasoning abilities, they can directly serve as the reasoning module $p_\phi$ without further training. In this case, RAG can be treated as a data-centric problem:

$$\max_{z \in \mathcal{K}} \log p(y|z,x) = \max_{z \in \mathcal{K}} \frac{p(z|x,y)}{p(z|x)} p(y|x), \tag{12}$$

where the goal is to find appropriate knowledge that supports the desired output. The *interpretability* of RAG-based models comes from the information in $z$: (1) $z$ usually elucidates or supplements the task-specific information in $x$; (2) $z$ could explain the generation of output $y$. Unlike other models that estimate $p(y|x)$ in an end-to-end manner, where the decision process is not comprehensible, the RAG process provides justification or rationale $z$ that supports the result.

### 7.2 Enhancing Decision-Making Control with Explicit Knowledge

*7.2.1 Reducing Hallucinations in Response.* "Hallucination" in the context of LLMs refers to instances where these models generate information that, while coherent and contextually appropriate, is not based on factual accuracy or real-world evidence [125]. This issue can lead to the production of misleading or entirely fabricated content, posing a significant challenge to the reliability and trustworthiness of LLMs' outputs. RAG offers a powerful solution to mitigate the problem of hallucinations in LLMs. By actively incorporating up-to-date, verified external knowledge at the point of generating responses, RAG ensures that the information produced by the model is anchored in reality. This process significantly enhances the factual basis of the model's outputs, thereby reducing the occurrence of hallucinations. [247] applies neural-retrieval-in-the-loop architectures to knowledge-grounded dialogue, which significantly reduces factual inaccuracies in chatbots, as confirmed by human evaluations. [250] introduces RAG-end2end, which joint trains retriever and generator components together. Their method demonstrates notable performance improvements across specialized domains like healthcare and news while reducing knowledge hallucination.

*7.2.2 Dynamic Responses to Knowledge Updating.* RAG empowers LLMs by incorporating the most up-to-date information, keeping their decision-making aligned with the latest developments. This feature is especially vital in fast-evolving fields such as medicine and technology, where the need for timely and accurate information is paramount [197]. For example, research by [129] demonstrates significant enhancements in output relevance and accuracy through real-time information retrieval. Similarly, [101] suggest using retrieved factual data to correct and update the knowledge within pre-trained LLMs efficiently. Additionally, [280] introduce a method for integrating newly retrieved knowledge from a multilingual database directly into the model prompts, facilitating updates in a multilingual context.

*7.2.3 Domain-specific Customization.* RAG enhances LLMs by incorporating knowledge from specialized sources, enabling the creation of models tailored to specific domains. Research by [100] illustrates how integrating databases specific to certain fields into the retrieval process can empower models to deliver expert-level responses, boosting their effectiveness in both professional and academic contexts. [244] have applied this concept in the medical domain with MedEdit, utilizing an in-context learning strategy to merge relevant medical knowledge into query prompts for more accurate medical advice. Moreover, recent research finds that LLMs struggle to capture specific knowledge that is not widely discussed in the pre-training data. Specifically, [192] observe that LLMs often fail to learn long-tail factual knowledge with relatively low popularity, finding that simply increasing model size does not significantly enhance the recall of such information. However, they note that retrieval-augmented LLMs surpass much larger models in accuracy. Similarly, [135] highlights LLMs' challenges with acquiring rare knowledge and proposes that retrieval augmentation offers a viable solution, minimizing reliance on extensive pre-training for capturing nuanced, less common information.

## 7.3 Challenges

We discuss the challenges in RAG that are relevant to its explainability aspects: (1) In the retrieval stage $p_\theta(z|x)$, does the retrieved information $z$ always elucidate the task-specific information contained in the input $x$? (2) In the generation stage $p_\phi(y|x, z)$, does $z$ effectively serve as an explanation for the generation of output $y$?

*7.3.1 Retrieval Accuracy Bottlenecks.* Existent RAG methods typically rely on similarity search to pinpoint relevant information [82, 153], which represents a substantial improvement over basic keyword searches [231]. However, these methods may struggle with complex queries that demand deeper comprehension and nuanced reasoning. The recent "lost-in-the-middle" phenomenon [175] has revealed that an ineffective retrieval can result in the accumulation of extraneous or conflicting information, negatively affecting the generation quality. Efficiently handling intricate and multi-hop questions remains a significant challenge, highlighting the need for ongoing research to enhance the capabilities of RAG systems.

*7.3.2 Controllable Generation Bottlenecks.* In-context learning stands out as the premier method for incorporating external knowledge to boost the capabilities of LLMs such as GPT-4 [14, 82]. Despite its effectiveness, there's no surefire way to ensure that these models consistently leverage the provided external knowledge within the prompts for their decision-making processes [154, 216, 295, 320]. The challenge of optimizing the use of external explanations to achieve more precise and controlled decision-making in LLMs is an ongoing issue that has yet to be fully addressed.

## 8 TRAINING DATA AUGMENTATION WITH EXPLANATION

This section explores the generation of synthetic data from explanations using LLMs, a technique poised to enhance various machine learning tasks. In machine learning, limited data availability often leads to data imbalance and scarcity,

which in turn constrains model performance in many domains. A viable solution is data augmentation, where LLMs, with advanced generative capabilities, can be utilized [256, 293]. Nevertheless, several challenges exist for effective text augmentation. First, for utility, the generated samples should exhibit diversity compared to the original data. Second, these samples should exhibit useful patterns relevant to downstream tasks. To this end, explanation tools could guide augmentations by providing supplemental context and rationales [37].

Explanations can be particularly beneficial in two scenarios. In the *first scenario*, explanation tools are used to to identify existing deficiencies, which effectively guides the data augmentation. The augmented data include samples and related label annotations. The *second scenario* employs LLMs to directly produce explanatory texts, such as rationales, as supplementary information to enrich the dataset.

## 8.1 Label-based Data Augmentation

For learning tasks, the training dataset includes a large volume of data points, defined as $\mathcal{D} = \{(x_1, y_1), \cdots, (x_n, y_n)\}$. We categorize learning tasks into classification tasks and generation tasks.

In traditional classification tasks, each data point $x_i$ has a label $y_i$ representing a class. Models are trained to predict the labels for new data points. However, these models are prone to make predictions with spurious correlations also known as shortcuts [85]. Thus, model's predictions are predominantly based on shortcut features, and the underlying mechanisms of the model are not interpretable from a human perspective. Post-hoc explanation techniques play a crucial role in detecting undesirable correlations between input and predictions [172, 173]. For example, Du et al. [66] adopt Integrated Gradient (IG) to attribute a model's predictions to its input features, showing that the model tends to treat functional words, numbers, and negation words as shortcuts and rely on them for prediction in natural language understanding tasks.

LLMs can be applied to resolve data set issues by augmenting various data. For example, explanatory information such as counterfactuals has been incorporated to improve model robustness [284] and out-of-distribution performance [234]. LLMs can also synthesize rare examples, which helps models to generalize better on unseen data [304]. Besides, another work utilizes attribution-based interpretability tool to identify vulnerable words and leverage LLMs to synthesize natural adversarial examples to enhance model safety [286]. Similarly, LLMs are also helpful in mitigating biases such as fairness issues in models. He et al. [111] claims that it automatically identifies underrepresented subgroups, and use LLMs to augment theese subgroups while avoids hurting other groups.

In text generation tasks, each data point consists of a prompt $x$ and a response $y$. The data used for model training is usually text including both prompt and response. Generative models are designed to predict next token based on existing input, which can be formulated as $y_m : \arg\max p_\theta(y_m | x, y_{(i\,:\,m-1)})$. To make models generate desired output, well-aligned models are indispensable. Recently, a few work employ LLMs to augment data and enhance model alignment. For example, Wang et al. [273] identifies the weakness of existing datasets and instructs LLMs to generate desired data to improve model safety. Another work use attribution-based method to identify key information from privacy-relevant data and replace the privacy information with synthetic information, which provides a way to avoid privacy leakage [330].

## 8.2 Rationale-based Data Enrichment

LLMs have been leveraged to directly generate natural language explanations as augmented data. This strategy makes use of LLMs' chain-of-thought reasoning abilities [144]. The training datasets contain data with annotated information, denotes as $\mathcal{D} = \{(x_1, r_1, y_1), \cdots, (x_m, r_m, y_m)\}$. Each data point includes not only prompt $x$ and response $y$, but

reasoning process or annotated information defined as $r$. As state-of-art LLMs are excellent in interpreting context, many studies have utilized them to provide annotations and enrich data.

One line of research uses generated explanations to assist and guide model training. Li et al. [164] introduce LLM explanations to facilitate the training of smaller models, which improves their reasoning capabilities and acquires explanation generation abilities. Hsieh et al [121] integrate LLM rationales as additional supervision to guide the training of smaller models. Experiments have shown that this approach not only requires fewer training data but also outperforms traditional fine-tuning and distillation methods. Another line of study employ annotated information to boost models' performance on specific tasks. For example, LLM explanations have also been utilized to generate explanations for sentiment labels of aspects in sentences to mitigate spurious correlations in aspect-based sentiment analysis tasks [278]. Besides, Code translation generation tasks incorporate explanations as an intermediate step, improving model performance by 12% on average [257]. The result shows that explanations are particularly useful in zero-shot settings. Wei et al. [291] instruct LLMs to explain the retrieval process from documents to produce rationales. Then they can be used to train models and further improve models ability on retrieval-augmented generation tasks. Similarly, Lupidi et al. [187] prompt LLMs to generate data explanations from real-world sources such as articles, which are then used to fine-tune models and improve their performance on multi-hop reasoning and information retrieval.

### 8.3 Challenges

*8.3.1 Computational Overhead.* Conventional post-hoc explanations, built on well-trained models, are often resource-intensive tasks. The first scenario mentioned above leverages interpretability techniques to accurately diagnose dataset issues. This process typically requires multiple rounds of model training and applying interpretability methods to develop fair and robust models. Consequently, the crafting process can be both time- and energy-consuming.

*8.3.2 Data Quality and Volume.* Despite their advanced capabilities, LLMs still have limitations when dealing with highly specialized or niche contexts. For example, one of the most prominent issues is "hallucination", where models generate plausible but incorrect or misleading responses. This could adversely affect the quality of augmented data, potentially introducing more biases to which LLMs are also vulnerable. Determining the precise amount of data required is also challenging, often leading to new dataset imbalances. Managing the quantity of LLM-generated data is an immense challenge, as augmented data can introduce other biases [341]. This stems from LLMs' limited ability to accurately control the quantity and distribution of generated data. Together, these factors underscore the complexities and challenges in fully harnessing LLMs' potential for data augmentation and related tasks.

## 9 GENERATING USER-FRIENDLY EXPLANATION WITH LLMS

Previous sections mainly focused on quantitative explanations with LLM via numerical values. For example, sample-based explanation discussed in Section 4 aims to assign each training sample an influence score (see Eqs.8-11) that measures the confidence that we can use that training sample to explain the prediction of a test sample. However, using numerical values for explanations is not intuitive, which can be difficult to understand by practitioners with little domain knowledge [41, 148, 150, 160]. User-friendly explanations, on the contrary, aim to generate human-understandable explanations, e.g., natural language-based descriptions, regarding why a model makes certain predictions or what role a neuron plays in the network, such that the explanations can be well-understood by both researchers and practitioners.

## 9.1 Explaining Small Models with LLMs

Recently, there has been growing interest in leveraging LLMs to generate free-text explanations for small models. For example, to explain black-box text classifiers, [26] propose a prompting-based strategy to identify keywords $K = \{k_1, k_2, ..., k_n\}$ in the input texts $x$ with pretrained LLMs that are informative for the label $y$, and ask LLMs to substitute them with another set of keywords $K' = \{k'_1, k'_2, ..., k'_n\}$, such that changed text $x'$ changes the label prediction to $y'$. They view the textual mapping rule "if we change $K$ into $K'$ in $x$, then $y$ will be classified as $y'$" as the counterfactual explanation for the model. In addition, to explain the neuron of a pretrained language model (e.g., GPT2), Bills et al. [27] propose to summarize the neuron activation patterns into *textual phrases* with a larger language model (e.g., GPT4), where the neuron activation patterns are expressed as a sequence of (token, attribution score) pairs. To verify the identified patterns, they generate activation patterns according to the phrases via the same LLM and compare their similarity with the true activation patterns of the neuron, where the phrases with high scores are considered more confident to serve as the explanation for the neuron.

The explaining ability of LLMs is not necessarily limited to text models. For example, [337] propose using pretrained vision-language models to generate explanations for a neuron $\theta_i$ of an image classification model. Specifically, for each class $y = y_c$, they first find regions in images with label $y_c$ that have maximum activation of the neuron $\theta_i$ as the surrogate explainees for $\theta_i$, and prompt LLMs such as ChatGPT to generate candidate explanations (words, short phrases) for the class label $y_c$. Then, they use the pretrained vision-language model CLIP [223] to match the candidate explanations with the surrogate explainees as the explanations for the neuron $\theta_i$. Recently, LLMs have also found applications in explaining recommender systems [349]. Specifically, [313] found that LLMs can well interpret the latent space of sequential recommendation model after alignment, whereas [151] propose to align user tokens of LLMs with the learned user embeddings of small recommendation model to generate explanations of user preferences encoded in the embeddings. Recently, [239] propose a unified framework to explain all models where inputs and outputs can be converted to textual strings. Specifically, the explainer LLM is used as an agent to interact with the explainee model by iteratively creating inputs and observing outputs from the model, where the textual explanations are generated by viewing all the interactions as the context.

## 9.2 Self-Explanation of LLMs

Due to the black-box nature of LLMs, it is promising to generate user-friendly explanations for the LLMs themselves, such that their operational mechanics and the predictions can be understood by human experts. Based on whether the LLM needs to be retrained to generate explanations for themselves, the self-explanation of LLM can be categorized into two classes: *fine-tuning based* approach and *in-context based* approach, which will be introduced in the following parts.

**Fine-tuning based approaches.** Given sufficient exemplar explanations on the labels of the training data (e.g., in recommendation datasets such as the Amazon Review datasets [110] or the Yelp dataset [347], users have provided explanations on why they have purchased certain items, which can be viewed as explanations for the ratings), LLMs can learn to generate explanations for their predictions as an *auxiliary task* through supervised learning. One exemplar method is P5 [86], which fine-tunes the pre-trained language model T5 [225] on both the rating and explanation data to generate an explanation alongside the recommendations. Recently, several works have improved upon P5 [49, 350], which fine-tunes different LLMs such as GPT2, LLaMA, Vicuna, etc., and propose different prompt learning strategies [161] with generating explanation as the auxiliary task. With explanations introduced as additional supervision signals to fine-tune pretrained LLMs for recommendations, the performance can be improved with good explainability.

**In-context based approaches.** In many applications, there is often a lack of sufficient exemplar explanations. However, the unique capability of modern LLMs to reason and provide answers through human-like prompts introduces the potential for in-context based explanations. Here, explanations for predictions are crafted solely based on the information within the prompt. A leading approach in this domain is the Chain-of-Thoughts (CoT) prompting [289], which provides few-shot examples (with or without explanations) in the prompt and asks the LLM to generate answers after reasoning step-by-step, where the intermediate reasoning steps that provide more context for generating the final answer can be viewed as explanations. However, CoT generates reasoning first and then based on which generates predictions, where the reasoning steps can influence prediction results [188]. If explanations are generated after the prediction, since the explanation is conditioned on the predicted label, it can provide a more faithful post-hoc explanation of why the model makes certain decisions [147]. The application of in-context based self-explanation of LLMs is broad. For example, [127] explore generating zero-shot self-explanation of sentiment analysis with LLMs by directly asking them to generate explanations alongside the predictions. In addition, [123] propose a chain-of-explanation strategy that aims to explain how LLMs can detect hate speech from the textual input. [182] find that CoT can generate well-supported explanations for question answering with scientific knowledge.

### 9.3 Challenges

*9.3.1 Usability v.s. Reliability.* Many existing methods rely on prompts to generate user-friendly explanations, which are not as reliable as numerical methods with good theoretical foundations. [316] find that the explanations by CoT may not be factually grounded in the inputs. Therefore, they believe that these explanations are more suitable as post-hoc explanations regarding why the LLM makes certain predictions. However, the validity of viewing CoT explanations as post-hoc justifications has been questioned by recent findings [264], which uses biased datasets (e.g., the few-shot examples in the prompt always answer "A" for multiple choice questions) to show that the generated explanations may be plausible, but systematically unfaithful to represent the true reasoning process of the LLMs. Therefore, there's a growing need for more theoretical scrutiny of user-friendly explanations to ensure faithfulness and credibility.

*9.3.2 Constrained Application Scenarios.* Currently, the utilization of LLMs to explain smaller black-box models is mainly limited to those that deal with data with rich textual information [26, 151]. Although [338] propose a strategy to explain image classifiers, the ability to match candidate textual explanations with image patterns still relies on the pretrained vision-language model CLIP. Therefore, there is a compelling need to development more versatile strategies for explaining models across a wider range of fields. This endeavor could depend on the fundamental research on combining LLM with other domain-specific tasks, such as the development of Graph-Language Models that are applicable to unseen graphs in a zero-shot manner.

## 10 INTERPRETABLE AI SYSTEM DESIGN WITH LLMS

An intriguing but challenging problem in XAI is creating model architectures or even AI systems that are inherently interpretable [232], where different model components represent clear and comprehensible concepts or functionalities that are easily distinguishable from one another. Machine learning models such as support vector machines [114] and tree-based models [251] were classical techniques for achieving model interpretability. In the deep learning era, typical research areas in this context include concept-bottleneck models [143, 327], disentangled representation learning [56, 118], and network dissection [19, 20]. Nevertheless, under the traditional deep learning setting, the usability of these techniques remains limited because of two major challenges. First, it is difficult to define the spectrum of

concepts or functionalities the model is expected to capture. Second, the efficacy of interpretable models often falls short compared to black-box models, thereby constraining their practical utility.

Large foundation models, such as large language models (LLMs) and vision language models (VLMs), provide opportunities to bridge the gap. By leveraging the common-sense knowledge embedded within them, foundation models can *design interpretable architectures* by providing cues that encourage creating and using the features or procedures within AI workflows. This is different from traditional deep learning pipelines, where the deep models automatically discover the features during the training process, which may not end up with model components with clear meanings. Furthermore, LLMs can decompose complex tasks into simpler and collaborative sub-tasks, enhancing both the system's interpretability and its overall performance.

### 10.1 Designing Interpretable Network Architectures with LLMs

Representative methods for developing interpretable deep architectures include Generalized Additive Models (GAMs) [3, 181, 351] and Concept Bottleneck Models (CBMs) [143, 242, 327]. These models map inputs into a human-understandable latent space, and then apply a linear transformation from this space to the target label. For example, to build a classifier that diagnoses arthritis, we can let the model identify features such as "bone spurs" and "sclerosis", and then use these interpretable features for the final decision. However, these approaches often require the involvement of experts to define the latent space, which can limit the learning capabilities of deep models. Some work tries to automate the discovery of semantic concepts during model training, such as by requiring independence between concepts [118, 325] or clustering data [89], but they lack direct control over the outcomes and does not ensure the clarity of the concepts. One promising strategy is to utilize LLMs to provide comprehensible concept candidates. [198] use human language as an internal representation for visual recognition, and create an interpretable concept bottleneck for downstream tasks. By basing the decision on those comprehensible concepts, the model architecture itself is provided with better transparency. Similarly, a recent approach *Labo* [312] constructs high-performance CBMs without manual concept annotations. This method controls the concept selection in bottlenecks by generating candidates from the LLMs, which contain significant world knowledge [217] that can be explored by prompting a string prefix. Human studies further indicate that those LLM-sourced bottlenecks are much factual and groundable, maintaining great inherent interpretability for model designs. Besides the concept-based models, another promising strategy is to employ LLMs to enhance the conventional architectures that are inherently interpretable, such as GAMs and Decision Trees (DTs). [248] leverages the knowledge captured in LLMs to enhance GAMs and DTs, where LLMs are only involved during the augmented model training instead of the inference process. For GAMs training, LLMs can provide decoupled embeddings for enhancement. For DTs training, LLMs are able to help generate improved features for splitting. The LLM-augmented GAMs and DTs enable full transparency, where only the summing coefficients and input key phrases are required for interpretation. With the extra information from LLMs, augmented GAMs and DTs are capable of achieving better generalization performance compared with non-augmented ones.

### 10.2 Designing Interpretable AI Workflows with LLM Agents

Traditional deep models are usually designed in an end-to-end manner. The internal workflows are not quite understandable to general users. By utilizing common-sense world knowledge, LLMs can break down complex problems into smaller ones and organize the workflows among them, leading to a more interpretable design of AI systems [74]. A recent example of interpretable AI workflow design comes from [243], where an LLM-powered agent leverages ChatGPT to integrate various off-the-shelf AI models (e.g., from Hugging Face [131]) to handle different downstream application

tasks. In order to handle complicated tasks in a transparent workflow, LLMs play a pivotal role in coordinating external models with language mediums to harness their powers. By planning the target task, selecting candidate models, executing decomposed subtasks and summarizing responses, LLMs can help disassemble tasks based on user requests, and assign appropriate models to the tasks based on the model descriptions. Similarly, to transparentize the workflow, [179] introduces a task decomposer to analyze the user prompts and break it down into a number of subtasks for solving using LLMs. Each subtask is well managed and attributed with *description*, *domain*, *inputs*, and *outputs*. In this way, the AI systems are then capable of handling intricate user prompts with a step-by-step understandable workflow. Under the prompting paradigm, [137] also employs LLMs to solve complex tasks by decomposition. Drawing inspiration from software libraries where the workflows are trackable, the decomposer and shared subtasks are designed in a modular manner. One step further, [283] introduces an interactive planning approach for complex tasks, which enhances the error correction on initial LLM-generated plans by integrating plan execution descriptions and providing self-explanation of the feedback. Such an interactive nature enables better workflow transparency in long-term planning and multi-step reasoning task scenarios.

## 10.3 Challenges

*10.3.1 Planning Feasibility in Complicated Scenarios.* Despite the task planning capability of LLMs, it is still challenging to apply to certain scenarios in real-world applications due to feasibility issues. One typical scenario is the few-shot planning cases [95], where acquiring large datasets for training is either impractical or cost-prohibitive, thus making feasible planning about unseen cases from sparse exemplars extremely challenging. To better assist the interpretable designs, LLM planning needs to generalize well without extensive supervision and is expected to have the ability to integrate information from prior experiences as well as knowledge. Besides, another important scenario lies in the dynamic planning settings [51], in which LLMs integrate feedback from the environment iteratively, letting the agent take thinking steps or augment its context with a reasoning trace. Dynamic scenarios urgently and frequently involve high computational costs resulting from the iterated invocations of LLMs, and still face challenges in dealing with the limits of the context window and recovering from hallucinations on planning.

*10.3.2 Assistance Reliability with Knowledge Gaps.* LLMs exhibit remarkable proficiency in encapsulating real-world knowledge within their parameters, but they resort to hallucinations and biases with high confidence when certain knowledge is missing or unreliable. As a result, a crucial research challenge keeps rising, *i.e.*, how to effectively detect and mitigate the LLM knowledge gaps from humans when employing LLMs for designs. We will need further research on evaluating and developing robust LLM mechanisms to address the knowledge-gapping problems, with the goal of helping improve LLM reliability, reducing hallucinations and mitigating biases. Furthermore, the intersections between the knowledge gaps and the safety aspects are also a great challenge to be solved, which may pose some security concerns, especially when using LLMs for downstream models or workflow designs.

## 11 EMULATING HUMANS WITH LLMS FOR XAI

This section discusses how LLMs can play the role of humans in serving XAI development. Building explainable models requires two main steps where humans are in the loop: (1) collecting a dataset with human-annotated rationales to train the models; (2) collecting human feedback on the quality of explanations produced by the models for evaluation. The significant cost and time required for human involvement raise the main challenge in scaling up this procedure. LLMs emerge as a promising solution to this challenge, thanks to their capability to emulate human reasoning and

produce responses that closely resemble human-generated content. In the following, we introduce the methods that demonstrate LLMs' ability to generate human-like annotations and feedback for XAI.

### 11.1 Emulating Human Annotators for Training Explainable Models

Incorporating human-understandable rationales into model development has shown its effectiveness in enhancing both the transparency and performance of the system for various NLP tasks, such as question answering [163, 294], sentiment analysis [9, 65], and common sense reasoning [33, 226]. We use the term *rationales* to describe supportive evidence that justifies the connection between inputs and outputs [99]. Traditionally, the rationales are collected by leveraging human annotations [34, 285] or applying expert-designed rules [7, 163], resulting in expensive costs or limited quality. Recently, researchers in **automatic annotation** [21, 61, 90] have begun to explore the potential of leveraging advanced LLMs to emulate human annotators in annotating the target labels of task-specific examples. These studies found that advanced LLMs show comparable annotation qualities against average crowd human annotators on most tasks with a lower cost, pointing out the scalability of using machine-emulated annotators. Inspired by these works, some studies [123, 124, 213, 296, 311, 335] attempt to leverage advanced LLMs to collect rationales by applying the **chain-of-thought** prompting. Specifically, researchers provide several input-rationale-output demonstrations within the input text to prompt the LLMs to generate rationale and output for an unlabeled input instance. The quality of such annotated rationales largely relies on the in-context learning capabilities of LLMs, leading to uncontrollable annotation quality on uncommon tasks. Other scholars [44, 186, 314] propose a **human-in-the-loop** LLM-based annotation framework based on the **active-learning** architecture. This framework initially collects a small seed dataset with human-annotated rationales and labels. This seed dataset is used to train an explainable classifier for this downstream task. Then, each unlabeled sample is passed through the trained explainable classifier. This is followed by a selection strategy that chooses representative samples according to metrics such as explanation plausibility, prediction uncertainty, and sample diversity. Finally, LLMs are leveraged to annotate the rationales and labels of these selected unlabeled samples. This procedure could be repeated multiple times, and the trained explainable classifier from the latest time is the final output of this framework. Compared with other methods, this approach balances the annotation quality and the cost budget in developing explainable models by using LLM-emulated annotators.

### 11.2 Emulating Human Feedback for Evaluating Explainable Models

The explanations generated by the explainable models could be classified into two categories: extractive and abstractive [99]. Extractive explanations derive directly from the input data, exemplified by attribution-based methods that emphasize specific segments of the input text. In contrast, abstractive explanations are generated in a free-form text manner, such as chain-of-thought (CoT) responses [289], offering a more nuanced interpretation. The quality of extractive explanations is typically assessed through their **agreement with human-annotated rationales** [58], such as accuracy, recall, and precision. However, evaluating abstractive explanations presents a significant challenge, as it is impractical to exhaustively evaluate all reasonable abstractive results comprehensively. To automatically assess abstractive explanations, early studies first collect some free-text rationales, and then compute **the cosine similarity between embeddings of explanations and rationales** [45, 157]. A higher similarity between the abstraction explanation and the annotated rationales indicates a more transparent model. Recently, some researchers have followed the **LLM-as-a-judge** framework to ask LLM to judge the model explanations without referring to any human-annotated rationales [27, 199, 235, 288, 345], emphasizing the potential of emulating human feedback with advanced LLMs.

## 11.3 Challenges

*11.3.1 Uncontrollable Credibility of Emulation.* While LLMs could assist in rationale collection and explanation evaluation, their provided results can be biased and may not always match human annotators, due to hallucinated responses in their unfamiliar domains [133] and that LLMs are still immature [106]. The quality of data gathered from this process is compromised, impacting the development of XAI systems. To improve the quality of annotations and feedback, future research could focus on incorporating hallucination detection [59] and retrieval augmented generation [228] techniques. These methods could enhance the reliability of LLM outputs, making them more comparable to human-generated content in the context of XAI development.

*11.3.2 Ethical Considerations in LLM Annotation.* When LLM annotators keep human annotators away from subjective scenarios, such as hate speech detection [124], LLMs also have a chance to inject unethical opinions into their annotated datasets. Although most advanced LLMs are fine-tuned to align with human values [211], such as being helpful, honest, and harmless, many studies have shown that this protection mechanism can be jailbroken [287, 353], causing the model to produce values-violating answers. Ensuring LLM annotators follow ethical guidelines is worth further exploration.

## 12 CONCLUSION

In this paper, we discuss a crucial yet frequently underappreciated aspect of Explainable AI (XAI) – *usability*. To this end, we present 10 strategies for advancing Usable XAI within the LLM paradigm, including (1) leveraging explanations to reciprocally enhance LLMs and general AI systems, and (2) enriching XAI approaches by integrating LLM capabilities. Unlocking the potential of XAI's usability can help address various challenges in LLM such as human alignment. We also provide case studies to several critical topics, aiming to provide resources for interested developers. We further discuss open challenges at the end of each strategy, suggesting directions for future work in this evolving area.

## REFERENCES

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv arXiv:2303.08774* (2023).

[2] Julius Adebayo, Melissa Hall, Bowen Yu, and Bobbie Chern. 2023. Quantifying and mitigating the impact of label errors on model disparity metrics. *arXiv arXiv:2310.02533* (2023).

[3] Rishabh Agarwal, Levi Melnick, Nicholas Frosst, Xuezhou Zhang, Ben Lengerich, Rich Caruana, and Geoffrey E Hinton. 2021. Neural additive models: Interpretable machine learning with neural nets. *Advances in neural information processing systems* 34 (2021), 4699–4711.

[4] Ayush Agrawal, Mirac Suzgun, Lester Mackey, and Adam Kalai. 2024. Do Language Models Know When They're Hallucinating References?. In *Findings of the Association for Computational Linguistics: EACL 2024.* 912–928.

[5] Gustaf Ahdritz, Tian Qin, Nikhil Vyas, Boaz Barak, and Benjamin L Edelman. 2024. Distinguishing the Knowable from the Unknowable with Language Models. *arXiv arXiv:2402.03563* (2024).

[6] Ekin Akyurek, Tolga Bolukbasi, Frederick Liu, Binbin Xiong, Ian Tenney, Jacob Andreas, and Kelvin Guu. 2022. Towards Tracing Knowledge in Language Models Back to the Training Data. In *Findings of EMNLP.* 2429–2446.

[7] Tariq Alhindi, Savvas Petridis, and Smaranda Muresan. 2018. Where is your evidence: Improving fact-checking by justification modeling. In *Proceedings of the first workshop on fact extraction and verification (FEVER).* 85–90.

[8] Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Mérouane Debbah, Étienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malartic, et al. 2023. The falcon series of open language models. *arXiv arXiv:2311.16867* (2023).

[9] Diego Matteo Antognini and Boi Faltings. 2021. Rationalization through Concepts. *Findings of ACL-IJCNLP 2021* (2021), 761–775.

[10] Andy Arditi, Oscar Obeso, Aaquib Syed, Daniel Paleka, Nina Panickssery, Wes Gurnee, and Neel Nanda. 2025. Refusal in language models is mediated by a single direction. *Advances in Neural Information Processing Systems* 37 (2025), 136037–136083.

[11] Walter Edwin Arnoldi. 1951. The principle of minimized iterations in the solution of the matrix eigenvalue problem. *Quarterly of applied mathematics* 9, 1 (1951), 17–29.

[12] Sanjeev Arora, Yuanzhi Li, Yingyu Liang, Tengyu Ma, and Andrej Risteski. 2018. Linear algebraic structure of word senses, with applications to polysemy. *Transactions of the Association for Computational Linguistics* 6 (2018), 483–495.

[13] Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2023. Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection. In *The Twelfth International Conference on Learning Representations.*

[14] Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2023. Self-rag: Learning to retrieve, generate, and critique through self-reflection. *arXiv arXiv:2310.11511* (2023).

[15] Amos Azaria and Tom Mitchell. 2023. The internal state of an llm knows when its lying. *arXiv arXiv:2304.13734* (2023).

[16] Andrew Bai, Chih-Kuan Yeh, Pradeep Ravikumar, Neil YC Lin, and Cho-Jui Hsieh. 2022. Concept gradient: Concept-based interpretation without linear assumption. *arXiv preprint arXiv:2208.14966* (2022).

[17] Yang Bai, Ge Pei, Jindong Gu, Yong Yang, and Xingjun Ma. 2024. Special characters attack: Toward scalable training data extraction from large language models. *arXiv preprint arXiv:2405.05990* (2024).

[18] Randall Balestriero, Romain Cosentino, and Sarath Shekkizhar. 2023. Characterizing large language model geometry solves toxicity detection and generation. *arXiv arXiv:2312.01648* (2023).

[19] David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. 2017. Network dissection: Quantifying interpretability of deep visual representations. In *CVPR.* 6541–6549.

[20] David Bau, Jun-Yan Zhu, Hendrik Strobelt, Bolei Zhou, Joshua B Tenenbaum, William T Freeman, and Antonio Torralba. 2018. GAN Dissection: Visualizing and Understanding Generative Adversarial Networks. In *ICLR.*

[21] Mohammad Belal, James She, and Simon Wong. 2023. Leveraging chatgpt as text annotation tool for sentiment analysis. *arXiv* (2023).

[22] Yonatan Belinkov, Nadir Durrani, Fahim Dalvi, Hassan Sajjad, and James Glass. 2017. What do neural machine translation models learn about morphology? *arXiv preprint arXiv:1704.03471* (2017).

[23] Nora Belrose, Zach Furman, Logan Smith, Danny Halawi, Igor Ostrovsky, Lev McKinney, Stella Biderman, and Jacob Steinhardt. 2023. Eliciting latent predictions from transformers with the tuned lens. *arXiv preprint arXiv:2303.08112* (2023).

[24] Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Michal Podstawski, Hubert Niewiadomski, Piotr Nyczyk, et al. 2023. Graph of thoughts: Solving elaborate problems with large language models. *arXiv* (2023).

[25] Amrita Bhattacharjee, Shaona Ghosh, Traian Rebedea, and Christopher Parisien. 2024. Towards Inference-time Category-wise Safety Steering for Large Language Models. In *Neurips Safe Generative AI Workshop 2024.*

[26] Amrita Bhattacharjee, Raha Moraffah, Joshua Garland, and Huan Liu. 2023. LLMs as Counterfactual Explanation Modules: Can ChatGPT Explain Black-box Text Classifiers? *arXiv arXiv:2309.13340* (2023).

[27] Steven Bills, Nick Cammarata, Dan Mossing, Henk Tillman, Leo Gao, Gabriel Goh, Ilya Sutskever, Jan Leike, Jeff Wu, and William Saunders. 2023. Language models can explain neurons in language models. *URL https://openaipublic. blob. core. windows. net/neuron-explainer/paper/index. html.(Date accessed: 14.05. 2023)* (2023).

[28] Bernd Bohnet, Vinh Q Tran, Pat Verga, Roee Aharoni, Daniel Andor, Livio Baldini Soares, Massimiliano Ciaramita, Jacob Eisenstein, Kuzman Ganchev, Jonathan Herzig, et al. 2022. Attributed question answering: Evaluation and modeling for attributed large language models. *arXiv preprint arXiv:2212.08037* (2022).

[29] Trenton Bricken, Jonathan Marcus, Siddharth Mishra-Sharma, Meg Tong, Ethan Perez, Mrinank Sharma, Kelley Rivoire, and Thomas Henighan. 2024. *Using Dictionary Learning Features as Classifiers.* https://transformer-circuits.pub/2024/features-as-classifiers/index.html

[30] Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec, Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter, Tom Henighan, and Christopher Olah. 2023. Towards Monosemanticity: Decomposing Language Models With Dictionary Learning. *Transformer Circuits Thread* (2023). https://transformer-circuits.pub/2023/monosemantic-features/index.html.

[31] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *NeurIPS* 33 (2020), 1877–1901.

[32] Collin Burns, Haotian Ye, Dan Klein, and Jacob Steinhardt. 2022. Discovering latent knowledge in language models without supervision. *arXiv preprint arXiv:2212.03827* (2022).

[33] Bodhisattwa Prasad Majumder1 Oana-Maria Camburu, Thomas Lukasiewicz, and Julian McAuley. 2021. Rationale-inspired natural language explanations with commonsense. *arXiv arXiv:2106.13876* (2021).

[34] Oana-Maria Camburu, Tim Rocktäschel, Thomas Lukasiewicz, and Phil Blunsom. 2018. e-snli: Natural language inference with natural language explanations. *NeurIPS* 31 (2018).

[35] Nick Cammarata, Shan Carter, Gabriel Goh, Chris Olah, Michael Petrov, Ludwig Schubert, Chelsea Voss, Ben Egan, and Swee Kiat Lim. 2020. Thread: circuits. *Distill* 5, 3 (2020), e24.

[36] James Campbell, Richard Ren, and Phillip Guo. 2023. Localizing Lying in Llama: Understanding Instructed Dishonesty on True-False Questions Through Prompting, Probing, and Patching. *arXiv arXiv:2311.15131* (2023).

[37] Samuel Carton, Surya Kanoria, and Chenhao Tan. 2021. What to learn, and how: Toward effective learning from rationales. *arXiv* (2021).

[38] Boli Chen, Yao Fu, Guangwei Xu, Pengjun Xie, Chuanqi Tan, Mosha Chen, and Liping Jing. 2021. Probing BERT in Hyperbolic Spaces. In *International Conference on Learning Representations.*

[39] Chao Chen, Kai Liu, Ze Chen, Yi Gu, Yue Wu, Mingyuan Tao, Zhihang Fu, and Jieping Ye. 2023. INSIDE: LLMs' Internal States Retain the Power of Hallucination Detection. In *The Twelfth ICLR.*

[40] Haozhe Chen, Carl Vondrick, and Chengzhi Mao. 2024. SelfIE: Self-Interpretation of Large Language Model Embeddings. In *Forty-first International Conference on Machine Learning*.

[41] Haozhe Chen, Carl Vondrick, and Chengzhi Mao. 2024. SelfIE: Self-Interpretation of Large Language Model Embeddings. *arXiv preprint arXiv:2403.10949* (2024).

[42] Qiguang Chen, Libo Qin, Jiaqi Wang, Jingxuan Zhou, and Wanxiang Che. 2024. Unlocking the capabilities of thought: A reasoning boundary framework to quantify and optimize chain-of-thought. *Advances in Neural Information Processing Systems* 37 (2024), 54872–54904.

[43] Shiqi Chen, Miao Xiong, Junteng Liu, Zhengxuan Wu, Teng Xiao, Siyang Gao, and Junxian He. [n. d.]. In-Context Sharpness as Alerts: An Inner Representation Perspective for Hallucination Mitigation. In *Forty-first International Conference on Machine Learning*.

[44] Wei-Lin Chen, An-Zi Yen, Hen-Hsen Huang, Cheng-Kuang Wu, and Hsin-Hsi Chen. 2023. ZARA: Improving Few-Shot Self-Rationalization for Small Language Models. *arXiv arXiv:2305.07355* (2023).

[45] Hao Cheng, Shuo Wang, Wensheng Lu, Wei Zhang, Mingyang Zhou, Kezhong Lu, and Hao Liao. 2023. Explainable Recommendation with Personalized Review Retrieval and Aspect Learning. *arXiv arXiv:2306.12657* (2023).

[46] Xiaoxue Cheng, Junyi Li, Wayne Xin Zhao, and Ji-Rong Wen. 2024. ChainLM: Empowering Large Language Models with Improved Chain-of-Thought Prompting. In *LREC/COLING*.

[47] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. *See http://vicuna.lmsys.org (accessed 14 April 2023)* (2023).

[48] Arthur Conmy, Augustine Mavor-Parker, Aengus Lynch, Stefan Heimersheim, and Adrià Garriga-Alonso. 2023. Towards automated circuit discovery for mechanistic interpretability. *Advances in Neural Information Processing Systems* 36 (2023), 16318–16352.

[49] Zeyu Cui, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. 2022. M6-rec: Generative pretrained language models are open-ended recommender systems. *arXiv arXiv:2205.08084* (2022).

[50] Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. 2023. Sparse autoencoders find highly interpretable features in language models. *arXiv preprint arXiv:2309.08600* (2023).

[51] Gautier Dagan, Frank Keller, and Alex Lascarides. 2023. Dynamic planning with a LLM. *arXiv arXiv:2308.06391* (2023).

[52] Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2022. Knowledge Neurons in Pretrained Transformers. In *ACL*.

[53] Fahim Dalvi, Hassan Sajjad, Nadir Durrani, and Yonatan Belinkov. 2020. Analyzing Redundancy in Pretrained Transformer Models. In *EMNLP*.

[54] Marina Danilevsky, Kun Qian, Ranit Aharonov, Yannis Katsis, Ban Kawas, and Prithviraj Sen. 2020. A Survey of the State of Explainable AI for Natural Language Processing. In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of ACL and the 10th IJCNLP*. 447–459.

[55] Guy Dar, Mor Geva, Ankit Gupta, and Jonathan Berant. 2023. Analyzing Transformers in Embedding Space. In *ACL*.

[56] Emily L Denton et al. 2017. Unsupervised learning of disentangled representations from video. *NeurIPS* 30 (2017).

[57] Erik Derner, Kristina Batistič, Jan Zahálka, and Robert Babuška. 2023. A security risk taxonomy for large language models. *arXiv* (2023).

[58] Jay DeYoung, Sarthak Jain, Nazneen Fatema Rajani, Eric Lehman, Caiming Xiong, Richard Socher, and Byron C Wallace. 2020. ERASER: A Benchmark to Evaluate Rationalized NLP Models. In *ACL*. 4443–4458.

[59] Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu, Roberta Raileanu, Xian Li, Asli Celikyilmaz, and Jason Weston. 2023. Chain-of-verification reduces hallucination in large language models. *arXiv arXiv:2309.11495* (2023).

[60] Alexander Yom Din, Taelin Karidi, Leshem Choshen, and Mor Geva. 2024. Jump to Conclusions: Short-Cutting Transformers with Linear Transformations. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*. 9615–9625.

[61] B. Ding, C. Qin, L. Liu, Y. K. Chia, S. Joty, B. Li, and L. Bing. 2022. Is gpt-3 a good data annotator? *arXiv* (2022).

[62] Finale Doshi-Velez and Been Kim. 2017. Towards a rigorous science of interpretable machine learning. *arXiv arXiv:1702.08608* (2017).

[63] Filip Karlo Došilović, Mario Brčić, and Nikica Hlupić. 2018. Explainable artificial intelligence: A survey. In *2018 41st International convention on information and communication technology, electronics and microelectronics (MIPRO)*. IEEE, 0210–0215.

[64] Mengnan Du, Ninghao Liu, and Xia Hu. 2019. Techniques for interpretable machine learning. *Commun. ACM* 63, 1 (2019), 68–77.

[65] Mengnan Du, Ninghao Liu, Fan Yang, and Xia Hu. 2019. Learning credible deep neural networks with rationale regularization. In *ICDM*.

[66] Mengnan Du, Varun Manjunatha, Rajiv Jain, Ruchi Deshpande, Franck Dernoncourt, Jiuxiang Gu, Tong Sun, and Xia Hu. 2021. Towards interpreting and mitigating shortcut learning behavior of NLU models. *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies* (2021).

[67] H. Duan, Yi Yang, and K. Y. Tam. 2024. Do LLMs Know about Hallucination? An Empirical Investigation of LLM's Hidden States. *arXiv* (2024).

[68] Clément Dumas, Chris Wendler, Veniamin Veselovsky, Giovanni Monea, and Robert West. 2024. Separating Tongue from Thought: Activation Patching Reveals Language-Agnostic Concept Representations in Transformers. *arXiv preprint arXiv:2411.08745* (2024).

[69] Esin Durmus, Alex Tamkin, Jack Clark, Jerry Wei, Jonathan Marcus, Joshua Batson, Kunal Handa, Liane Lovitt, Meg Tong, Miles McCain, Oliver Rausch, Saffron Huang, Sam Bowman, Stuart Ritchie, Tom Henighan, and Deep Ganguli. 2024. *Evaluating Feature Steering: A Case Study in Mitigating Social Biases*. https://anthropic.com/research/evaluating-feature-steering

[70] Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. 2018. HotFlip: White-Box Adversarial Examples for Text Classification. In *Proceedings of the 56th ACL (Volume 2: Short Papers)*. 31–36.

[71] Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, et al. 2021. A mathematical framework for transformer circuits. *Transformer Circuits Thread* 1 (2021).

[72] Siqi Fan, Xin Jiang, Xiang Li, Xuying Meng, Peng Han, Shuo Shang, Aixin Sun, Yequan Wang, and Zhongyuan Wang. 2024. Not All Layers of LLMs Are Necessary During Inference. arXiv:2403.02181 [cs.CL] https://arxiv.org/abs/2403.02181

[73] Thomas Fel, Victor Boutin, Louis Béthune, Rémi Cadène, Mazda Moayeri, Léo Andéol, Mathieu Chalvidal, and Thomas Serre. 2023. A holistic approach to unifying automatic concept extraction and concept importance estimation. *Advances in Neural Information Processing Systems* 36 (2023), 54805–54818.

[74] Shangbin Feng, Weijia Shi, Yuyang Bai, Vidhisha Balachandran, Tianxing He, and Yulia Tsvetkov. 2023. Knowledge Card: Filling LLMs' Knowledge Gaps with Plug-in Specialized Language Models. In *ICLR*.

[75] Javier Ferrando, Gerard I Gállego, Ioannis Tsiamas, and Marta R Costa-jussà. 2023. Explaining How Transformers Use Context to Build Predictions. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 5486–5513.

[76] Javier Ferrando, Gabriele Sarti, Arianna Bisazza, and Marta R Costa-jussà. 2024. A Primer on the Inner Workings of Transformer-based Language Models. *arXiv preprint arXiv:2405.00208* (2024).

[77] Yi-Fu Fu, Yu-Chieh Tu, Tzu-Ling Cheng, Cheng-Yu Lin, Yi-Ting Yang, Heng-Yi Liu, Keng-Te Liao, Da-Cheng Juan, and Shou-De Lin. 2024. Think or Remember? Detecting and Directing LLMs Towards Memorization or Generalization. *arXiv preprint arXiv:2412.18497* (2024).

[78] Isabel O Gallegos, Ryan A Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K Ahmed. 2023. Bias and fairness in large language models: A survey. *arXiv arXiv:2309.00770* (2023).

[79] Luyu Gao, Zhuyun Dai, Panupong Pasupat, Anthony Chen, Arun Tejasvi Chaganty, Yicheng Fan, Vincent Y Zhao, Ni Lao, Hongrae Lee, Da-Cheng Juan, et al. 2023. RARR: Researching and Revising What Language Models Say, Using Language Models. In *The 61st Annual Meeting Of The Association For Computational Linguistics*.

[80] Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever, Jan Leike, and Jeffrey Wu. 2024. Scaling and evaluating sparse autoencoders. *arXiv preprint arXiv:2406.04093* (2024).

[81] Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen. 2023. Enabling Large Language Models to Generate Text with Citations. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. 6465–6488.

[82] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. *arXiv arXiv:2312.10997* (2023).

[83] Yuan Ge, Yilun Liu, Chi Hu, Weibin Meng, Shimin Tao, Xiaofeng Zhao, Hongxia Ma, Li Zhang, Boxing Chen, Hao Yang, Bei Li, Tong Xiao, and Jingbo Zhu. 2024. Clustering and Ranking: Diversity-preserved Instruction Selection through Expert-aligned Quality Estimation. arXiv:2402.18191 [cs.CL] https://arxiv.org/abs/2402.18191

[84] Atticus Geiger, Hanson Lu, Thomas Icard, and Christopher Potts. 2021. Causal abstractions of neural networks. *Advances in Neural Information Processing Systems* 34 (2021), 9574–9586.

[85] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. 2020. Shortcut learning in deep neural networks. *Nature Machine Intelligence* 2, 11 (2020), 665–673.

[86] Shijie Geng, Shuchang Liu, Zuohui Fu, Yingqiang Ge, and Yongfeng Zhang. 2022. Recommendation as language processing (rlp): A unified pretrain, personalized prompt & predict paradigm (p5). In *RecSys*. 299–315.

[87] Mor Geva, Avi Caciularu, Kevin Wang, and Yoav Goldberg. 2022. Transformer Feed-Forward Layers Build Predictions by Promoting Concepts in the Vocabulary Space. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. 30–45.

[88] Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2021. Transformer Feed-Forward Layers Are Key-Value Memories. In *EMNLP*.

[89] Amirata Ghorbani, James Wexler, James Y Zou, and Been Kim. 2019. Towards automatic concept-based explanations. *NeurIPS* 32 (2019).

[90] Fabrizio Gilardi, Meysam Alizadeh, and Maël Kubli. 2023. Chatgpt outperforms crowd-workers for text-annotation tasks. *arXiv* (2023).

[91] Andrey Gromov, Kushal Tirumala, Hassan Shapourian, Paolo Glorioso, and Daniel A. Roberts. 2025. The Unreasonable Ineffectiveness of the Deeper Layers. arXiv:2403.17887 [cs.CL] https://arxiv.org/abs/2403.17887

[92] Roger Grosse, Juhan Bae, Cem Anil, Nelson Elhage, Alex Tamkin, Amirhossein Tajdini, Benoit Steiner, Dustin Li, Esin Durmus, Ethan Perez, et al. 2023. Studying large language model generalization with influence functions. *arXiv arXiv:2308.03296* (2023).

[93] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948* (2025).

[94] Han Guo, Nazneen Rajani, Peter Hase, Mohit Bansal, and Caiming Xiong. 2021. FastIF: Scalable Influence Functions for Efficient Model Interpretation and Debugging. In *EMNLP*. 10333–10350.

[95] Qing Guo, Prashan Wanigasekara, Skyler Zheng, Jacob Zhiyuan Fang, Xinwei Deng, and Chenyang Tao. 2023. How do multimodal LLMs really fare in classical vision few-shot challenges? A deep dive. (2023).

[96] Abhijeet Gupta, Gemma Boleda, Marco Baroni, and Sebastian Padó. 2015. Distributional vectors encode referential attributes. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*. 12–21.

[97] Akshat Gupta, Dev Sajnani, and Gopala Anumanchipalli. 2024. A Unified Framework for Model Editing. arXiv:2403.14236 [cs.LG] https://arxiv.org/abs/2403.14236

[98] Yoav Gur-Arieh, Roy Mayan, Chen Agassy, Atticus Geiger, and Mor Geva. 2025. Enhancing Automated Interpretability with Output-Centric Feature Descriptions. *arXiv preprint arXiv:2501.08319* (2025).

[99] Sai Gurrapu, Ajay Kulkarni, Lifu Huang, Ismini Lourentzou, and Feras A Batarseh. 2023. Rationalization for explainable nlp: A survey. *Frontiers in Artificial Intelligence* 6 (2023).

[100] Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training. In *ICML*.

[101] Xiaoqi Han, Ru Li, Hongye Tan, Wang Yuanlong, Qinghua Chai, and Jeff Pan. 2023. Improving sequential model editing with fact retrieval. In *Findings of the Association for Computational Linguistics: EMNLP 2023*. 11209–11224.

[102] Xiaochuang Han and Yulia Tsvetkov. 2022. Orca: Interpreting prompted language models via locating supporting data evidence in the ocean of pretraining data. *arXiv arXiv:2205.12600* (2022).

[103] Xiaochuang Han, Byron C Wallace, and Yulia Tsvetkov. 2020. Explaining Black Box Predictions and Unveiling Data Artifacts through Influence Functions. In *ACL*. 5553–5563.

[104] Michael Hanna, Ollie Liu, and Alexandre Variengien. 2023. How does GPT-2 compute greater-than?: Interpreting mathematical abilities in a pre-trained language model. *Advances in Neural Information Processing Systems* 36 (2023), 76033–76060.

[105] Michael Hanna, Sandro Pezzelle, and Yonatan Belinkov. [n. d.]. Have Faith in Faithfulness: Going Beyond Circuit Overlap When Finding Model Mechanisms. In *ICML 2024 Workshop on Mechanistic Interpretability*.

[106] Jacqueline Harding, William D'Alessandro, NG Laskowski, and Robert Long. 2023. AI language models cannot replace human research participants. *Ai & Society* (2023), 1–3.

[107] Peter Hase, Mohit Bansal, Been Kim, and Asma Ghandeharioun. 2024. Does localization inform editing? surprising differences in causality-based localization vs. knowledge editing in language models. *NeurIPS* 36 (2024).

[108] Jinghan He, Kuan Zhu, Haiyun Guo, Junfeng Fang, Zhenglin Hua, Yuheng Jia, Ming Tang, Tat-Seng Chua, and Jinqiao Wang. 2024. Cracking the Code of Hallucination in LVLMs with Vision-aware Head Divergence. *arXiv preprint arXiv:2412.13949* (2024).

[109] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *CVPR*. 770–778.

[110] Ruining He, Wang-Cheng Kang, and Julian McAuley. 2017. Translation-based recommendation. In *RecSys*. 161–169.

[111] Zexue He, Marco Tulio Ribeiro, and Fereshte Khani. 2023. Targeted Data Generation: Finding and Fixing Model Weaknesses. *arXiv* (2023).

[112] Zirui He, Haiyan Zhao, Yiran Qiao, Fan Yang, Ali Payani, Jing Ma, and Mengnan Du. 2025. SAIF: A Sparse Autoencoder Framework for Interpreting and Steering Instruction Following of Language Models. *arXiv preprint arXiv:2502.11356* (2025).

[113] Thomas Heap, Tim Lawson, Lucy Farnik, and Laurence Aitchison. 2025. Sparse Autoencoders Can Interpret Randomly Initialized Transformers. *arXiv preprint arXiv:2501.17727* (2025).

[114] Marti A. Hearst, Susan T Dumais, Edgar Osuna, John Platt, and Bernhard Scholkopf. 1998. Support vector machines. *IEEE Intelligent Systems and their applications* 13, 4 (1998), 18–28.

[115] Juyeon Heo, Christina Heinze-Deml, Oussama Elachqar, Shirley Ren, Udhay Nallasamy, Andy Miller, Kwan Ho Ryan Chan, and Jaya Narain. 2024. Do LLMs" know" internally when they follow instructions? *arXiv preprint arXiv:2410.14516* (2024).

[116] John Hewitt and Percy Liang. 2019. Designing and Interpreting Probes with Control Tasks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. 2733–2743.

[117] John Hewitt and Christopher D Manning. 2019. A structural probe for finding syntax in word representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. 4129–4138.

[118] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. 2016. beta-vae: Learning basic visual concepts with a constrained variational framework. In *ICLR*.

[119] Namgyu Ho, Laura Schmid, and Se-Young Yun. 2022. Large language models are reasoning teachers. *arXiv arXiv:2212.10071* (2022).

[120] Bairu Hou, Yang Zhang, Jiabao Ji, Yujian Liu, Kaizhi Qian, Jacob Andreas, and Shiyu Chang. 2025. Thinkprune: Pruning long chain-of-thought of llms via reinforcement learning. *arXiv preprint arXiv:2504.01296* (2025).

[121] Cheng-Yu Hsieh, Chun-Liang Li, Chih-kuan Yeh, Hootan Nakhost, Yasuhisa Fujii, Alex Ratner, Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. 2023. Distilling Step-by-Step! Outperforming Larger Language Models with Less Training Data and Smaller Model Sizes. In *Findings of the Association for Computational Linguistics: ACL 2023*. Association for Computational Linguistics, 8003–8017. https://doi.org/10.18653/v1/2023.findings-acl.507

[122] Yebowen Hu, Kaiqiang Song, Sangwoo Cho, Xiaoyang Wang, Hassan Foroosh, and Fei Liu. 2023. Decipherpref: Analyzing influential factors in human preference judgments via gpt-4. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. 8344–8357.

[123] Fan Huang, Haewoon Kwak, and Jisun An. 2023. Chain of explanation: New prompting method to generate quality natural language explanation for implicit hate speech. In *Companion Proceedings of the WWW 2023*. 90–93.

[124] Fan Huang, Haewoon Kwak, and Jisun An. 2023. Is ChatGPT better than Human Annotators? Potential and Limitations of ChatGPT in Explaining Implicit Hate Speech. In *Companion Proceedings of the WWW 2023*. 294–297.

[125] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. 2023. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *arXiv* (2023).

[126] Ruixuan Huang. 2024. Steering LLMs' Behavior with Concept Activation Vectors. Draft manuscript. Available on LessWrong forum..

[127] Shiyuan Huang, Siddarth Mamidanna, Shreedhar Jangam, Yilun Zhou, and Leilani H Gilpin. 2023. Can Large Language Models Explain Themselves? A Study of LLM-Generated Self-Explanations. *arXiv arXiv:2310.11207* (2023).

[128] Youcheng Huang, Jingkun Tang, Duanyu Feng, Zheng Zhang, Wenqiang Lei, Jiancheng Lv, and Anthony G Cohn. 2024. Dishonesty in helpful and harmless alignment. *arXiv preprint arXiv:2406.01931* (2024).

[129] Gautier Izacard and Edouard Grave. 2020. Leveraging passage retrieval with generative models for open domain question answering. *arXiv* (2020).

[130] S. Jain, R. Kirk, E. S. Lubana, R. P Dick, H. Tanaka, E. Grefenstette, T. Rocktäschel, and D. S. Krueger. 2023. Mechanistically analyzing the effects of fine-tuning on procedurally defined tasks. *arXiv* (2023).

[131] Shashank Mohan Jain. 2022. Introduction to Transformers for NLP: With the Hugging Face Library and Models to Solve Problems. Springer.

[132] Stanisław Jastrzebski, Devansh Arpit, Nicolas Ballas, Vikas Verma, Tong Che, and Yoshua Bengio. 2018. Residual Connections Encourage Iterative Inference. In *International Conference on Learning Representations*.

[133] Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *Comput. Surveys* 55, 12 (2023), 1–38.

[134] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7B. *arXiv arXiv:2310.06825* (2023).

[135] Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. 2023. Large language models struggle to learn long-tail knowledge. In *ICML*. 15696–15707.

[136] Jikun Kang, Xin Zhe Li, Xi Chen, Amirreza Kazemi, Qianyi Sun, Boxing Chen, Dong Li, Xu He, Quan He, Feng Wen, et al. 2024. Mindstar: Enhancing math reasoning in pre-trained llms at inference time. *arXiv preprint arXiv:2405.16265* (2024).

[137] Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. 2022. Decomposed Prompting: A Modular Approach for Solving Complex Tasks. In *The Eleventh ICLR*.

[138] Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, et al. 2018. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). In *International conference on machine learning*. PMLR, 2668–2677.

[139] Jaehyung Kim, Jaehyun Nam, Sangwoo Mo, Jongjin Park, Sang-Woo Lee, Minjoon Seo, Jung-Woo Ha, and Jinwoo Shin. 2023. SuRe: Improving Open-domain Question Answering of LLMs via Summarized Retrieval. In *The Twelfth ICLR*.

[140] Seungone Kim, Se June Joo, Doyoung Kim, Joel Jang, Seonghyeon Ye, Jamin Shin, and Minjoon Seo. 2023. The CoT Collection: Improving Zero-shot and Few-shot Learning of Language Models via Chain-of-Thought Fine-Tuning. *arXiv arXiv:2305.14045* (2023).

[141] Pieter-Jan Kindermans, Sara Hooker, Julius Adebayo, Maximilian Alber, Kristof T Schütt, Sven Dähne, Dumitru Erhan, and Been Kim. [n. d.]. The (Un) reliability of Saliency Methods. In *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*. 267–280.

[142] Pang Wei Koh and Percy Liang. 2017. Understanding black-box predictions via influence functions. In *ICML*. PMLR, 1885–1894.

[143] Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Mussmann, Emma Pierson, Been Kim, and Percy Liang. 2020. Concept bottleneck models. In *ICML*. PMLR, 5338–5348.

[144] Takeshi Kojima, Shixiang (Shane) Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large Language Models are Zero-Shot Reasoners. In *Advances in Neural Information Processing Systems*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (Eds.), Vol. 35. Curran Associates, Inc., 22199–22213. https://proceedings.neurips.cc/paper_files/paper/2022/file/8bb0d291acd4acf06ef112099c16f326-Paper-Conference.pdf

[145] Enja Kokalj, Blaž Škrlj, Nada Lavrač, Senja Pollak, and Marko Robnik-Šikonja. 2021. BERT meets shapley: Extending SHAP explanations to transformer-based classifiers. In *Proceedings of the EACL Hackashop on News Media Content Analysis and Automated Report Generation*. 16–21.

[146] Satyapriya Krishna, Jiaqi Ma, Dylan Slack, Asma Ghandeharioun, Sameer Singh, and Himabindu Lakkaraju. 2024. Post hoc explanations of language models can improve language models. *Advances in Neural Information Processing Systems* 36 (2024).

[147] Tamera Lanham, Anna Chen, Ansh Radhakrishnan, Benoit Steiner, Carson Denison, Danny Hernandez, Dustin Li, Esin Durmus, Evan Hubinger, Jackson Kernion, et al. 2023. Measuring faithfulness in chain-of-thought reasoning. *arXiv arXiv:2307.13702* (2023).

[148] Ehsan Latif and Xiaoming Zhai. 2024. Fine-tuning chatgpt for automatic scoring. *Computers and Education: Artificial Intelligence* (2024), 100210.

[149] Andrew Lee, Xiaoyan Bai, Itamar Pres, Martin Wattenberg, Jonathan K Kummerfeld, and Rada Mihalcea. 2024. A mechanistic understanding of alignment algorithms: A case study on dpo and toxicity. *arXiv arXiv:2401.01967* (2024).

[150] Gyeong-Geon Lee, Lehong Shi, Ehsan Latif, Yizhu Gao, Arne Bewersdorf, Matthew Nyaaba, Shuchen Guo, Zihao Wu, Zhengliang Liu, Hui Wang, et al. 2023. Multimodality of ai for education: Towards artificial general intelligence. *arXiv arXiv:2312.06037* (2023).

[151] Yuxuan Lei, Jianxun Lian, Jing Yao, Xu Huang, Defu Lian, and Xing Xie. 2023. RecExplainer: Aligning Large Language Models for Recommendation Model Interpretability. *arXiv arXiv:2311.10947* (2023).

[152] Chak Tou Leong, Yi Cheng, Jiashuo Wang, Jian Wang, and Wenjie Li. 2023. Self-detoxifying language models via toxification reversal. *arXiv* (2023).

[153] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *NeurIPS* 33 (2020), 9459–9474.

[154] Daliang Li, Ankit Singh Rawat, Manzil Zaheer, Xin Wang, Michal Lukasik, Andreas Veit, Felix Yu, and Sanjiv Kumar. 2022. Large language models with controllable working memory. *arXiv arXiv:2211.05110* (2022).

[155] Haoran Li, Dadi Guo, Wei Fan, Mingshi Xu, and Yangqiu Song. 2023. Multi-step jailbreaking privacy attacks on chatgpt. *arXiv* (2023).

[156] Jiwei Li, Xinlei Chen, Eduard Hovy, and Dan Jurafsky. 2016. Visualizing and Understanding Neural Models in NLP. In *NAACL*.

[157] Jiacheng Li, Zhankui He, Jingbo Shang, and Julian McAuley. 2023. Ucepic: Unifying aspect planning and lexical constraints for generating explanations in recommendation. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 1248–1257.

[158] Jiwei Li, Will Monroe, and Dan Jurafsky. 2016. Understanding neural networks through representation erasure. *arXiv arXiv:1612.08220* (2016).

[159] Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2024. Inference-time intervention: Eliciting truthful answers from a language model. *NeurIPS* 36 (2024).

[160] Lei Li, Yongfeng Zhang, and Li Chen. 2020. Generate neural template explanations for recommendation. In *CIKM*. 755–764.

[161] Lei Li, Yongfeng Zhang, and Li Chen. 2023. Personalized prompt learning for explainable recommendation. *ACM Trans. on Info. Sys.* (2023).

[162] Ming Li, Yong Zhang, Zhitao Li, Jiuhai Chen, Lichang Chen, Ning Cheng, Jianzong Wang, Tianyi Zhou, and Jing Xiao. 2024. From Quantity to Quality: Boosting LLM Performance with Self-Guided Data Selection for Instruction Tuning. arXiv:2308.12032 [cs.CL] https://arxiv.org/abs/2308.12032

[163] Q. Li, Q. Tao, S. Joty, J. Cai, and J. Luo. 2018. Vqa-e: Explaining, elaborating, and enhancing your answers for visual questions. In *ECCV*.

[164] Shiyang Li, Jianshu Chen, Yelong Shen, Zhiyu Chen, Xinlu Zhang, Zekun Li, Hong Wang, Jing Qian, Baolin Peng, Yi Mao, et al. 2022. Explanations from large language models make small reasoners better. *arXiv arXiv:2210.06726* (2022).

[165] Tianlong Li, Xiaoqing Zheng, and Xuanjing Huang. 2024. Open the Pandora's Box of LLMs: Jailbreaking LLMs through Representation Engineering. *arXiv arXiv:2401.06824* (2024).

[166] Wei Li, Wenhao Wu, Moye Chen, Jiachen Liu, Xinyan Xiao, and Hua Wu. 2022. Faithfulness in natural language generation: A systematic survey of analysis, evaluation and optimization methods. *arXiv arXiv:2203.05227* (2022).

[167] Xiaopeng Li, Shasha Li, Shezheng Song, Jing Yang, Jun Ma, and Jie Yu. 2024. PMET: Precise Model Editing in a Transformer. arXiv:2308.08742 [cs.CL] https://arxiv.org/abs/2308.08742

[168] Yingji Li, Mengnan Du, Rui Song, Xin Wang, and Ying Wang. 2023. A survey on fairness in large language models. *arXiv arXiv:2308.10149* (2023).

[169] Yingji Li, Mengnan Du, Xin Wang, and Ying Wang. 2023. Prompt Tuning Pushes Farther, Contrastive Learning Pulls Closer: A Two-Stage Approach to Mitigate Social Biases. In *Proceedings of the 61st ACL (Volume 1: Long Papers)*. 14254–14267.

[170] Tom Lieberum, Matthew Rahtz, János Kramár, Geoffrey Irving, Rohin Shah, and Vladimir Mikulik. 2023. Does circuit analysis interpretability scale? evidence from multiple choice capabilities in chinchilla. *arXiv arXiv:2307.09458* (2023).

[171] Guangliang Liu, Haitao Mao, Bochuan Cao, Zhiyu Xue, Xitong Zhang, Rongrong Wang, Jiliang Tang, and Kristen Johnson. 2024. On the intrinsic self-correction capability of LLMs: Uncertainty and latent concept. *arXiv preprint arXiv:2406.02378* (2024).

[172] Ninghao Liu, Mengnan Du, Ruocheng Guo, Huan Liu, and Xia Hu. 2021. Adversarial attacks and defenses: An interpretation perspective. *ACM SIGKDD Explorations Newsletter* 23, 1 (2021), 86–99.

[173] Ninghao Liu, Hongxia Yang, and Xia Hu. 2018. Adversarial detection with model interpretation. In *SIGKDD*. 1803–1811.

[174] Nelson F Liu, Matt Gardner, Yonatan Belinkov, Matthew E Peters, and Noah A Smith. 2019. Linguistic knowledge and transferability of contextual representations. *arXiv preprint arXiv:1903.08855* (2019).

[175] Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics* 12 (2024), 157–173.

[176] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *Comput. Surveys* 55, 9 (2023), 1–35.

[177] Yi Liu, Gelei Deng, Zhengzi Xu, Yuekang Li, Yaowen Zheng, Ying Zhang, Lida Zhao, Tianwei Zhang, and Yang Liu. 2023. Jailbreaking chatgpt via prompt engineering: An empirical study. *arXiv arXiv:2305.13860* (2023).

[178] Yang Liu, Yuanshun Yao, Jean-Francois Ton, Xiaoying Zhang, Ruocheng Guo Hao Cheng, Yegor Klochkov, Muhammad Faaiz Taufiq, and Hang Li. 2023. Trustworthy LLMs: a Survey and Guideline for Evaluating Large Language Models' Alignment. *arXiv arXiv:2308.05374* (2023).

[179] Zhaoyang Liu, Zeqiang Lai, Zhangwei Gao, Erfei Cui, Zhiheng Li, Xizhou Zhu, Lewei Lu, Qifeng Chen, Yu Qiao, Jifeng Dai, et al. 2023. Controlllm: Augment language models with tools by searching on graphs. *arXiv arXiv:2310.17796* (2023).

[180] Ilya Loshchilov and Frank Hutter. 2018. Decoupled Weight Decay Regularization. In *ICLR*.

[181] Yin Lou, Rich Caruana, and Johannes Gehrke. 2012. Intelligible models for classification and regression. In *SIGKDD*. 150–158.

[182] Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. *NeurIPS* 35 (2022), 2507–2521.

[183] Scott M Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. *NeurIPS* 30 (2017).

[184] Haotian Luo, Li Shen, Haiying He, Yibo Wang, Shiwei Liu, Wei Li, Naiqiang Tan, Xiaochun Cao, and Dacheng Tao. 2025. O1-Pruner: Length-Harmonizing Fine-Tuning for O1-Like Reasoning Pruning. *arXiv preprint arXiv:2501.12570* (2025).

[185] Haoyan Luo and Lucia Specia. 2024. From Understanding to Utilization: A Survey on Explainability for Large Language Models. *arXiv* (2024).

[186] Yun Luo, Zhen Yang, Fandong Meng, Yingjie Li, Fang Guo, Qinglin Qi, Jie Zhou, and Yue Zhang. 2023. XAL: EXplainable Active Learning Makes Classifiers Better Low-resource Learners. *arXiv arXiv:2310.05502* (2023).

[187] Alisia Lupidi, Carlos Gemmell, Nicola Cancedda, Jane Dwivedi-Yu, Jason Weston, Jakob Foerster, Roberta Raileanu, and Maria Lomeli. 2024. Source2synth: Synthetic data generation and curation grounded in real data sources. *arXiv preprint arXiv:2409.08239* (2024).

[188] Qing Lyu, Shreya Havaldar, Adam Stein, Li Zhang, Delip Rao, Eric Wong, Marianna Apidianaki, and Chris Callison-Burch. 2023. Faithful chain-of-thought reasoning. *arXiv arXiv:2301.13379* (2023).

[189] Weicheng Ma, Henry Scheible, Brian Wang, Goutham Veeramachaneni, Pratim Chowdhary, Alan Sun, Andrew Koulogeorge, Lili Wang, Diyi Yang, and Soroush Vosoughi. 2023. Deciphering Stereotypes in Pre-Trained Language Models. In *EMNLP*. 11328–11345.

[190] Monte MacDiarmid, Timothy Maxwell, Nicholas Schiefer, Jesse Mu, Jared Kaplan, David Duvenaud, Sam Bowman, Alex Tamkin, Ethan Perez, Mrinank Sharma, et al. 2024. Simple probes can catch sleeper agents. *Anthropic Research Updates* (2024).

[191] L. C. Magister, J. Mallinson, J. Adamek, E. Malmi, and A. Severyn. 2022. Teaching small language models to reason. *arXiv* (2022).

[192] Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. In *Proceedings of the 61st ACL (Volume 1: Long Papers)*. 9802–9822.

[193] Samuel Marks, Can Rager, Eric J Michaud, Yonatan Belinkov, David Bau, and Aaron Mueller. 2024. Sparse Feature Circuits: Discovering and Editing Interpretable Causal Graphs in Language Models. *CoRR* (2024).

[194] Samuel Marks and Max Tegmark. 2023. The geometry of truth: Emergent linear structure in large language model representations of true/false datasets. *arXiv arXiv:2310.06824* (2023).

[195] Xin Men, Mingyu Xu, Qingyu Zhang, Bingning Wang, Hongyu Lin, Yaojie Lu, Xianpei Han, and Weipeng Chen. 2024. ShortGPT: Layers in Large Language Models are More Redundant Than You Expect. arXiv:2403.03853 [cs.CL] https://arxiv.org/abs/2403.03853

[196] Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in gpt. *Advances in neural information processing systems* 35 (2022), 17359–17372.

[197] Kevin Meng, Arnab Sen Sharma, Alex J Andonian, Yonatan Belinkov, and David Bau. 2022. Mass-Editing Memory in a Transformer. In *ICLR*.

[198] Sachit Menon and Carl Vondrick. 2022. Visual Classification via Description from Large Language Models. In *The Eleventh ICLR*.

[199] Ning Miao, Yee Whye Teh, and Tom Rainforth. 2023. Selfcheck: Using llms to zero-shot check their own step-by-step reasoning. *arXiv* (2023).

[200] Timothee Mickus, Denis Paperno, and Mathieu Constant. 2022. How to dissect a Muppet: The structure of transformer embedding spaces. *Transactions of the Association for Computational Linguistics* 10 (2022), 981–996.

[201] Beren Millidge and Sid Black. 2022. The singular value decompositions of transformer weight matrices are highly interpretable. https://www.alignmentforum.org/posts/mkbGjzxD8d8XqKHzA/the-singular-value-decompositions-of-transformer-weight

[202] Hosein Mohebbi, Ali Modarressi, and Mohammad Taher Pilehvar. 2021. Exploring the Role of BERT Token Representations to Explain Sentence Probing Results. In *EMNLP*. 792–806.

[203] Hosein Mohebbi, Willem Zuidema, Grzegorz Chrupała, and Afra Alishahi. 2023. Quantifying Context Mixing in Transformers. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*. 3378–3400.

[204] Grégoire Montavon, Alexander Binder, Sebastian Lapuschkin, Wojciech Samek, and Klaus-Robert Müller. 2019. Layer-wise relevance propagation: an overview. *Explainable AI: interpreting, explaining and visualizing deep learning* (2019), 193–209.

[205] Grégoire Montavon, Sebastian Lapuschkin, Alexander Binder, Wojciech Samek, and Klaus-Robert Müller. 2017. Explaining nonlinear classification decisions with deep taylor decomposition. *Pattern recognition* 65 (2017), 211–222.

[206] W James Murdoch, Chandan Singh, Karl Kumbier, Reza Abbasi-Asl, and Bin Yu. 2019. Definitions, methods, and applications in interpretable machine learning. *Proceedings of the National Academy of Sciences* 116, 44 (2019), 22071–22080.

[207] Milad Nasr, Nicholas Carlini, Jonathan Hayase, Matthew Jagielski, A Feder Cooper, Daphne Ippolito, Christopher A Choquette-Choo, Eric Wallace, Florian Tramèr, and Katherine Lee. 2023. Scalable extraction of training data from (production) language models. *arXiv arXiv:2311.17035* (2023).

[208] Angus Nicolson, Yarin Gal, and J Alison Noble. 2024. TextCAVs: Debugging vision models using text. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 99–109.

[209] nostalgebraist. 2020. Interpreting GPT: the Logit Lens. *LESSWRONG* (2020).

[210] Bruno A Olshausen and David J Field. 1997. Sparse coding with an overcomplete basis set: A strategy employed by V1? *Vision research* 37, 23 (1997), 3311–3325.

[211] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *NeurIPS* 35 (2022), 27730–27744.

[212] Alexander Pan, Lijie Chen, and Jacob Steinhardt. 2024. LatentQA: Teaching LLMs to Decode Activations Into Natural Language. *arXiv preprint arXiv:2412.08686* (2024).

[213] Richard Yuanzhe Pang, Weizhe Yuan, He He, Kyunghyun Cho, Sainbayar Sukhbaatar, and Jason Weston. 2024. Iterative reasoning preference optimization. *Advances in Neural Information Processing Systems* 37 (2024), 116617–116637.

[214] Seongheon Park, Xuefeng Du, Min-Hsuan Yeh, Haobo Wang, and Yixuan Li. 2025. How to Steer LLM Latents for Hallucination Detection? *arXiv preprint arXiv:2503.01917* (2025).

[215] Gonçalo Paulo and Nora Belrose. 2025. Sparse Autoencoders Trained on the Same Data Learn Different Features. *arXiv preprint arXiv:2501.16615* (2025).

[216] Fabio Petroni, Patrick Lewis, Aleksandra Piktus, Tim Rocktäschel, Yuxiang Wu, Alexander H Miller, and Sebastian Riedel. 2020. How context affects language models' factual predictions. *arXiv arXiv:2005.04611* (2020).

[217] Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language Models as Knowledge Bases?. In *EMNLP-IJCNLP*. 2463–2473.

[218] Tiago Pimentel, Josef Valvoda, Rowan Hall Maudslay, Ran Zmigrod, Adina Williams, and Ryan Cotterell. 2020. Information-Theoretic Probing for Linguistic Structure. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 4609–4622.

[219] G. Pruthi, F. Liu, S. Kale, and M. Sundararajan. 2020. Estimating training data influence by tracing gradient descent. *NeurIPS* (2020).

[220] Jirui Qi, Gabriele Sarti, Raquel Fernández, and Arianna Bisazza. 2024. Model Internals-based Answer Attribution for Trustworthy Retrieval-Augmented Generation. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 6037–6053.

[221] Chen Qian, Dongrui Liu, Jie Zhang, Yong Liu, and Jing Shao. 2024. Dean: Deactivating the coupled neurons to mitigate fairness-privacy conflicts in large language models. *arXiv preprint arXiv:2410.16672* (2024).

[222] Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiaao Chen, Michihiro Yasunaga, and Diyi Yang. 2023. Is ChatGPT a general-purpose natural language processing task solver? *arXiv arXiv:2302.06476* (2023).

[223] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *ICML*. PMLR, 8748–8763.

[224] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog* 1, 8 (2019), 9.

[225] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research* 21, 1 (2020), 5485–5551.

[226] Nazneen Fatema Rajani, Bryan McCann, Caiming Xiong, and Richard Socher. 2019. Explain Yourself! Leveraging Language Models for Commonsense Reasoning. In *Proceedings of the 57th ACL*. 4932–4942.

[227] Aishik Rakshit, Smriti Singh, Shuvam Keshari, Arijit Ghosh Chowdhury, Vinija Jain, and Aman Chadha. 2024. From Prejudice to Parity: A New Approach to Debiasing Large Language Model Word Embeddings. *arXiv arXiv:2402.11512* (2024).

[228] Ruiyang Ren, Yuhao Wang, Yingqi Qu, Wayne Xin Zhao, Jing Liu, Hao Tian, Hua Wu, Ji-Rong Wen, and Haifeng Wang. 2023. Investigating the factual knowledge boundary of large language models with retrieval augmentation. *arXiv arXiv:2307.11019* (2023).

[229] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. " Why should i trust you?" Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. 1135–1144.

[230] Nina Rimsky, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Turner. 2024. Steering Llama 2 via Contrastive Activation Addition. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 15504–15522.

[231] Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: BM25 and beyond. *Foundations and Trends® in Information Retrieval* 3, 4 (2009), 333–389.

[232] Cynthia Rudin. 2019. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature machine intelligence* 1, 5 (2019), 206–215.

[233] Cynthia Rudin, Chaofan Chen, Zhi Chen, Haiyang Huang, Lesia Semenova, and Chudi Zhong. 2022. Interpretable machine learning: Fundamental principles and 10 grand challenges. *Statistic Surveys* 16 (2022), 1–85.

[234] Rachneet Sachdeva, Martin Tutek, and Iryna Gurevych. 2023. CATfOOD: Counterfactual Augmented Training for Improving Out-of-Domain Performance and Calibration. *arXiv arXiv:2309.07822* (2023).

[235] Swarnadeep Saha, Xian Li, Marjan Ghazvininejad, Jason Weston, and Tianlu Wang. 2025. Learning to Plan & Reason for Evaluation with Thinking-LLM-as-a-Judge. *arXiv preprint arXiv:2501.18099* (2025).

[236] Gabriele Sarti, Nils Feldhus, Ludwig Sickert, and Oskar Van Der Wal. 2023. Inseq: An Interpretability Toolkit for Sequence Generation Models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*. 421–435.

[237] A. Scherlis, K. Sachan, A. S Jermyn, J. Benton, and B. Shlegeris. 2022. Polysemanticity and capacity in neural networks. *arXiv* (2022).

[238] Andrea Schioppa, Polina Zablotskaia, David Vilar, and Artem Sokolov. 2022. Scaling up influence functions. In *AAAI*, Vol. 36. 8179–8186.

[239] Sarah Schwettmann, Tamar Shaham, Joanna Materzynska, Neil Chowdhury, Shuang Li, Jacob Andreas, David Bau, and Antonio Torralba. 2024. FIND: A Function Description Benchmark for Evaluating Interpretability Methods. *NeurIPS* 36 (2024).

[240] Marco Scialanga, Thibault Laugel, Vincent Grari, and Marcin Detyniecki. 2025. SAKE: Steering Activations for Knowledge Editing. *arXiv preprint arXiv:2503.01751* (2025).

[241] Ramprasaath R Selvaraju, Abhishek Das, Ramakrishna Vedantam, Michael Cogswell, Devi Parikh, and Dhruv Batra. 2016. Grad-CAM: Why did you say that? *arXiv arXiv:1611.07450* (2016).

[242] Chenming Shang, Shiji Zhou, Yujiu Yang, Hengyuan Zhang, Xinzhe Ni, and Yuwang Wang. 2024. Incremental Residual Concept Bottleneck Models. *arXiv preprint arXiv:2404.08978* (2024).

[243] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. 2024. Hugginggpt: Solving ai tasks with chatgpt and its friends in hugging face. *NeurIPS* 36 (2024).

[244] Yucheng Shi, Shaochen Xu, Zhengliang Liu, Tianming Liu, Xiang Li, and Ninghao Liu. 2023. MedEdit: Model Editing for Medical Question Answering with External Knowledge Bases. *arXiv arXiv:2309.16035* (2023).

[245] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. 2017. Learning important features through propagating activation differences. In *ICML*. PMLR, 3145–3153.

[246] Dong Shu, Xuansheng Wu, Haiyan Zhao, Daking Rai, Ziyu Yao, Ninghao Liu, and Mengnan Du. 2025. A Survey on Sparse Autoencoders: Interpreting the Internal Mechanisms of Large Language Models. arXiv:2503.05613 [cs.LG] https://arxiv.org/abs/2503.05613

[247] K. Shuster, S. Poff, M. Chen, D. Kiela, and J. Weston. 2021. Retrieval augmentation reduces hallucination in conversation. *arXiv* (2021).

[248] Chandan Singh, Armin Askari, Rich Caruana, and Jianfeng Gao. 2023. Augmenting interpretable models with large language models during training. *Nature Communications* 14, 1 (2023), 7913.

[249] Shashwat Singh, Shauli Ravfogel, Jonathan Herzig, Roee Aharoni, Ryan Cotterell, and Ponnurangam Kumaraguru. 2024. Mimic: Minimally modified counterfactuals in the representation space. *arXiv e-prints* (2024), arXiv–2402.

[250] Shamane Siriwardhana, Rivindu Weerasekera, Elliott Wen, Tharindu Kaluarachchi, Rajib Rana, and Suranga Nanayakkara. 2023. Improving the domain adaptation of retrieval augmented generation (RAG) models for open domain question answering. *TACL* 11 (2023), 1–17.

[251] Yan-Yan Song and LU Ying. 2015. Decision tree methods: applications for classification and prediction. *Shanghai archives of psychiatry* (2015).

[252] Alessandro Stolfo, Vidhisha Balachandran, Safoora Yousefi, Eric Horvitz, and Besmira Nushi. 2024. Improving instruction-following in language models through activation steering. *arXiv preprint arXiv:2410.12877* (2024).

[253] Zhiqing Sun, Xuezhi Wang, Yi Tay, Yiming Yang, and Denny Zhou. 2023. Recitation-Augmented Language Models. In *The Eleventh International Conference on Learning Representations*.

[254] Manan Suri, Nishit Anand, and Amisha Bhaskar. 2025. Mitigating Memorization in LLMs using Activation Steering. *arXiv preprint arXiv:2503.06040* (2025).

[255] Yiming Tan, Dehai Min, Yu Li, Wenbo Li, Nan Hu, Yongrui Chen, and Guilin Qi. 2023. Can ChatGPT replace traditional KBQA models? An in-depth analysis of the question answering performance of the GPT LLM family. In *International Semantic Web Conference*. Springer, 348–367.

[256] Zhen Tan, Alimohammad Beigi, Song Wang, Ruocheng Guo, Amrita Bhattacharjee, Bohan Jiang, Mansooreh Karami, Jundong Li, Lu Cheng, and Huan Liu. 2024. Large Language Models for Data Annotation: A Survey. *arXiv preprint arXiv:2402.13446* (2024).

[257] Zilu Tang, Mayank Agarwal, Alex Shypula, Bailin Wang, Derry Wijaya, Jie Chen, and Yoon Kim. 2023. Explain-then-translate: an analysis on improving program translation with self-generated explanations. *arXiv arXiv:2311.07070* (2023).

[258] Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen, Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L Turner, Callum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers, Edward Rees, Joshua Batson, Adam Jermyn, Shan Carter, Chris Olah, and Tom Henighan. 2024. Scaling Monosemanticity: Extracting Interpretable Features from Claude 3 Sonnet. *Transformer Circuits Thread* (2024). https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html

[259] Himanshu Thakur, Atishay Jain, Praneetha Vaddamanu, Paul Pu Liang, and Louis-Philippe Morency. 2023. Language Models Get a Gender Makeover: Mitigating Gender Bias with Few-Shot Data Interventions. *arXiv arXiv:2306.04597* (2023).

[260] Erico Tjoa and Cuntai Guan. 2020. A survey on explainable artificial intelligence (xai): Toward medical xai. *IEEE TNNLS* 32, 11 (2020), 4793–4813.

[261] E. Todd, M. L Li, A. Sen Sharma, A. Mueller, B. C Wallace, and D. Bau. 2023. Function vectors in large language models. *arXiv* (2023).

[262] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv arXiv:2302.13971* (2023).

[263] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv arXiv:2307.09288* (2023).

[264] Miles Turpin, Julian Michael, Ethan Perez, and Samuel Bowman. 2024. Language models don't always say what they think: unfaithful explanations in chain-of-thought prompting. *NeurIPS* 36 (2024).

[265] E. Voita, R. Sennrich, and I. Titov. 2020. Analyzing the source and target contributions to predictions in neural machine translation. *arXiv* (2020).

[266] Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. 2019. Analyzing Multi-Head Self-Attention: Specialized Heads Do the Heavy Lifting, the Rest Can Be Pruned. In *Proceedings of the 57th ACL*. 5797–5808.

[267] Elena Voita and Ivan Titov. 2020. Information-Theoretic Probing with Minimum Description Length. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 183–196.

[268] Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. 2023. Transformers learn in-context by gradient descent. In *ICML*. 35151–35174.

[269] Sandra Wachter, Brent Mittelstadt, and Chris Russell. 2017. Counterfactual explanations without opening the black box: Automated decisions and the GDPR. *Harv. JL & Tech.* 31 (2017), 841.

[270] David Wadden, Shanchuan Lin, Kyle Lo, Lucy Lu Wang, Madeleine van Zuylen, Arman Cohan, and Hannaneh Hajishirzi. 2020. Fact or Fiction: Verifying Scientific Claims. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 7534–7550.

[271] walkerspider. 2022. Dan is my new friend. https://www.reddit.com/r/ChatGPT/comments/zlcyr9/dan_is_my_new_friend/. [Accessed 27-02-2024].

[272] Ben Wang and Aran Komatsuzaki. 2021. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. Github.

[273] Fei Wang, Ninareh Mehrabi, Palash Goyal, Rahul Gupta, Kai-Wei Chang, and Aram Galstyan. 2024. Data advisor: Dynamic data curation for safety alignment of large language models. *arXiv preprint arXiv:2410.05269* (2024).

[274] Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. [n. d.]. Interpretability in the Wild: a Circuit for Indirect Object Identification in GPT-2 Small. In *The Eleventh International Conference on Learning Representations*.

[275] Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. 2022. Interpretability in the Wild: a Circuit for Indirect Object Identification in GPT-2 Small. In *The Eleventh ICLR*.

[276] Mengru Wang, Ningyu Zhang, Ziwen Xu, Zekun Xi, Shumin Deng, Yunzhi Yao, Qishen Zhang, Linyi Yang, Jindong Wang, and Huajun Chen. 2024. Detoxifying Large Language Models via Knowledge Editing. arXiv:2403.14472 [cs.CL] https://arxiv.org/abs/2403.14472

[277] Mengru Wang, Ningyu Zhang, Ziwen Xu, Zekun Xi, Shumin Deng, Yunzhi Yao, Qishen Zhang, Linyi Yang, Jindong Wang, and Huajun Chen. 2024. Detoxifying large language models via knowledge editing. *arXiv preprint arXiv:2403.14472* (2024).

[278] Qianlong Wang, Keyang Ding, Bin Liang, Min Yang, and Ruifeng Xu. 2023. Reducing Spurious Correlations in Aspect-based Sentiment Analysis with Explanation from Large Language Models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*. 2930–2941.

[279] S. Wang, Y. Zhu, H. Liu, Z. Zheng, and C. Chen. 2023. Knowledge Editing for Large Language Models: A Survey. *arXiv* (2023).

[280] Weixuan Wang, Barry Haddow, and Alexandra Birch. 2023. Retrieval-augmented Multilingual Knowledge Editing. *arXiv arXiv:2312.13040* (2023).

[281] Xiaozhi Wang, Kaiyue Wen, Zhengyan Zhang, Lei Hou, Zhiyuan Liu, and Juanzi Li. 2022. Finding Skill Neurons in Pre-trained Transformer-based Language Models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. 11132–11152.

[282] Xuezhi Wang and Denny Zhou. 2024. Chain-of-Thought Reasoning Without Prompting. *arXiv arXiv:2402.10200* (2024).

[283] Zihao Wang, Shaofei Cai, Guanzhou Chen, Anji Liu, Xiaojian Shawn Ma, and Yitao Liang. 2024. Describe, explain, plan and select: interactive planning with LLMs enables open-world multi-task agents. *NeurIPS* 36 (2024).

[284] Zhao Wang and Aron Culotta. 2021. Robustness to spurious correlations in text classification via automatically generated counterfactuals. In *AAAI*, Vol. 35. 14024–14031.

[285] Ziqi Wang, Yujia Qin, Wenxuan Zhou, Jun Yan, Qinyuan Ye, Leonardo Neves, Zhiyuan Liu, and Xiang Ren. 2019. Learning from Explanations with Neural Execution Tree. In *ICLR*.

[286] Zimu Wang, Wei Wang, Qi Chen, Qiufeng Wang, and Anh Nguyen. 2023. Generating Valid and Natural Adversarial Examples with Large Language Models. *arXiv* (2023).

[287] Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. 2023. Jailbroken: How does llm safety training fail? *arXiv arXiv:2307.02483* (2023).

[288] Hui Wei, Shenghua He, Tian Xia, Fei Liu, Andy Wong, Jingyang Lin, and Mei Han. 2024. Systematic evaluation of llm-as-a-judge in llm alignment tasks: Explainable metrics and diverse prompt templates. *arXiv preprint arXiv:2408.13006* (2024).

[289] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *NeurIPS* 35 (2022), 24824–24837.

[290] Xizi Wei, Mark JF Gales, and Kate M Knill. 2021. Analysing bias in spoken language assessment using concept activation vectors. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 7753–7757.

[291] Zhepei Wei, Wei-Lin Chen, and Yu Meng. 2024. InstructRAG: Instructing Retrieval Augmented Generation via Self-Synthesized Rationales. In *Adaptive Foundation Models: Evolving AI for Personalized and Efficient Learning*.

[292] Orion Weller, Marc Marone, Nathaniel Weir, Dawn Lawrie, Daniel Khashabi, and Benjamin Van Durme. 2024. "According to...": Prompting Language Models Improves Quoting from Pre-Training Data. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*. 2288–2301.

[293] C. Whitehouse, M. Choudhury, and A. F. Aji. 2023. LLM-powered Data Augmentation for Enhanced Crosslingual Performance. *arXiv* (2023).

[294] Jialin Wu, Liyan Chen, and Raymond J Mooney. 2020. Improving VQA and its Explanations by Comparing Competing Explanations. *arXiv* (2020).

[295] K. Wu, E. Wu, and J. Zou. 2024. How faithful are RAG models? Quantifying the tug-of-war between RAG and LLMs' internal prior. *arXiv* (2024).

[296] Tianhao Wu, Janice Lan, Weizhe Yuan, Jiantao Jiao, Jason Weston, and Sainbayar Sukhbaatar. 2024. Thinking llms: General instruction following with thought generation. *arXiv preprint arXiv:2410.10630* (2024).

[297] Xinwei Wu, Junzhuo Li, Minghui Xu, Weilong Dong, Shuangzhi Wu, Chao Bian, and Deyi Xiong. 2023. DEPN: Detecting and Editing Privacy Neurons in Pretrained Language Models. arXiv:2310.20138 [cs.CR] https://arxiv.org/abs/2310.20138

[298] Xuansheng Wu, Wenlin Yao, Jianshu Chen, Xiaoman Pan, Xiaoyang Wang, Ninghao Liu, and Dong Yu. 2023. From language modeling to instruction following: Understanding the behavior shift in llms after instruction tuning. *arXiv arXiv:2310.00492* (2023).

[299] Xuansheng Wu, Wenhao Yu, Xiaoming Zhai, and Ninghao Liu. 2025. Self-Regularization with Latent Space Explanations for Controllable LLM-based Classification. arXiv:2502.14133 [cs.CL] https://arxiv.org/abs/2502.14133

[300] Xuansheng Wu, Jiayi Yuan, Wenlin Yao, Xiaoming Zhai, and Ninghao Liu. 2025. Interpreting and Steering LLMs with Mutual Information-based Explanations on Sparse Autoencoders. *arXiv preprint arXiv:2502.15576* (2025).

[301] Zhengxuan Wu, Aryaman Arora, Atticus Geiger, Zheng Wang, Jing Huang, Dan Jurafsky, Christopher D Manning, and Christopher Potts. 2025. AXBENCH: Steering LLMs? Even Simple Baselines Outperform Sparse Autoencoders. *arXiv preprint arXiv:2501.17148* (2025).

[302] Zhiyong Wu, Yun Chen, Ben Kao, and Qun Liu. 2020. Perturbed Masking: Parameter-free Probing for Analyzing and Interpreting BERT. In *ACL*.

[303] Zhengxuan Wu and Desmond C Ong. 2021. On explaining your explanations of bert: An empirical study with sequence classification. *arXiv* (2021).

[304] Albert Xu, Xiang Ren, and Robin Jia. 2023. Contrastive Novelty-Augmented Learning: Anticipating Outliers with Large Language Models. In *Proceedings of the 61st ACL (Volume 1: Long Papers)*. 11778–11801.

[305] Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023. WizardLM: Empowering Large Language Models to Follow Complex Instructions. arXiv:2304.12244 [cs.CL] https://arxiv.org/abs/2304.12244

[306] Weijia Xu, Sweta Agrawal, Eleftheria Briakou, Marianna J Martindale, and Marine Carpuat. 2023. Understanding and Detecting Hallucinations in Neural Machine Translation via Model Introspection. *Transactions of the Association for Computational Linguistics* 11 (2023), 546–564.

[307] Ziwei Xu, Sanjay Jain, and Mohan Kankanhalli. 2024. Hallucination is inevitable: An innate limitation of large language models. *arXiv* (2024).

[308] Fan Yang, Mengnan Du, and Xia Hu. 2019. Evaluating explanation without ground truth in interpretable machine learning. *arXiv* (2019).

[309] Junjie Yang, Ke Lin, and Xing Yu. 2025. Think When You Need: Self-Adaptive Chain-of-Thought Learning. *arXiv preprint arXiv:2504.03234* (2025).

[310] Xianjun Yang, Shaoliang Nie, Lijuan Liu, Suchin Gururangan, Ujjwal Karn, Rui Hou, Madian Khabsa, and Yuning Mao. 2025. Diversity-driven Data Selection for Language Model Tuning through Sparse Autoencoder. arXiv:2502.14050 [cs.CL] https://arxiv.org/abs/2502.14050

[311] Yuqing Yang, Yan Ma, and Pengfei Liu. 2024. Weak-to-strong reasoning. *arXiv preprint arXiv:2407.13647* (2024).

[312] Yue Yang, Artemis Panagopoulou, Shenghao Zhou, Daniel Jin, Chris Callison-Burch, and Mark Yatskar. 2023. Language in a bottle: Language model guided concept bottlenecks for interpretable image classification. In *CVPR*. 19187–19197.

[313] Zhengyi Yang, Jiancan Wu, Yanchen Luo, Jizhi Zhang, Yancheng Yuan, An Zhang, Xiang Wang, and Xiangnan He. 2023. Large language model can interpret latent space of sequential recommender. *arXiv arXiv:2310.20487* (2023).

[314] Bingsheng Yao, Ishan Jindal, Lucian Popa, Yannis Katsis, Sayan Ghosh, Lihong He, Yuxuan Lu, Shashank Srivastava, Yunyao Li, James Hendler, et al. 2023. Beyond Labels: Empowering Human Annotators with Natural Language Explanations through a Novel Active-Learning Architecture. In *Findings of the Association for Computational Linguistics: EMNLP 2023*. 11629–11643.

[315] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. *arXiv arXiv:2305.10601* (2023).

[316] Xi Ye and Greg Durrett. 2022. The unreliability of explanations in few-shot prompting for textual reasoning. *NeurIPS* 35 (2022), 30378–30392.

[317] Kayo Yin and Graham Neubig. 2022. Interpreting Language Models with Contrastive Explanations. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*. 184–198.

[318] Qingyu Yin, Chak Tou Leong, Hongbo Zhang, Minjun Zhu, Hanqi Yan, Qiang Zhang, Yulan He, Wenjie Li, Jun Wang, Yue Zhang, and Linyi Yang. 2024. Direct Preference Optimization Using Sparse Feature-Level Constraints. arXiv:2411.07618 [cs.AI] https://arxiv.org/abs/2411.07618

[319] Xunjian Yin, Xu Zhang, Jie Ruan, and Xiaojun Wan. 2024. Benchmarking Knowledge Boundary for Large Language Model: A Different Perspective on Model Evaluation. *arXiv arXiv:2402.11493* (2024).

[320] O. Yoran, T. Wolfson, O. Ram, and J. Berant. 2023. Making retrieval-augmented language models robust to irrelevant context. *arXiv* (2023).

[321] Ashkan Yousefpour, Taeheon Kim, Ryan S Kwon, Seungbeen Lee, Wonje Jeung, Seungju Han, Alvin Wan, Harrison Ngan, Youngjae Yu, and Jonghyun Choi. 2025. Representation Bending for Large Language Model Safety. *arXiv preprint arXiv:2504.01550* (2025).

[322] Charles Yu, Sullam Jeoung, Anish Kasi, Pengfei Yu, and Heng Ji. 2023. Unlearning bias in language models by partitioning gradients. In *Findings of the Association for Computational Linguistics: ACL 2023*. 6032–6048.

[323] Lei Yu, Meng Cao, Jackie Chi Kit Cheung, and Yue Dong. 2024. Mechanistic understanding and mitigation of language model non-factual hallucinations. *arXiv preprint arXiv:2403.18167* (2024).

[324] Lei Yu, Virginie Do, Karen Hambardzumyan, and Nicola Cancedda. 2024. Robust LLM safeguarding via refusal feature adversarial training. *arXiv preprint arXiv:2409.20089* (2024).

[325] Yaodong Yu, Kwan Ho Ryan Chan, Chong You, Chaobing Song, and Yi Ma. 2020. Learning diverse and discriminative representations via the principle of maximal coding rate reduction. *NeurIPS* 33 (2020), 9422–9434.

[326] Hao Yuan, Haiyang Yu, Shurui Gui, and Shuiwang Ji. 2022. Explainability in graph neural networks: A taxonomic survey. *IEEE transactions on pattern analysis and machine intelligence* 45, 5 (2022), 5782–5799.

[327] Mert Yuksekgonul, Maggie Wang, and James Zou. 2022. Post-hoc Concept Bottleneck Models. In *The Eleventh ICLR*.

[328] Eric Zelikman, Georges Harik, Yijia Shao, Varuna Jayasiri, Nick Haber, and Noah D Goodman. 2024. Quiet-STaR: Language Models Can Teach Themselves to Think Before Speaking. *CoRR* (2024).

[329] Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. 2022. Star: Bootstrapping reasoning with reasoning. *Advances in Neural Information Processing Systems* 35 (2022), 15476–15488.

[330] Shenglai Zeng, Jiankun Zhang, Pengfei He, Jie Ren, Tianqi Zheng, Hanqing Lu, Han Xu, Hui Liu, Yue Xing, and Jiliang Tang. 2024. Mitigating the privacy issues in retrieval-augmented generation (rag) via pure synthetic data. *arXiv preprint arXiv:2406.14773* (2024).

[331] Yi Zeng, Hongpeng Lin, Jingwen Zhang, Diyi Yang, Ruoxi Jia, and Weiyan Shi. 2024. How johnny can persuade llms to jailbreak them: Rethinking persuasion to challenge ai safety by humanizing llms. *arXiv arXiv:2401.06373* (2024).

[332] Quan-shi Zhang and Song-Chun Zhu. 2018. Visual interpretability for deep learning: a survey. *Frontiers of Info. Tech. & EE* 19, 1 (2018), 27–39.

[333] Xuan Zhang, Chao Du, Tianyu Pang, Qian Liu, Wei Gao, and Min Lin. 2024. Chain of preference optimization: Improving chain-of-thought reasoning in llms. *Advances in Neural Information Processing Systems* 37 (2024), 333–356.

[334] Xiaofeng Zhang, Yihao Quan, Chaochen Gu, Chen Shen, Xiaosong Yuan, Shaotian Yan, Hao Cheng, Kaijie Wu, and Jieping Ye. 2024. Seeing clearly by layer two: Enhancing attention heads to alleviate hallucination in lvlms. *arXiv preprint arXiv:2411.09968* (2024).

[335] Zhihan Zhang, Tao Ge, Zhenwen Liang, Wenhao Yu, Dian Yu, Mengzhao Jia, Dong Yu, and Meng Jiang. 2024. Learn Beyond The Answer: Training Language Models with Reflection for Mathematical Reasoning. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. 14720–14738.

[336] Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. 2023. Multimodal chain-of-thought reasoning in language models. *arXiv arXiv:2302.00923* (2023).

[337] Chenxu Zhao, Wei Qian, Yucheng Shi, Mengdi Huai, and Ninghao Liu. 2023. Automated Natural Language Explanation of Deep Visual Neurons with Large Models. *arXiv arXiv:2310.10708* (2023).

[338] H. Zhao, H. Chen, F. Yang, N. Liu, H. Deng, H. Cai, S. Wang, D. Yin, and M. Du. 2023. Explainability for LLMs: A survey. *TIST* (2023).

[339] H. Zhao, F. Yang, H. Lakkaraju, and M. Du. 2024. Opening the Black Box of LLMs: Two Views on Holistic Interpretability. *arXiv* (2024).

[340] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223* (2023).

[341] Xinyan Zhao and VG Vinod Vydiswaran. 2021. Lirex: Augmenting language inference with relevant explanations. In *AAAI*, Vol. 35. 14532–14539.

[342] Z. Zhao, Y. Shi, S. Wu, F. Yang, W. Song, and N. Liu. 2023. Interpretation of Time-Series Deep Models: A Survey. *arXiv* (2023).

[343] Chujie Zheng, Fan Yin, Hao Zhou, Fandong Meng, Jie Zhou, Kai-Wei Chang, Minlie Huang, and Nanyun Peng. [n. d.]. On Prompt-Driven Safeguarding for Large Language Models. In *ICLR 2024 Workshop on Secure and Trustworthy Large Language Models*.

[344] Chujie Zheng, Fan Yin, Hao Zhou, Fandong Meng, Jie Zhou, Kai-Wei Chang, Minlie Huang, and Nanyun Peng. 2024. On Prompt-Driven Safeguarding for Large Language Models. In *Forty-first International Conference on Machine Learning*.

[345] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems* 36 (2023), 46595–46623.

[346] Zexuan Zhong, Zhengxuan Wu, Christopher D Manning, Christopher Potts, and Danqi Chen. 2023. MQuAKE: Assessing Knowledge Editing in Language Models via Multi-Hop Questions. *arXiv arXiv:2305.14795* (2023).

[347] Kun Zhou, Hui Wang, Wayne Xin Zhao, Yutao Zhu, Sirui Wang, Fuzheng Zhang, Zhongyuan Wang, and Ji-Rong Wen. 2020. S3-Rec: Self-supervised learning for sequential recommendation with mutual information maximization. In *CIKM*. 1893–1902.

[348] Zijian Zhou, Xiaoqiang Lin, Xinyi Xu, Alok Prakash, Daniela Rus, and Bryan Kian Hsiang Low. 2024. DETAIL: Task DEmonsTration Attribution for Interpretable In-context Learning. *arXiv preprint arXiv:2405.14899* (2024).

[349] Yaochen Zhu, Jing Ma, and Jundong Li. 2023. Causal Inference in Recommender Systems: A Survey of Strategies for Bias Mitigation, Explanation, and Generalization. *arXiv arXiv:2301.00910* (2023).

[350] Yaochen Zhu, Liang Wu, Qi Guo, Liangjie Hong, and Jundong Li. 2024. Collaborative large language model for recommender systems. In *WWW*.

[351] Honglei Zhuang, Xuanhui Wang, Michael Bendersky, Alexander Grushetsky, Yonghui Wu, Petr Mitrichev, Ethan Sterling, Nathan Bell, Walker Ravina, and Hai Qian. 2021. Interpretable ranking with generalized additive models. In *WSDM*.

[352] Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, et al. 2023. Representation Engineering: A Top-Down Approach to AI Transparency. *CoRR* (2023).

[353] Andy Zou, Z. Wang, J Z. Kolter, and M. Fredrikson. 2023. Universal and transferable adversarial attacks on aligned language models. *arXiv* (2023).