
A Practical Review of Mechanistic Interpretability for Transformer-Based Language Models

Daking Rai*

George Mason University

drai2@gmu.edu

Yilun Zhou

Salesforce Research

yilun.zhou@salesforce.com

Shi Feng

George Washington University

shi.feng@gwu.edu

Abulhair Saparov

Purdue University

asaparov@purdue.edu

Ziyu Yao*

George Mason University

ziyuyao@gmu.edu

GitHub Paper Collection: <https://github.com/Dakingrai/awesome-mechanistic-interpretability-lm-papers>

Abstract

Mechanistic interpretability (MI) is an emerging sub-field of interpretability that seeks to understand a neural network model by reverse-engineering its internal computations. Recently, MI has garnered significant attention for interpreting transformer-based language models (LMs), resulting in many novel insights yet introducing new challenges. However, there has not been work that comprehensively reviews these insights and challenges, particularly as a guide for newcomers to this field. To fill this gap, we provide a comprehensive survey from a **task-centric** perspective, organizing the taxonomy of MI research around specific research questions or tasks. We outline the fundamental objects of study in MI, along with the techniques, evaluation methods, and key findings for each task in the taxonomy. In particular, we present a *task-centric taxonomy* as a *roadmap for beginners* to navigate the field by helping them quickly identify impactful problems in which they are most interested and leverage MI for their benefit. Finally, we discuss the current gaps in the field and suggest potential future directions for MI research.

*Corresponding authors

Contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 4 |
| 2 | Background: Transformer-based LMs | 5 |
| 3 | What is Mechanistic Interpretability? | 6 |
| 3.1 | Fundamental Objects of Study in MI | 6 |
| 3.1.1 | Features | 6 |
| 3.1.2 | Circuits | 6 |
| 3.1.3 | Universality | 7 |
| 3.2 | How is MI Different From Other Sub-fields of Interpretability? | 7 |
| 4 | Techniques and Evaluation Methods | 8 |
| 4.1 | Vocabulary Projection Methods | 8 |
| 4.1.1 | Basic Concepts | 8 |
| 4.1.2 | Technique-specific Interpretation | 9 |
| 4.1.3 | Technical Advancements | 9 |
| 4.2 | Intervention-based Methods | 10 |
| 4.2.1 | Basic Concepts | 10 |
| 4.2.2 | Technique-specific Interpretation | 12 |
| 4.2.3 | Technical Advancements | 12 |
| 4.3 | Sparse Autoencoder (SAE) | 15 |
| 4.3.1 | Basic Concepts | 15 |
| 4.3.2 | Technique-specific Interpretation | 16 |
| 4.3.3 | Technical Advancements | 16 |
| 4.4 | Others | 17 |
| 4.4.1 | Probing | 17 |
| 4.4.2 | Visualization | 17 |
| 5 | A Beginner’s Roadmap to MI | 18 |
| 5.1 | Feature Study | 18 |
| 5.1.1 | General Workflow for Targeted Feature Study | 18 |
| 5.1.2 | General Workflow for Open-ended Feature Study | 19 |
| 5.1.3 | Evaluation of Feature Study | 20 |
| 5.2 | Circuit Study | 21 |
| 5.2.1 | General Workflow for Circuit Study | 22 |
| 5.3 | Study of Universality | 24 |

| | | |
|----------|--|-----------|
| 5.3.1 | General Workflow for Universality Study | 24 |
| 6 | Case Studies of Beginner’s Roadmap | 25 |
| 6.1 | Case Study for Feature Study | 25 |
| 6.1.1 | Targeted Feature Study with Probing | 25 |
| 6.1.2 | Open-ended Feature Study with SAEs | 26 |
| 6.2 | Case Study for Circuit Study | 26 |
| 7 | Findings and Applications | 27 |
| 7.1 | Feature Study | 27 |
| 7.2 | Circuit Study | 28 |
| 7.2.1 | Interpreting LM Behaviors | 28 |
| 7.2.2 | Interpreting LM components | 29 |
| 7.3 | Universality | 30 |
| 7.3.1 | Universality of Features | 30 |
| 7.3.2 | Universality of Circuits | 30 |
| 7.4 | Findings on Model Capabilities | 30 |
| 7.5 | Findings on Learning Dynamics | 32 |
| 7.6 | Applications of MI | 32 |
| 7.6.1 | Model Enhancement | 32 |
| 7.6.2 | AI Safety | 33 |
| 7.6.3 | Others | 33 |
| 8 | Discussion and Future Work | 33 |
| 8.1 | Advancement in MI Techniques | 33 |
| 8.2 | Practical Utility of MI Studies | 34 |
| 8.3 | Standardized Benchmarks and Metrics | 34 |
| 8.4 | Automated Hypothesis Generation in MI Practices | 34 |
| 8.5 | Generalizing Beyond MI for More Intuitive Applications | 35 |
| 8.6 | A New Paradigm of Human-AI Collaboration Driven by MI | 35 |

| Related Survey | Format | Technical Focus | Findings | Applications | Discussion of Future Research |
|-------------------------|-------------------|---|--------------------------------|---|-------------------------------|
| Zhao et al. (2024) | Technique-centric | Non MI-focused (e.g., CoT explanations) | Explanation methods comparison | Model debugging and improvement | Explainability-focused |
| Ferrando et al. (2024) | Technique-centric | MI-focused | LM properties | Not a focus | Brief |
| Bereska & Gavves (2024) | Technique-centric | MI-focused | LM properties | AI safety | Comprehensive |
| Ours | Task-centric | MI-focused | LM properties | Knowledge editing, LM steering, AI safety, and others | Comprehensive |

Table 1: Comparison of our survey with others on LM interpretability. Our survey brings a unique perspective of *task-centric* survey and includes a comprehensive review of MI techniques, findings, and applications.

1 Introduction

In recent years, transformer-based language models (LMs) have achieved remarkable success in a wide range of natural language processing (NLP) tasks (Radford et al., 2019; Brown et al., 2020; Achiam et al., 2023; Touvron et al., 2023; Bubeck et al., 2023). Alongside these advancements, there are growing concerns over the safety, reliability, generalizability, and robustness of their usage and development (Bengio et al., 2024; Chang et al., 2024; Yao et al., 2024; Weidinger et al., 2022), especially as they are increasingly implemented in real-world applications. These concerns primarily stem from our limited understanding of these LMs and the difficulty in interpreting their behavior.

Recently, mechanistic interpretability (MI) has emerged as a promising technique that fills the gap in the research field of interpretability. This line of methods attempts to interpret LM by reverse-engineering the underlying computation into human-understandable mechanisms (Olah et al., 2020; Elhage et al., 2021). It has shown promise in providing insights into the functions of LM components (e.g., neurons, attention heads), offering mechanistic explanations for various LM behaviors, and enabling users to leverage the explanations to address LM shortcomings and improve their behavior (Wang et al., 2022a; Marks et al., 2024b; Templeton et al., 2024). Despite the promise, however, there are concerns about the scalability and generalizability of MI findings, as well as its practical applications in tackling critical problems such as AI safety (Räuker et al., 2023; Casper, 2023a).

Observing these promises and challenges, we aim to provide a comprehensive review of MI in its applications to interpret transformer-based LMs. In particular, our survey offers a “*task-centric*” perspective, organizing the taxonomy of MI research with specific research questions or tasks and enabling readers new to the field to quickly identify problems of their interest. Specifically, this is achieved by structuring our taxonomy as a *Beginner’s Roadmap to MI* (Section 5), where each category and subcategory are centered on overarching research questions and accompanied by actionable workflows. These workflows outline key steps and applicable techniques, guiding practitioners through the entire process from problem formulation to practical implementation. In addition, we discuss the pros and cons of applicable techniques (e.g., vocabulary project methods [nostalgebraist 2020] vs. sparse autonecoders [Bricken et al. 2023]) and present readers with case studies of how prior works fit within the workflow, aiming to aid them in understanding the practices and choosing the best-fit approaches to meet their need. By doing so, our survey serves as a friendly and practical guide for beginning researchers (e.g., new doctoral students) to quickly grasp the field of MI, and for experienced practitioners to more easily organize their study plan.

There exist surveys on relevant topics, but they differ from ours in a number of key aspects (Table 1). For example, Zhao et al. (2024) presented a broad survey of LM interpretability (e.g., generating textual explanations for LMs) but did not focus on MI. Concurrent with our work, Ferrando et al. (2024) and Bereska & Gavves (2024) reviewed techniques developed for understanding the inner workings of transformer-based LMs; however, they present literature from a technique-centric perspective and have a limited scope (e.g.,

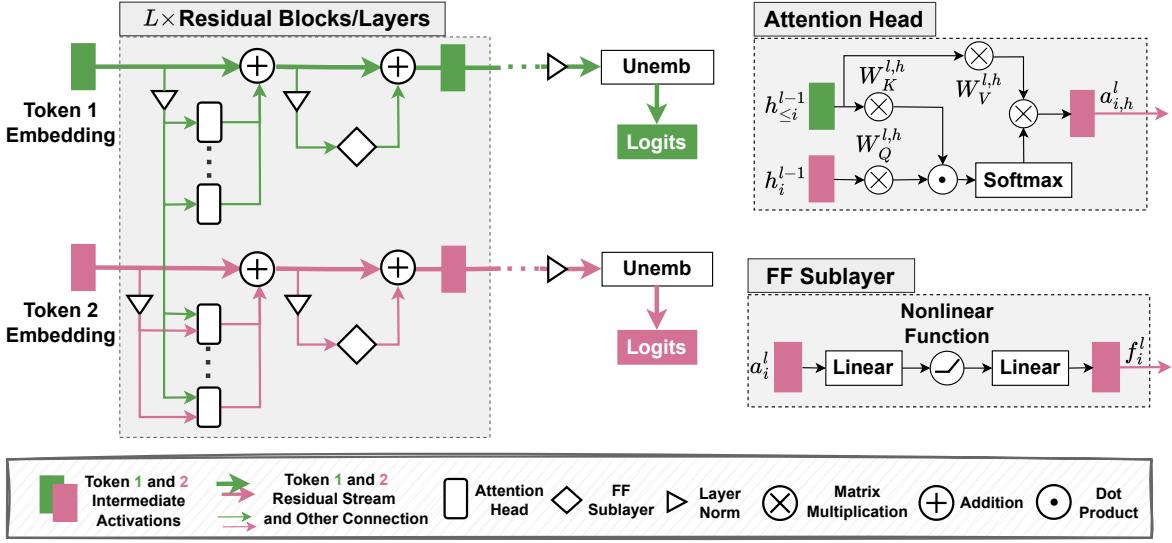


Figure 1: Architecture of transformer-based LMs.

mainly focusing on AI safety and without a focus on the applications of MI). In comparison, we present a unique perspective of a *task-centric* review and provide a comprehensive scope of MI advancements, aiming to enable beginners of MI to quickly pick up the field and leverage MI in their tasks.

Our paper is organized as follows. We first present the background knowledge needed to understand this survey (Section 2). Following this, we introduce the three *fundamental objects of study* of MI, which conceptually structure the research explorations in MI (Olah et al., 2020) and also discuss the current (ambiguous) uses of the term “MI” in literature (Section 3). Next, we introduce various techniques used in MI research (Section 4). For each technique, we introduce (1) its basic concepts and connections to other sub-fields of interpretability or to work that do not target transformer-based LMs, (2) recent advancements, particularly the limitations of the technique and what has or has not been addressed by further research, and (3) technique-specific evaluations. Then, we introduce *Beginner’s Roadmap*, which provides a *task-centric* overview of MI (Section 5). Corresponding to the fundamental objects of study, the roadmap summarizes the actionable workflows along three main lines, i.e., the *study of features*, the *study of circuits*, and the *study of universality*. Subsequently, we present case studies that map the research activities of prior works to the workflow presented in the beginner’s roadmap (Section 6). Next, we present an overview of findings and applications for each study in the roadmap (Section 7). This section includes an in-depth discussion of the findings and novel contributions of various existing works. Finally, we provide careful discussion on the challenges and issues associated with the current development of MI, as well as suggestions for future work (Section 8).

2 Background: Transformer-based LMs

A transformer-based LM (Vaswani et al., 2017) M takes input tokens $X = (x_1, \dots, x_n)$ and outputs a vector in $\mathbb{R}^{|\mathcal{V}|}$, a probability distribution over the vocabulary \mathcal{V} , to predict the next token x_{n+1} . The model refines the representation of each token x_i layer by layer (Figure 1). In the first layer, h_i^0 is an embedding vector of x_i , resulting from a lookup operation in an embedding matrix $W_E \in \mathbb{R}^{|\mathcal{V}| \times d}$.¹

This representation is then updated layer-by-layer through the calculations of *multi-head attention (MHA)* and *feed-forward (FF)* sublayers in each layer, i.e.,

$$h_i^l = h_i^{l-1} + a_i^l + f_i^l, \quad (1)$$

¹For brevity, we omit components such as position embedding and layer normalization in transformer, as they will not affect our discussion of MI. Readers should refer to Vaswani et al. (2017) for a complete description.

where h_i^l denotes the representation of token x_i at layer l , a_i^l is the attention output from the MHA sublayer, and f_i^l is the output from the FF sublayer. The sequence of h_i^l across the layers is also referred to as the *residual stream* (RS) of the transformer in literature (Elhage et al., 2021).

Briefly, the MHA sublayer with H attention heads is implemented via

$$\begin{aligned} a_i^l &= \text{concat}(a_{i,0}^l, \dots, a_{i,H}^l)W_O^l, \\ a_{i,h}^l &= \text{softmax}\left(\frac{(h_{i-1}^{l-1}W_Q^{l,h})(h_{\leq i}^{l-1}W_K^{l,h})^\top}{\sqrt{d_k}}\right) \cdot (h_{\leq i}^{l-1}W_V^{l,h}), \end{aligned} \quad (2)$$

where $a_{i,h}^l$ is the attention output from the h -th head, $W_Q^{l,h}$, $W_K^{l,h}$, $W_V^{l,h}$ are the query, key, and value (learned) projection matrices, respectively, and W_O^l projects the concatenated attention outputs from all heads to the model dimension d .

The FF sublayer then performs two linear transformations over $h_i^{l-1} + a_i^l$ with an element-wise non-linear function σ between them, i.e.,

$$f_i^l = W_v^l \sigma(W_k^l a_i^l + b_k^l) + b_v^l, \quad (3)$$

where W_v^l , W_k^l , b_k^l , and b_v^l are learned parameter matrices and biases. Finally, the RS of x_n at the final layer, h_n^L , is projected into a probability distribution over \mathcal{V} by applying an *unembedding* matrix $W_U \in \mathbb{R}^{d \times |\mathcal{V}|}$ and a softmax operation.

3 What is Mechanistic Interpretability?

The goal of MI is to reverse engineer the detailed computations performed by a model into human-understandable algorithms, similar to how a programmer might try to reverse engineer complicated binaries into human-readable source code. To achieve this goal, MI takes a bottom-up approach by decomposing the model into smaller components and more elementary computations (Zou et al., 2023). By understanding these smaller components and their interactions, MI aims to build a comprehensive understanding of the full model. Below, we first present the three fundamental objects of study in MI, and then discuss its connections and differences with respect to other sub-fields of AI interpretability.

3.1 Fundamental Objects of Study in MI

Following Olah et al. (2020), one of the earliest studies in MI, we categorize research of MI into three areas: the study of features, circuits, and their universality.

3.1.1 Features

A *feature* is a human-interpretable input property that is encoded in LM activations.² For example, when we provide an LM with the input token “*dog*”, it may extract features like “*animal*”, “*pet*”, “*has four legs*”, and other relevant features learned during pre-training, embedding these features in its activations. MI aims to interpret the LM representations by decoding the features encoded in them.

3.1.2 Circuits

While the study of features helps us to understand what information is encoded in a model’s activations, it does not inform us of how these features are extracted from the input, their interactions, or how they are used by the model to enable specific LM behaviors (e.g., reasoning). MI bridges this gap by investigating *circuits* – meaningful computational pathways that connect features and facilitate specific LM behaviors.

More formally, if we view an LM M as a computational graph with features as nodes and the weighted connections between them as edges, a circuit is a sub-graph of M responsible for implementing specific LM behaviors (Olah et al., 2020). Additionally, although circuits were initially defined as connections

²We use the terms “representations” and “activations” interchangeably.

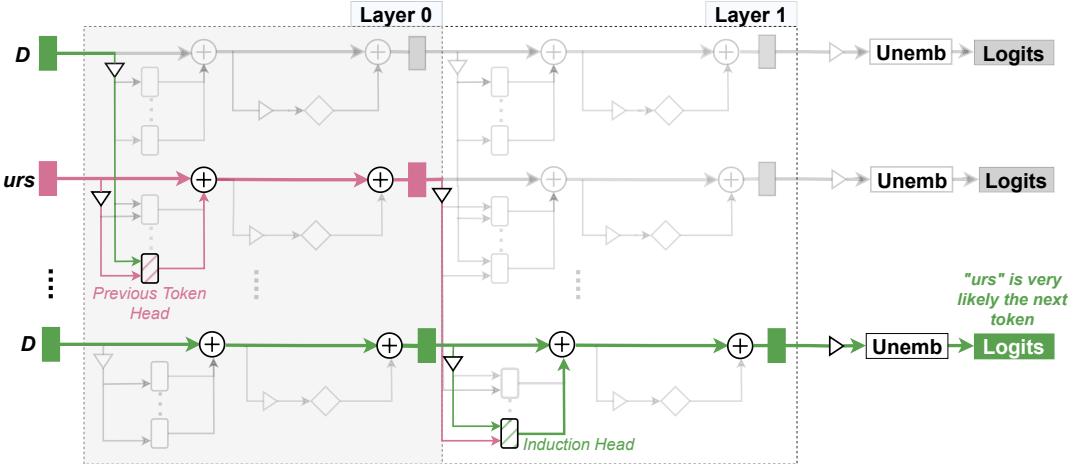


Figure 2: An example of an induction circuit discovered by Elhage et al. (2021) that consists of outputs of two attention heads (previous token head and induction head) as nodes and connection between them as edges.

between features (Olah et al., 2020), subsequent studies have generalized them as connections between the activation outputs of *transformer components* (Olsson et al., 2022; Wang et al., 2022a), where interpreting individual transformer components becomes part of the circuit interpretation. Therefore, we include research on interpreting individual transformer components as part of the study of circuits as well. As we will introduce in Section 5.2, prior works have defined circuits with various assumptions and at different levels of granularity. Some works even generalized the definitions of circuits to cover general information flows in an LM (Geva et al., 2023; Nikankin et al., 2024). In this survey, we will consider all of these works as studies of circuits, as long as the goal is to find computational pathways connecting transformer components or features in an LM.

An example circuit discovered by Elhage et al. (2021) in a toy LM is shown in Figure 2. This is an induction circuit consisting of two attention heads (previous token head and induction head) as nodes and input/output activations between them as edges of the circuit. The circuit implements the task of detecting and continuing repeated subsequences in the input (e.g., “Mr D urs ley was thin and bold. Mr D” -> “urs”), where the *previous token head* encodes the information that “*urs*” follows the “*D*” token in the RS, which is then read by the *induction head* to promote “*urs*” as the next token prediction.

3.1.3 Universality

For any feature or circuit that we have identified in an LM in one task, the critical question arises: *Do similar features and circuits exist in other LMs or tasks?* The investigation into this question has then given rise to the notion of *universality*, i.e., the extent to which similar features and circuits are formed across different LMs and tasks (Olah et al., 2020; Gurnee et al., 2024). The implications of universal features and circuits can be significant. For instance, many studies on features and circuits (Olsson et al., 2022; Elhage et al., 2022a;b) were performed with only toy or small LMs. If these features and circuits are universal, the insights from these studies can be transferred to other unexamined LMs and potentially state-of-the-art large LMs (LLMs). However, if they are not universal, then a significant amount of independent effort will be required to interpret each uninterpreted LM in each task.

3.2 How is MI Different From Other Sub-fields of Interpretability?

Although MI has become a popular sub-field within interpretability, it is still arguable what exactly makes an interpretability study *mechanistic*. Historically, the term *Mechanistic Interpretability* was coined to distinguish the field from interpretability approaches, such as saliency methods (Ribeiro et al., 2016; Sundararajan

et al., 2017; Lundberg, 2017), that generated explanations solely by analyzing inputs and outputs and provided little to no insight into the internal mechanisms of the model (Olah et al., 2020; Saphra & Wiegreffe, 2024). While this distinction clarifies that MI involves investigating the internals of a model, it does not address how MI is distinct from other interpretability research that also explores model internals. For instance, even prior to MI there was already a significant body of work, such as probing (Belinkov, 2022) and attention analysis (Vig, 2019), that investigated the internal representation of models. Relevant to this discussion, Zou et al. (2023) categorized interpretability work that investigates the internals of models into two broad categories, i.e., *top-down* and *bottom-up* approaches, with MI classified as bottom-up. Specifically, the *bottom-up* approach begins by breaking down the model into its smallest units of analysis, such as neurons, aiming to first understand these components before piecing them together to understand how they combine to produce the model’s overall behavior. In contrast, the *top-down* approach starts with higher-level model behaviors or coarser components and works downward to explore finer-grained mechanisms. Despite these distinctions, in practice, the boundaries between the two approaches often blur, as both ultimately involve studying neurons and representations within models. This lack of clarity has led to confusion within the community, with several researchers publicly expressing concerns that many MI studies are “not new” and often disregard prior literature, reinvent or rediscover existing techniques, or overlook important baseline work (Saphra & Wiegreffe, 2024). To provide a clearer historical context of MI, Saphra & Wiegreffe (2024) examined the historical evolution of the field and proposed four definitions of MI. In our survey, we adopt the *broad technical definition of MI* from Saphra & Wiegreffe (2024) and define *MI as any work that seeks to describe the internals of a model*. Additionally, we also urge the MI community to bridge the existing disconnect by more actively integrating studies from related non-MI research that explores similar topics.

4 Techniques and Evaluation Methods

In this section, we review major techniques that have been developed to study the fundamental objects described in Section 3.1 for understanding transformer-based LMs. During the introduction of each technique, we will present its basic concepts, technique-specific evaluation, and recent advancements.

4.1 Vocabulary Projection Methods

4.1.1 Basic Concepts

This class of techniques aims to decode the information within the LM representation by projecting it to the model’s vocabulary space. Typically, during the LM inference, only the activation of the final layer (h_i^L) is multiplied by the unembedding matrix (W_U) to calculate the logit distribution for the next token prediction. The vocabulary projection method, however, proposes to multiply W_U with intermediate representations (e.g., $h_i^l, 0 \leq l < L$), generating logit distributions from these intermediate layers as well. By examining tokens with the highest logit values, one can infer the candidate tokens the LM is considering for the next prediction within that specific intermediate representation.

This technique can be viewed as a practical application of the *iterative inference perspective* (Jastrzebski et al., 2017; Elhage et al., 2021; Geva et al., 2021) on how an LM makes predictions, where each layer of an LM is seen as progressively refining a latent prediction of the next token, and the vocabulary projection technique enables us to examine how these predictions are refined and evolve across the model’s intermediate layers.

Logit lens (nostalgebraist, 2020) was the first technique that proposed to employ W_U for projecting the intermediate activations of an LM (GPT-2 [Radford et al. 2019] in the original work) to the vocabulary space. We illustrate the method in Figure 3. Specifically, given an activation h_i^l in the residual stream, logit lens calculates the following:

$$\text{LogitLens}(h_i^l) = \text{LayerNorm}(h_i^l)W_U \quad (4)$$

Similarly, the logit lens was also applied to project the output activations of MHA (a_i^l) and FF sub-layers (f_i^l) to the vocabulary space before they are added to the residual stream, so as to interpret the contribution of MHA and FF for the next token prediction.

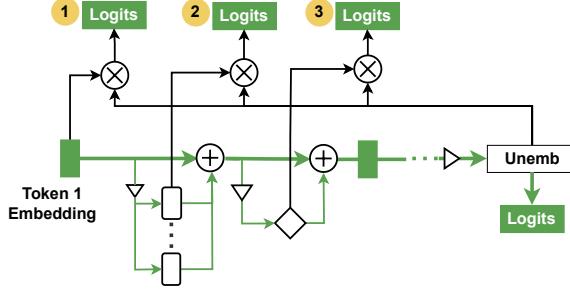


Figure 3: Logit lens implementation at (1) RS, (2) attention head, and (3) FF sublayer.

4.1.2 Technique-specific Interpretation

As described earlier, once the logit distribution projected from the intermediate activation of interest is obtained, one can infer some of the information encoded within activation by examining the tokens with the highest logit values. This interpretation process is typically performed by human evaluators. For example, if the top tokens primarily relate to items that are associated with breakfast (e.g., pancakes, coffee, bread), it suggests that the activation encodes the “breakfast” feature. To automate this manual evaluation process, Bills et al. (2023) proposed using a large LM (e.g., GPT-4) to automatically generate textual explanations based on the output of the logit lens. Similarly, Rai & Yao (2024) proposed using an LM to automatically classify whether the logit distribution likely encodes a pre-defined feature category (e.g., arithmetic addition). Lastly, it is important to note that the explanation generated by vocabulary projection methods is only correlational rather than causal in nature (Katz et al., 2024), requiring caution when drawing inferences from these methods.

4.1.3 Technical Advancements

Since the introduction of the logit lens, several advancements have been made to enhance its capabilities or to address its limitations, leading to the development of various vocabulary projection methods. Specifically, the advancements can be categorized into the following three categories – (1) Reliability: improving the faithfulness of the technique’s output and enabling more reliable vocabulary projection (2) Decoding position: adapting and applying the technique to decode features from various intermediate representations and model weights; and (3) Decoding expressivity: enhancing the degree to which the technique decodes features beyond immediate next-token predictions.

Reliability Belrose et al. (2023); Din et al. (2023) noted that the logit lens may produce unreliable projection for certain LMs (e.g., BLOOM, OPT 125M), particularly in the earlier layers. For example, the top-1 projected token for BLOOM (Scao et al., 2022) in more than half of the earlier layers is often the input token, rather than a plausible continuation towards the final output token. The authors posit that the logit lens makes an oversimplified assumption that all layers operate in the same embedding space, which may not be true. Consequently, W_U , which is only pre-trained to project an LM’s final-layer activation, cannot reliably project all the intermediate activations to the vocabulary space. This observation inspired Belrose et al. (2023, i.e., *tuned lens*) and Din et al. (2023) to train translators that transform the intermediate activations to an output representation that better aligns with the representation space of the final layer before the logit lens is applied. Specifically, both approaches linearly transformed the intermediate activation (e.g., h_i^l) to an activation that matches more closely with the basis of the final layer activation (h_i^L). However, Din et al. (2023) achieved it by directly minimizing the difference between the transformed activation and the final-layer activation, while Belrose et al. (2023) designed the loss function to encourage the matching between the logit distributions of the transformed activation and the final-layer activation. Finally, Sakarvadia et al. (2023) proposed the *attention lens*, which similarly involved training translators for projecting output activations of attention heads with the logit lens.

Decoding Position While the logit lens was initially proposed to decode an LM’s prediction information from *activations*, subsequent work has applied or adapted the technique to decode information from *model weights* and *gradients*. We summarize the various decoding positions as follows. (1) **Activations:** In addition to the logit lens (nostalgia, 2020), the *tuned lens* (Belrose et al., 2023) and the approach of Din et al. (2023) were both proposed to project the activations of RS, FF and MHA sub-layers into the vocabulary space. In contrast, the *attention lens* (Sakarvadia et al., 2023) was proposed specifically for projecting the output activations of individual attention heads within MHA sub-layers. (2) **Model weights:** Geva et al. (2022) showed that the logit lens can also be used to interpret the *weights* of transformer components. Specifically, Geva et al. (2022) used it to understand the role of the FF sublayer in the LM’s prediction by projecting columns of the parameter matrix W_v^l to the vocabulary space. Dar et al. (2023) further proposed to apply the technique to project the FF sub-layers (W_k^l and W_v^l) and the query-key ($W_Q^{l,h}W_K^{l,h}$) and value-output ($W_V^{l,h}W_O^l$) interaction parameter matrices of individual attention heads to the vocabulary space, aiming to study how they transfer and mix information from source tokens to the target token. (3) **Gradients:** As a representative work, Katz et al. (2024) proposed *backward lens*, which projects the gradient of FF sub-layers during back-propagation to the vocabulary space, enabling the investigation of how new information is stored in the FF sub-layers of LMs during training.

Decoding Expressivity Several methods have extended the logit lens to decode information beyond the immediate next-token prediction. For instance, Pal et al. (2023) proposed the *linear model approximation* technique that can be used to determine if an intermediate activation h_i^l at i^{th} decoding timestep has already encoded information about the future output tokens h_{i+n}^L , where $n \geq 2$. Specifically, Pal et al. (2023) learned a linear model that transforms h_i^l to h_{i+n}^L and then used the unembedding matrix W_U to project the transformed activations to vocabulary space. Pal et al. (2023) further proposed two other approaches integrating vocabulary projection with intervention, including the *future lens* approach, and Ghandeharioun et al. (2024) concurrently proposed a generalized framework called *PatchScope*, which similarly performs intervention while projecting an activation to the vocabulary space; we will discuss both approaches in detail when we introduce intervention-based methods (Section 4.2). Similarly, Hernandez et al. (2023) proposed the *attribute lens* that can be used to decode which objects (e.g., “Seattle”) are encoded in the activation for a relation (e.g., “is located”), given a subject (e.g., “The space needle”) as the model input. Specifically, they proposed to learn a linear function for a given relation, which can transform the intermediate activation when the subject is used as input into an object representation. Subsequently, the transformed object representation can be decoded into a vocabulary space or object-token distribution using the unembedding matrix W_U . Finally, Cancedda (2024) introduced *Logit Spectrology*, a method that applies singular value decomposition to W_U . Particularly, they found that the logit lens projection is dominated by the larger singular value vectors, while it fails to capture information projected by the other singular value vectors. Building on this observation, they proposed spectral filtering, a method that applies filters to selectively pass information encoded in specific bands of the representation. As a result, this approach allows one to discover information that is typically overlooked by the logit lens, such as how an LM avoids using outputs from irrelevant LM components during next-token predictions.

4.2 Intervention-based Methods

4.2.1 Basic Concepts

Intervention-based methods investigate an LM’s behaviors by directly altering the value of its intermediate representation during the forward pass (i.e., *intervention*) and observing the resulting change of the model output. In this approach, we view an LM as a computational or a causal graph (Pearl, 2009), where its transformer components $\{C_0,..C_n\} \in \mathcal{C}$ are the nodes, and their connections represent the edges. By intervening on these nodes and edges, we can explore the causal relationships between different components and their influence on the model’s output. In practice, the component C_i can be defined at varying levels of granularity, such as a single neuron, an attention head, or an entire transformer layer.

There are mainly two types of intervention-based methods: *Noising-* and *Denoising*-based intervention methods.

| Intervention Type | Corruption Technique | Clean Run | Corrupt Run | Patch Run |
|-------------------|---|----------------------|---|--|
| Noising-based | Zero Ablation (Olsson et al., 2022) | Run with X_{clean} | N/A | Perform clean run but replace activation of C_i with zero-vector |
| | Random-noise Ablation (Rai & Yao, 2024) | Run with X_{clean} | N/A | Perform clean run but add random noise to activation of C_i |
| | Mean Ablation (Wang et al., 2022a) | Run with X_{clean} | Run with multiple $X_{corrupt}$ | Perform clean run but replace the activation of C_i with the mean activation of multiple $X_{corrupt}$ |
| | Resampling (Chan et al., 2022a) | Run with X_{clean} | Run with one $X_{corrupt}$ | Perform clean run but replace activation of C_i with the activation from corrupt run |
| Denoising-based | Random-noise (Meng et al., 2022) | Run with X_{clean} | Add noise to the embedding of X_{clean} | Perform corrupt run but replace the activation of C_i with its activation from the clean run |
| | Resampling (Hanna et al., 2024a) | Run with X_{clean} | Run with $X_{corrupt}$ | Perform corrupt run but replace activation of C_i with its activation from the clean run |

Table 2: A summary of corruption techniques in intervention-based methods.

Noising-based Interventions This line of methods involves removing (also referred to as *noising*) the contribution of a specific component C_i during the LM’s forward pass and observing the resulting change in the model’s output Y . The idea is that if the component C_i is important to the original model output, there should be a significant change in the model output after the contribution of C_i is removed from the model inference. In other words, noising-based methods aim to discover components that are *necessary* for an LM to exhibit a certain behavior as removing the contributions of these components can break the LM’s behavior (Heimersheim & Nanda, 2024). Noising-based interventions are also known as *ablation* (Chan et al., 2022a) or *knockout* (Wang et al., 2022a).

Denoising-based Interventions This line of methods typically involves two interventions to measure the importance of a component C_i for the model output Y . First, we perform *noising* to all LM components \mathcal{C} such that the model does not predict Y anymore. Next, only C_i is *denoised* to see how much C_i alone restores the likelihood (alternatively, logit; see Section 4.2.2) of the model predicting y . If we observe a substantial increase in the likelihood of predicting Y as the next token, it shows that C_i plays an important role in the model’s original prediction. In other words, denoising-based methods identify components that are *sufficient* for restoring an LM’s behavior to some degree. Denoising-based interventions have also been termed as *causal tracing* or *causal mediation analysis* (Meng et al., 2022).

Intervention Procedure The typical procedure for intervention-based methods, both noising and denoising, for measuring the contribution of component C_i is outlined below. Conceptually, both types of intervention involve three runs.

- **Clean run:** Run the model M with a prompt X_{clean} (e.g., “The capital city of France is”) that showcases the LM behavior Y_{clean} (e.g., “Paris”) and cache the activations of all LM components.
- **Corrupt run:** Run the model M with a corrupted prompt $X_{corrupt}$ by selecting a prompt distinct from X_{clean} (e.g., “The capital city of Italy is”) which leads to model predicting $Y_{corrupt}$ (e.g., “Rome”) instead of Y_{clean} . Cache the activations of all LM components from this corrupt run. In practice, one may implement the corrupt run in various ways. When patching C_i with a zero vector, a corrupt prompt is not needed. One may not need a discrete prompt of $X_{corrupt}$ either and can instead apply continuous noise to the clean activation of C_i . We summarize the various corruption techniques in Table 2.
- **Patch run:** This run can be performed in two ways based on whether you want to perform *noising* or *denoising* intervention: (a) **Noising:** Run the model M with clean prompt X_{clean} but replace or *patch* the activation of C_i from the corrupt run. This patch removes C_i ’s contribution towards the model predicting Y_{clean} when it is provided with the input X_{clean} . (b) **Denoising:** Run the model M with corrupted prompt $X_{corrupt}$ but replace or *patch* the activation of C_i from the *clean run*.

This patch includes only C_i 's contribution towards the model predicting Y_{clean} when it is provided with the input X_{clean} .

4.2.2 Technique-specific Interpretation

Once we perform the three runs, we measure the patching effect to evaluate the importance of a component C_i for the LM output Y . For noising, we analyze the difference between the output of the clean and patch run to observe the resulting changes after removing the contribution of C_i from the clean run. Alternatively, for denoising, we analyze the difference between the output of the corrupted and the patch run to observe how much patching C_i recovers the clean run's output Y_{clean} . Specifically, we can measure the change in the LM output using the following metrics:

- **Probability or Logit:** Measure the change in the softmax probability or logit of Y_{clean} before and after the patching.
- **Probability or Logit Difference:** Measure the change in the difference between the probability or logit of Y_{clean} and that of $Y_{corrupt}$ before and after the patching.
- **KL Divergence:** Measure the Kullback-Leibler (KL) divergence between the logit or probability distributions before and after the patching. This metric compares the full output distributions rather than focusing solely on the change in the logit or probability of a single token.

Which Metric to Select? In general, logit difference is recommended because it allows us to control things we do not want to measure (e.g. components that promote both Y_{clean} and $Y_{corrupt}$), which is not possible using the absolute probability or logit. In addition, probability as a metric may fail to detect *negative model components* that suppress the correct output, as empirically observed by Zhang & Nanda (2023). While KL divergence does not focus on specific token logits or probabilities, it is also a reasonable metric for assessing the effect of patching. For a more detailed comparison of the reliability of these metrics across different scenarios, we refer readers to Heimersheim & Nanda (2024); Zhang & Nanda (2023).

4.2.3 Technical Advancements

Intervention-based techniques have undergone several refinements to ensure their correct usage, including various adaptations to address diverse investigative objectives. Specifically, we summarize these technical advancements as follows: (1) Reliability of corruption: ensuring that the corrupt run does not take the model out of distribution (OOD); (2) Intervention for localization: adapting the techniques to identify important components associated with specific LM behavior; (3) Scaling up and automating localization techniques: automating the intervention process to alleviate the need for human effort; (4) Intervention for validating hypotheses: applying intervention to validate a hypothesized interpretation of the LM; and (5) Integrating intervention with vocabulary projection for activation decoding: devising intervention methods for more expressive vocabulary projections of activations.

Reliability of Corruption The corruption of activations during the corrupt run can push the model out of distribution (OOD) as the modified activations may not resemble anything the model encountered during training. As a result, any decline in model performance may be due to the off-distribution behavior of the model, rather than the removal of C_i 's contribution to the behavior of interest (Zhang & Nanda, 2023). Specifically, Zero and Random-noise ablations can be unreliable because both replacement values are arbitrary choices and the model may not have encountered them during training. Mean ablation attempts to address this issue by replacing the activation of C_i with the average activation of C_i computed from multiple samples drawn from the same data distribution. This approach partially mitigates the OOD concern as it obtains the replacement value from in-distribution activations, making it more reliable than Zero and Random-noise ablation. However, the mean value could still push the LM off distribution if the activation distribution is non-linear (Chan et al., 2022a; Geiger et al., 2021). For example, if the activation distribution of C_i consists of points along the circumference of a circle, the mean of these points would be at the center, rather than at another point along the circumference (Chan et al., 2022a). Resampling ablation addresses

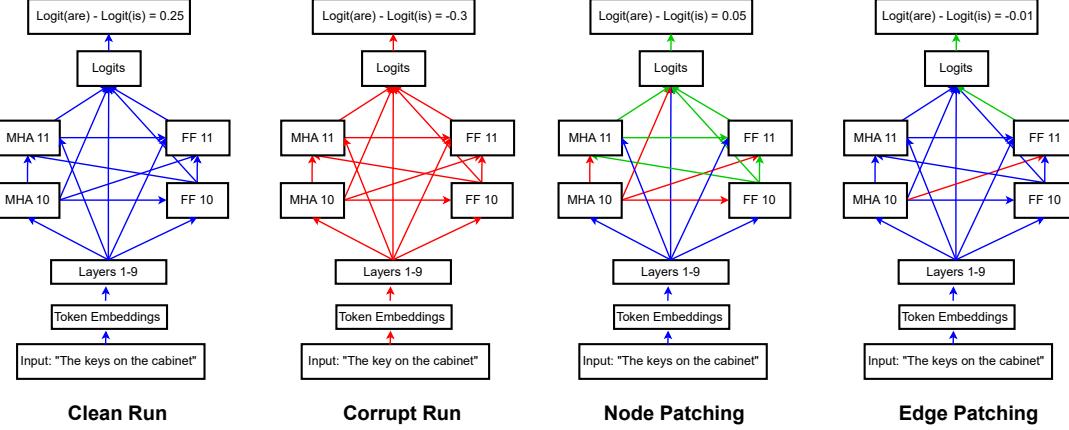


Figure 4: An example of *noising-based* intervention on GPT-2 at the *node* and the *edge* levels. In the example, a corrupt prompt, with the singular “key” replaced by a plural “keys”, was introduced to discover the circuit in GPT-2 for associating the verb with different forms of the noun. (1) Clean Run: The computational graph of GPT-2 when run on clean input. (2) Corrupt Run: The computational graph of GPT-2 when run on corrupt input. (3) Node Patching: The computational graph of GPT-2 when run on clean input, with the output activation of MHA10 replaced by its activation from the corrupt run (i.e., patching the node of “MHA10”). (4) Edge Patching: The computational graph of GPT-2 when run on clean input, with activations along the edge of (MHA10, FF11) replaced by activations from the corrupt run (i.e., patching the edge of “(MHA10, FF11)”). Colored connections indicate activations calculated from the clean run (Blue), the corrupt run (Red), and the re-calculation after patching (Green).

this issue by simply replacing the activation of C_i with the activation from a counterfactual example sampled from the same distribution (Chan et al., 2022a).

Interventions for Localization As we will introduce in our roadmap (Section 5), intervention is commonly used to localize important components (i.e., *nodes* of circuits) and computational pathways (i.e., *edges* of circuits) when one tries to discover circuits in an LM. Depending on the goal, the scope of intervention varies. **(1) Patching nodes vs. edges:** *Activation patching* is a widely used technique employed to localize key components (Meng et al., 2022). This method evaluates the importance of a component C_i by patching the activation of C_i from the *clean run* into the *corrupt run*, one component at a time, to observe its impact on the model’s output, as illustrated in Figure 4. The higher impact implies C_i is an important component of the behavior. *Path Patching* is an extension of activation patching to localize important edges between components (Wang et al., 2022a; Goldowsky-Dill et al., 2023). For instance, to assess whether the connection between two components, C_1 and C_2 , is significant, path patching applies activation patching to the output of C_1 *but, notably, only along paths serving as input to C_2* . In other words, when performing path patching, we allow for the effect of a patched activation only along the computational path we want to investigate, while any other directed nodes will not receive the patched activation despite the connection. Similar to patching nodes, if a change in the LM’s behavior is observed, the connection between the two components is considered important. **(2) Patching at various granularity levels:** As discussed in Section 3.1.1, it is hypothesized that LMs encode features within a linear subspace of the representation space, potentially represented in a non-standard basis, as indicated by the discovery of polysemantic neurons (see Section 7.1 for further discussion). However, activation patching focuses on patching the entire representation. To address this, Geiger et al. (2024) proposed *Distributed Interchange Interventions*, also known as *Distributed Alignment Search* (DAS), which targets interventions within specific subspaces instead of patching full representations.

Scaling Up and Automating Interventions for Localization Performing interventions to localize important components or connections requires humans in the loop, as it involves iteratively patching each component or connection within an LM to assess its significance, which can be time-consuming. To automate

this process, Conmy et al. (2023) introduced *ACDC* (*Automatic Circuit DisDiscovery*), which iteratively knocks out edges in the computational graph and removes any edge whose effect on the target metric is less than a specified threshold. However, ACDC is not scalable to large LMs as it requires independent inferences for every iteration. To address the issue, Syed et al. (2023) proposed *Edge Attribution Patching (EAP)* for locating important edges, where it employs *Attribution Patching* (Nanda et al., 2022) to approximate activation patching for locating important components, which requires only two forward passes and one backward pass for measuring all model components. In addition, Hanna et al. (2024b) recently augmented this approach with Integrated Gradient (Sundararajan et al., 2017) to address the concern of zero gradients, leading to the approach of *EAP-IG*.

Interventions for Validating Hypotheses Intervention-based methods are widely used to causally validate hypotheses about how specific LM components contribute to its behavior. For instance, *Causal Scrubbing* (Chan et al., 2022a) is a popular technique for formalizing and validating a hypothesis about the function of an LM component for the given behaviors. Specifically, causal scrubbing involves formalizing a hypothesis (G, I, c) , where G is the model’s original computational graph, I is an interpretation graph that reflects the hypothesized roles of the model’s components, and c is a function that maps between nodes of I to G . For example, when an LM is tasked to perform a two-digit addition, the interpretation graph I could specify the hypothesized computational pathway, where only certain transformer components are used for implementing this addition while others do not contribute to this task. To validate that I is a faithful interpretation of G for a given behavior, *resampling ablations* of component activations in G is performed. In our example, it means that replacing activations of model components that do not contribute to the addition task should not change the LM’s prediction behavior. This provides a more rigorous way to check our intuitions about how models work, compared to just looking at the model or manually changing its parts. Similarly, *Interchange Intervention* (Geiger et al., 2021) is another hypothesis verification technique that similarly formulates the hypothesis by representing the model M as a high-level conceptual model and performing resampling ablation on the low-level (i.e., the original) model to verify the hypothesis.

Integrating Intervention with Vocabulary Projection for Activation Decoding The intervention technique has also been adapted by Pal et al. (2023); Ghandeharioun et al. (2024) to decode features from intermediate activations, which can be viewed as augmented vocabulary projection (Section 4.1) with more powerful expressivity. For instance, to decode features from the end-position activation h_n^l , the *fixed prompt causal intervention* (Pal et al., 2023) runs another LM inference on a fixed generic prompt (e.g., “Tell me something about”) with its activation at the same location being patched by the target one, and then let the LM continue the generation. This generation can then be viewed as the projected features encoded in the activation h_n^l . When this fixed prompt is replaced by one learned for effective information elicitation, it results in the variant of *learned prompt causal intervention* (or *future lens*, as the initial goal of Pal et al. 2023 was to probe the future tokens from the end-position activation). Concurrently, *Patchscopes* (Ghandeharioun et al., 2024) provides a unifying framework that similarly combines vocabulary projection with intervention methods but is more generalized than the approaches of Pal et al. (2023). Specifically, Patchscopes defines a source prompt and a source LM, which produces the activation to be interpreted, as well as a target prompt and a target LM, which facilitates the activation interpretation. Patchscopes works by patching the activation derived from the source to the corresponding activation in the target. By carefully designing the target prompt (also termed as “inspection prompt”), Patchscopes allows one to decode the interested feature information. For instance, to check if an LM’s activation obtained from the end position of the source prompt “Amazon’s former CEO” encodes the feature of “Jeff Bezos”, an inspection prompt (or $X_{corrupt}$) can be “cat → cat; 135 → 135; hello → hello; ?”, which comprises few-shot demonstrations encouraging the model to explain its internal representation, with the final “?” serving as a placeholder. If the LM with the activation at the position of “?” patched by the source activation can generate “Jeff Bezos” as the continuation, we can infer that the activation contains information about this person’s name.

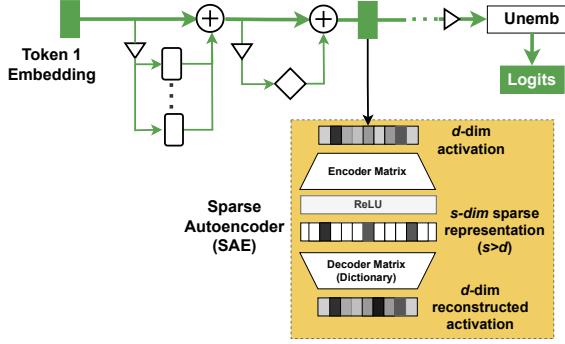


Figure 5: Sparse Autoencoder (SAE) applied to activation on RS.

4.3 Sparse Autoencoder (SAE)

4.3.1 Basic Concepts

Sparse Autoencoders (SAEs) are employed in feature discovery to tackle the problem of *superposition*, where LM activations encode more features than their dimensions, leading to *polysemantic* neurons that activate for multiple unrelated concepts. SAEs tackle this issue by mapping a d -dimensional activation into an s -dimensional sparse representation ($s > d$), where this new s -dimensional representation consists of neurons that are more *monosemantic*, each associated with only a single feature, making them more interpretable (Bricken et al., 2023; Gao et al., 2024; Rajamanoharan et al., 2024b). In other words, the primary goal of SAEs is to transform the less interpretable LM activation into a more interpretable sparse activation.

SAE Architecture An SAE (Figure 5) consists of an *encoder* that maps the d -dimensional input activation into an s -dimensional sparse activation and a *decoder* which reconstructs the original input activation from the sparse representation, defined as follows:

$$\begin{aligned} f(h) &= g(W_{enc}h + b_{enc}) \\ \hat{h} &= W_{dec}f(h) + b_{dec} \end{aligned} \tag{5}$$

Here, the encoder consisting of weights $W_{enc} \in \mathbb{R}^{d \times s}$, a bias term $b_{enc} \in \mathbb{R}^s$, and a non-linear activation function g (e.g., ReLU (Agarap, 2018)), maps the input activation $h \in \mathbb{R}^d$ to sparse activation $f(h) \in \mathbb{R}^s$. Subsequently, the decoder consisting of weights $W_{dec} \in \mathbb{R}^{s \times d}$ and a bias term $b_{dec} \in \mathbb{R}^d$, takes the sparse activation $f(h)$ from the encoder and reconstructs the input activation \hat{h} .

Loss Function An SAE is trained in an unsupervised approach using a dataset \mathcal{X} with a loss function defined as:

$$\mathcal{L}(h, \hat{h}) = \frac{1}{|\mathcal{X}|} \sum_{X \in \mathcal{X}} \left(\underbrace{\|h(X) - \hat{h}(X)\|_2^2}_{\mathcal{L}_{reconstruct}} + \underbrace{\lambda \|f(h(X))\|_1}_{\mathcal{L}_{sparsity}} \right) \tag{6}$$

where $\mathcal{L}_{reconstruct}$ ensures the faithfulness or fidelity of the reconstructed activation to the input, while the sparsity penalty $\mathcal{L}_{sparsity}$ restricts the encoder to activate only a small number of neurons for a given h , encouraging the SAE to learn a sparse interpretable representation for h during training. Specifically, without a sparsity penalty, the encoder could simply memorize the input activations since $f(h)$ has more dimensions than h . However, introducing a sparsity penalty forces the encoder to reconstruct h while activating only a small subset of neurons. This constraint encourages the SAE to capture the most important features of h in $f(h)$, as doing so minimizes the overall loss. As a result, the SAE learns to disentangle and represent the most salient features of h more effectively. In addition, these two objectives (fidelity and sparsity) are balanced by the L_1 coefficient λ .

Decoding Positions SAEs can be trained on activations from any LM component. For example, Templeton et al. (2024); Bricken et al. (2023); Templeton et al. (2024); Sharkey et al. (2022) trained SAEs on intermediate activations from the residual stream, Kissane et al. (2024); Krzyzanowski et al. (2024) trained SAEs on attention outputs, while Braun et al. (2024) trained SAEs on the outputs of FF sub-layers.

4.3.2 Technique-specific Interpretation

The fidelity of the SAE is measured using the reconstruction loss, where a low reconstruction loss implies high fidelity. However, evaluating the interpretability or quality of the features learned by an SAE remains an open research challenge. While increased sparsity can enhance interpretability, it does not necessarily guarantee highly interpretable features. Existing work has focused on analyzing the activation patterns of features with particular emphasis paid to sequences that a feature activates most strongly (Bills et al., 2023; Bricken et al., 2023; Cunningham et al., 2023; Templeton et al., 2024; Rajamanoharan et al., 2024a; Gao et al., 2024; Rajamanoharan et al., 2024b). The rating of a feature’s interpretability is usually either done by human raters (Bricken et al., 2023) or automatically by prompting a language model to do the evaluation (Bills et al., 2023).

4.3.3 Technical Advancements

SAEs have emerged as a dominant approach for conducting open-ended feature studies (Section 5.1.2) and have garnered significant attention within the MI community. The technical advancements in SAEs can be broadly categorized into two categories: (1) Balancing the trade-off between reconstruction and sparsity loss; (2) Discovering features that are functionally useful for the LM to implement downstream applications

Balancing the Trade-off between Reconstruction and Sparsity Loss The loss function of SAEs consists of two objectives: *reconstruction loss*, i.e., the sparse representation should accurately preserve information from input activation, and *sparsity loss*, i.e., only a small number of elements in the sparse representation should be active for any given input activation. These two objectives can conflict, as greater sparsity can reduce reconstruction loss (Gurnee, 2024). For instance, Wright & Sharkey (2024) noted that L_1 penalty in sparsity loss can lead to *feature suppression*, a systematic underestimation of feature activation magnitudes, particularly those with weak activations but high frequency, negatively impacting reconstruction loss. To address this, a line of research has focused on improving SAE architectures and training methods to better balance the trade-off between sparsity and reconstruction loss (Gao et al., 2024; Rajamanoharan et al., 2024b; Erichson et al., 2019). For instance, Gao et al. (2024) proposed *TopK SAE* that enforces sparsity by selecting only the K most active features and setting the rest to zero, where the sparse representation is obtained by $f(h) = \text{TopK}(g(W_{enc}h + b_{enc}))$. This approach eliminates the need for an L_1 penalty to achieve sparsity, mitigating the issue of feature suppression. Similarly, Rajamanoharan et al. (2024b) introduced *JumpReLU SAEs*, which leverage the JumpReLU function (Erichson et al., 2019) in place of the ReLU function. The JumpReLU function sets activations below a positive threshold to zero, effectively removing false positives and increasing sparsity, while leaving activations above the threshold intact, thus preventing feature suppression and improving fidelity. Finally, *Gated SAEs* (Rajamanoharan et al., 2024a) address the issue of feature suppression by using a gated ReLU encoder that decouples the detection of active features from the estimation of their magnitudes, and applying the L_1 penalty only to the feature detection. This decoupling allows the Gated SAE to achieve sparsity without underestimating feature magnitudes.

Discovering Functionally Important Features The primary goal of training SAEs is to discover functionally-important features used by LMs to implement their behavior. Specifically, a feature is considered functionally important if it facilitates the explanation of the model behavior on the training distribution. However, standard SAEs trained to minimize the reconstruction and sparsity loss may not strongly correlate with the goal of discovering functionally important features (Braun et al., 2024; Makelov et al., 2024). To address this, Braun et al. (2024) proposed *end-to-end (e2e) SAEs*, which are trained to minimize the KL divergence between the output distributions of the original model and the model with SAE-generated activations inserted, instead of minimizing the reconstruction loss. This approach ensures that e2e SAEs are optimized to identify features that influence the model’s predictions, rather than focusing on features that only accurately reconstruct activations.

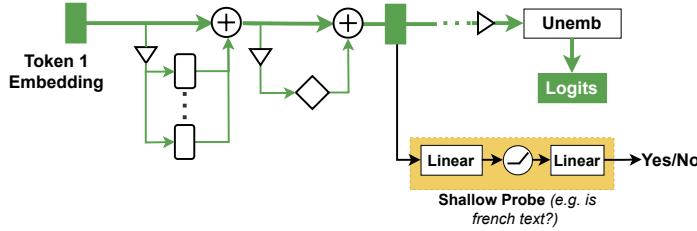


Figure 6: Probing on RS to detect whether it encodes a “French text” feature.

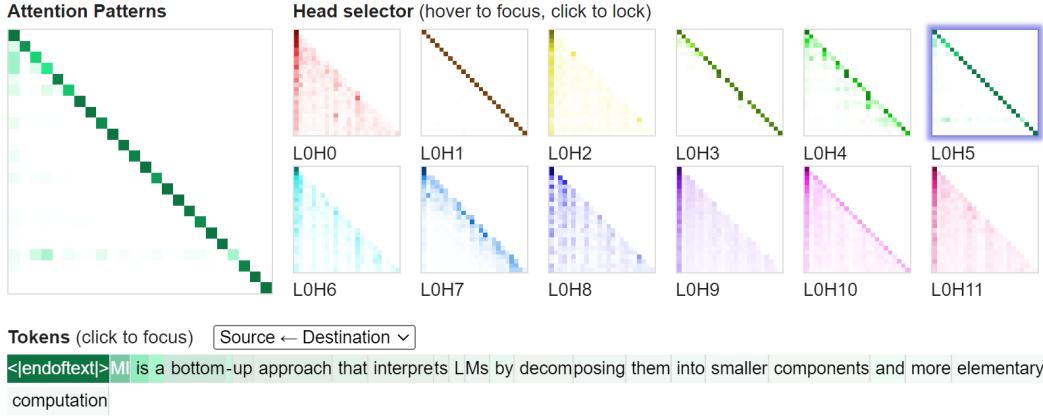


Figure 7: Attention visualization, created using the tool of Cooney & Nanda (2023).

4.4 Others

In addition to the three methods we present above, i.e., vocabulary projection, intervention, and SAEs, two other techniques are also widely used to aid the mechanistic interpretation of LMs. Both techniques have a long-standing history in the broad research topic of interpretable AI and machine learning. Below, we provide a brief introduction to them and their applications to interpreting LMs.

4.4.1 Probing

The probing technique (Figure 6) was developed and used extensively prior to the introduction of the MI field for investigating whether model activations encode various linguistic features such as part-of-speech (Conneau et al., 2018; Tenney et al., 2019a;b; Antverg & Belinkov, 2021). MI studies have also adopted probing as an important tool for investigating whether a pre-defined feature is present in intermediate activations. Specifically, it involves training a shallow (or linear) classifier, known as a *probe*, to predict whether a feature is present in those activations (Gurnee et al., 2023; Nikankin et al., 2024; Nanda et al., 2023b). However, it is important to note that the results of probing analyses only indicate a correlation, not a causal relation, between the feature and activations. We refer readers to Belinkov (2022) for further study on various types of probing.

4.4.2 Visualization

Visualization (Figure 7) is employed across various stages of an MI investigation, from generating initial hypotheses to refining them, conducting qualitative analyses, and validating results. For instance, attention patterns are often visualized to understand attention heads (Lieberum et al., 2023; Olsson et al., 2022), and a neuron activation across the input text is visualized to identify its functionality (Elhage et al., 2022a; Bricken et al., 2023). While visualization can be highly useful, it requires human effort to interpret results and carries

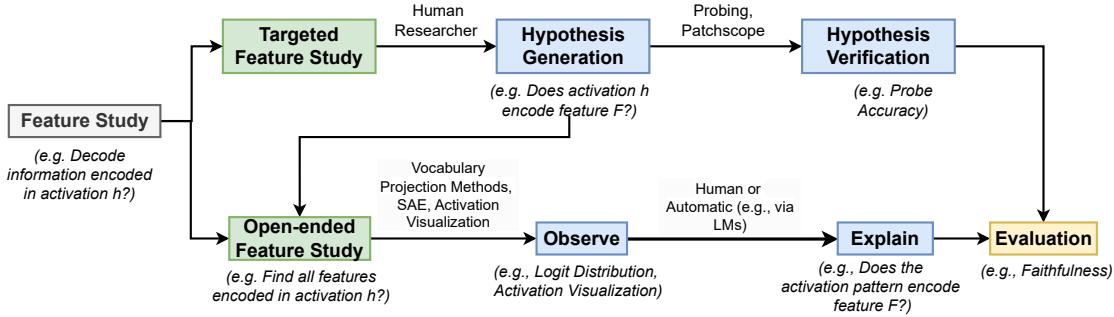


Figure 8: A task-centric beginner’s roadmap for feature study.

the risks of overgeneralization. Thus, any claims need to be substantiated with further experimentation and analysis.

5 A Beginner’s Roadmap to MI

A key motivation for this survey is to provide a friendly guide for researchers and developers interested in MI to quickly pick up the field. To this end, we provide a *beginner’s roadmap* in Figure 8, 11, and 13, where we categorize MI research into three categories – feature study, circuit study, and the study of universality, corresponding to the three objects of studies (Section 3.1). Each category is further divided into sub-categories with their corresponding MI workflow and associated techniques.

5.1 Feature Study

Feature studies can be broadly categorized into *targeted* and *open-ended* feature studies.

Targeted Feature Study A targeted feature study aims to determine whether a *pre-defined feature* is encoded in the representation of an LM. These pre-defined features are typically intuitive guesses or hypotheses made by humans, and we perform a targeted feature study when we want to verify the presence or absence of these features. Therefore, a general research question in a targeted feature study is, “*Does this LM representation (or a subset of its neurons) encode the feature F?*”

Open-ended Feature Study While targeted feature discovery is useful for determining whether a specific pre-defined feature is present in an activation, it relies on human intuition to hypothesize which features might be encoded. This could be problematic, as there may be features encoded within LM activations that do not align with human intuition. Consequently, targeted feature discovery alone is insufficient for fully interpreting an LM’s activations. To address this issue, open-ended feature discovery aims to enumerate all the features encoded within the activations of an LM. A general research question in an open-ended feature study is thus, “*What are all the features encoded in this LM representation?*”

5.1.1 General Workflow for Targeted Feature Study

Hypothesis Generation A targeted feature study begins with a hypothesis about the presence or absence of a pre-defined feature F in an LM representation. For example, suppose we are interested to know whether an LM activation encodes the “`is_python_code`” feature in its activations when the input contains a Python code snippet. Then, our hypothesis would be: “*Does this LM activation (e.g., h_i^1) encode the ‘is_python_code’ feature?*” We can refine our hypothesis to localize the feature encoding to a more granular level: “*Does this neuron in this activation (e.g., Neuron 1 in h_t^1) encode “is_python_code” feature?*”. Note that we can iteratively pose this question to all the LM activations when we want to identify a complete set of activations that encode the targeted feature.

Hypothesis Verification To verify whether or not a pre-defined feature exists in a given LM representation, there are typically two workflows to consider. The first workflow directly probes the LM representation and inquires the presence of the targeted feature F , whereas the second workflow follows a similar procedure as the open-ended feature study, which first identifies all features encoded in the representation and then confirms whether the feature F is among them. Below, we mainly present the first workflow. Details of the second workflow will be discussed in Section 5.1.2. For the first workflow, one option is to use the *Probing* method (Section 4.4.1) to verify the hypothesis, where a probe determining the existence of the targeted feature F will be trained. Alternatively, *PatchScope* (Ghandeharioun et al., 2024) is a flexible framework integrating vocabulary projection and intervention, and can also be used for verifying the above hypothesis. This method is especially useful to verify whether certain attributes associated with a given subject (e.g., attributes associated with the “largest city of” from the representation of “United States”) are encoded in the LM representation. In comparison to Probing, PatchScope enjoys the benefit of needing *no* additional training. However, it relies on humans to craft an expressive inspection prompt capable of eliciting the desired feature, which can be particularly challenging for abstract features. Therefore, the choice of these two techniques often depends on the type of feature being investigated and the availability of computing and data resources.

5.1.2 General Workflow for Open-ended Feature Study

Open-ended feature study typically involves two steps: (1) **Observe**: generate an *intermediate explanation* (e.g., logit distribution, activation visualization) for each activation; and (2) **Explain**: interpret the intermediate explanation to discover features encoded within the activations.

Observe Vocabulary projection methods (Section 4.1), neuron activation visualization (Section 4.4.2), and SAE (Section 4.3) are commonly used techniques for generating intermediate explanations. They are referred to as “intermediate explanations” because they require further interpretation during the Explain step to uncover features.

The vocabulary projection method generates a logit distribution as the intermediate explanation. As we elaborated in Section 4.1, one can then infer the encoded features by examining tokens with the highest logits from the distribution. However, it is important to note that the vocabulary projection method can only discover features in the prediction space. In other words, it only discovers features that are directly used for the next token prediction. For example, while completing the following input “*English translation for nourriture is*” with “*food*”, the LM might have extracted the *is_french_text* feature from the input and encoded it in its activations. However, the vocabulary projection method is unlikely to reveal the presence of the *is_french_text* feature. Consequently, intermediate explanations from vocabulary projection methods are not sufficient for discovering all features encoded in a given activation.

Neuron activation visualization (Sec 4.4.2) is another technique that involves developing an interface to visually highlight the text or tokens which elicit activity of a specific neuron, to create an intermediate explanation, as illustrated in Figure 9. The highlighted text snippets are typically sampled from a large corpus, focusing on instances where the neuron exhibits the highest activation. The idea here is that, for example, if a neuron consistently activates for text written in French, one can then label the neuron as encoding a French text detection feature. It is important to note that neuron activation visualization is based on an important assumption that features are represented by neurons in a one-to-one correspondence. However, subsequent studies (Elhage et al., 2022a; Marks et al., 2024b) have found *polysemantic neurons*, i.e., neurons that activate for multiple unrelated features, rendering the assumption of one-to-one correspondence to be false. The discovery of polysemantic neurons also complicates the Explain step, since polysemantic neurons can easily confuse the explainer when they activate for multiple unrelated features simultaneously.

To address this challenge, subsequent studies have proposed training SAEs (Section 4.3) that will transform the n -dimensional original activation into a higher k -dimensional sparse activation. The goal of this projection is to generate a new representation that faithfully retains all the features from the original activation but encodes them in a way that each neuron corresponds to a single feature, resulting in a representation composed of *monosemantic neurons*. The neuron activation visualization can then be applied to the sparse activation, rather than the original activation, for generating more interpretable intermediate explanations.

BASE64 NEURON (IN ONE-LAYER MODEL)

Dataset Examples

| |
|---|
| <EOT>8efMXtnduFiZVfa rzKbO9NszTKzaVMq51cVSAE0otA3QRy/I0zIU+bPRNineP195g98bBZx5TC5 <u>Bd</u> gSEQQkE61KQ 0x+azVlqduRYacXVphybbLrebP0a6otoXxCQAIkBYGaCQTyuP1Os5MfROdC1z+YerfRSLRISpN QiAcBGgahaHbQuQRJXngJvrQi/6sCotYpXW5way2febRccxxRGoiUAY84ii21Kskuh9ZEoIgwc Y8W4Bx98MO8s+3ChQvd0g2kVRK+zs4yD/4r7T/96U/dwleVyrVSOivw8Re1zlkzx3bs2BG12pbQ |
| "I'm not going to hold back." Brantley takes important step in comeback: https://t.co/O88CkxFQhw pic.twitter.com/UAYex7cPqk — Jordan Bastian (@MLBastian) March 17, 2016 |
| Clear acrylic pipe for filtration: https://amzn.to/30kNQGs Clear acrylic pipe products: https://amzn.to/30hJjB Protein Skimmers of choice: https://amzn.to/2LDu1Xn Finnex LED Aquarium light: https://amzn.to/2wOlbie |
| How is "Economic security for all who are unable or unwilling to work" vague? Progressives are pretty clear about what "economic security" means to them. https://t.co/CaHliQxszw pic.twitter.com/GfxG7ZK8D4 — Jeryl Bier (@JerylBier) February 9, 2019 |

Logit Weights

| | |
|--------------|-----------------------|
| +0.17 'Wl' | -0.60 ' section' |
| +0.15 'zA' | -0.60 ' segment' |
| +0.14 'même' | -0.60 ' of' |
| +0.14 'Rp' | -0.61 ' exam' |
| +0.13 'oi' | -0.61 ' hatch' |
| +0.13 'Cc' | -0.61 ' dusk' |
| +0.13 'Tx' | -0.61 ' eye' |
| +0.13 'Hg' | -0.62 ' count' |
| +0.13 'Hz' | -0.62 ' time' |
| +0.12 'mV' | -0.62 'levance' |
| +0.12 'zx' | -0.62 ' cent' |
| +0.12 'bw' | -0.63 'circumference' |
| +0.12 'GQ' | -0.63 ' balances' |
| +0.12 'Cb' | -0.63 ' operand' |
| +0.12 'Nz' | -0.64 ' volumes' |
| +0.12 'ru' | -0.64 ' quadrant' |
| +0.12 'fq' | -0.65 ' surface' |
| +0.12 '+/' | -0.66 ' volume' |
| +0.12 'Dt' | -0.70 ' end' |
| +0.12 'Jy' | -0.72 ' compartment' |

Figure 9: The neuron activation visualization developed by Elhage et al. (2022a) for labeling neurons with a feature by analyzing the text on which they activate.

Research on SAEs is still in its early stages, but it has shown promising results by discovering features that were not identified through neuron activation visualizations (Bricken et al., 2023; Templeton et al., 2024), which we will summarize in Section 7.1.

Although each of these techniques – vocabulary projection methods, neuron activation visualization, or SAEs, can be used individually to generate the intermediate explanation in the Observe step, in practice they are often combined to provide a more complete explanation of the LM activation. A user interface incorporating intermediate explanations from multiple techniques is shown in Figure 10. This may also include other explanations, such as *activation density*, which indicates how sparsely a neuron or feature activates for a given dataset. The availability of all these explanations can help the explainer more confidently identify and explain the features encoded in the corresponding activations.

Explain The intermediate explanation generated in the Observe step is then interpreted by explainers to uncover features encoded in the representation. For instance, when a vocabulary projection method is used in the Observe step, a human explainer views the top tokens in the projected logit distribution and labels the feature implied by the tokens (e.g., if top tokens are all associated with arithmetic addition, then a human evaluator infers that the activation is encoding features relevant to additive operation in arithmetic [Rai & Yao 2024]). Similarly, for neuron and SAE activation visualization, a human explainer is instructed to examine the activation visualization for each neuron, and then indicate whether they have found a plausible theory to explain the activations. For instance, Elhage et al. (2022a) provided the following instructions to human explainer – “mark INTERPRETABLE if 80% or more of the strongest firings can be explained by a single rule or category (e.g. the word ‘apple’, or any phrase relating to music), and NOT INTERPRETABLE otherwise”. While the explainers are often humans, we can also automate this Explain step using machines, such as *LLMs as explainers*. This involves providing LLMs with intermediate explanations and prompting them to label the corresponding features, if applicable. For instance, Bills et al. (2023) provided intermediate explanations from both the vocabulary projection method and the neuron activation visualization to a GPT-4 model for feature labeling; the authors reported a strong correlation between human and GPT-4-generated explanations, supporting the feasibility of leveraging LLMs for feature explanation.

5.1.3 Evaluation of Feature Study

The evaluation of a feature study is typically performed in two dimensions: *faithfulness* and *interpretability*. Specifically, faithfulness measures whether the discovered feature truly exists in the LM representation, while

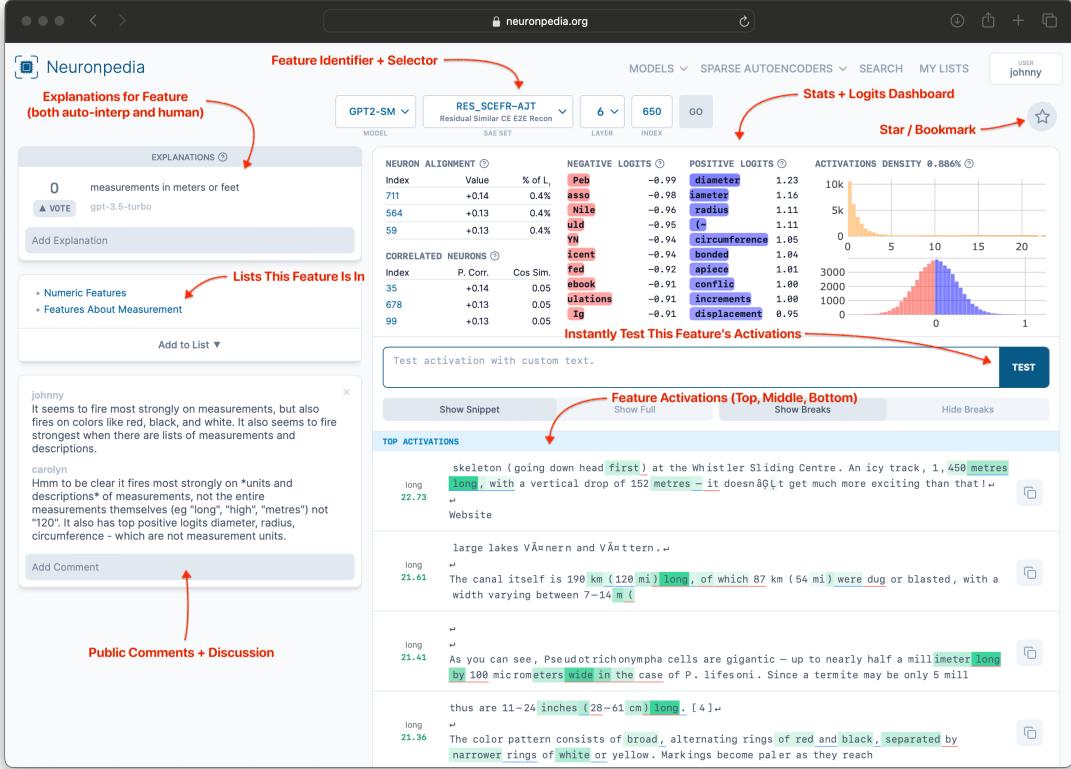


Figure 10: The interface prepared for the open-ended feature discovery using Neuronpedia (Lin, 2023). A human labeler or explainer uses this interface to label or explain the feature encoded by a neuron, based on the information presented within the interface.

interpretability assesses how easily a human can understand the feature description. Measuring faithfulness is inherently challenging due to the absence of ground truth. To this end, most studies often rely on technique-specific proxy measures, such as probing accuracy on a held-out test set for targeted feature studies and reconstruction loss when using SAEs for open-ended feature studies, where a higher accuracy or a lower reconstruction loss indicates greater faithfulness. In addition, some studies (Templeton et al., 2024; Marks et al., 2024b) also proposed manually altering the feature value, referred to as *activation steering*, during model inference to measure the direct effect of a feature on the next-token distribution of the model and causally measure the faithfulness of the feature study. On the other hand, interpretability of discovered features is typically conducted manually by a human or automatically by an LM. For instance, in the work of Bricken et al. (2023), a human is instructed to provide a score based on how interpretable the discovered features are, following a guideline “on a scale of 0–3, rate your confidence in this interpretation”. Alternatively, one can also use LMs to calculate an automatic explanation score (Bills et al., 2023) as a proxy for manual human evaluation.

5.2 Circuit Study

Similarly, the study of circuits (Figure 11) can be broadly divided into two categories, i.e., interpreting *an LM behavior* and interpreting *an LM component*.

Interpreting an LM Behavior This category of study involves identifying the circuit responsible for a specific LM behavior. For example, Elhage et al. (2021) investigated the circuits to explain the two-layer toy LM’s capability in sequence completion tasks and Wang et al. (2022a) discovered the circuit in GPT-2 in the task of indirect object identification. The two key questions in the circuit include: (1) *Localization*:

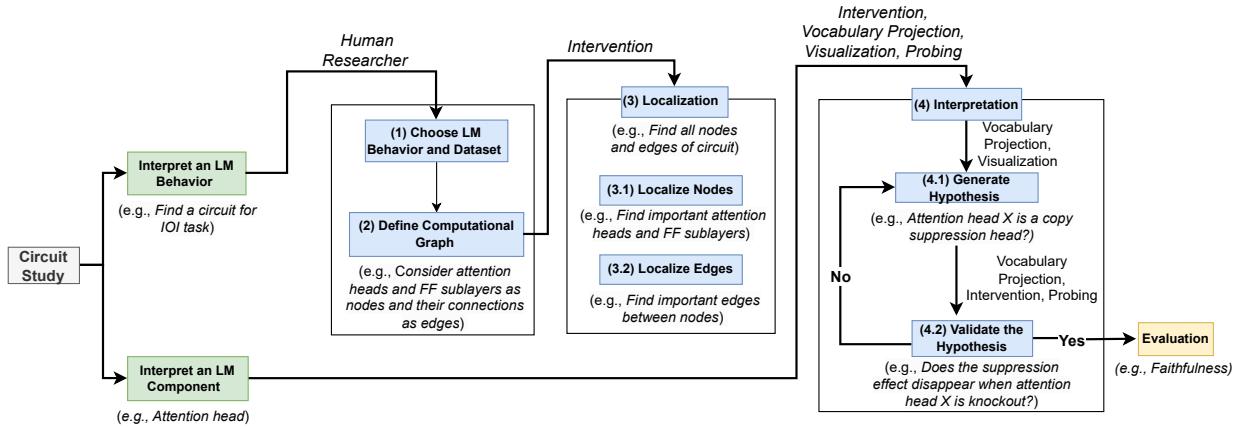


Figure 11: A task-centric beginner’s roadmap for circuit study.

Which LM components are responsible for this behavior? and (2) *Interpretation: How do these components implement the behavior?* Answering these two questions provides an algorithmic-level explanation of how the LM implements the observed behavior. Additionally, it is crucial to focus on behaviors where the LM performs with high accuracy or exhibits the capability clearly, as we can only investigate mechanisms that are already reliably present in the model. In Section 8, we will provide a more careful discussion about the feasibility of circuit studies.

Interpreting an LM Component This second category of circuit study aims to develop a thorough understanding of a specific LM component (e.g., attention head or FF sub-layer). The goal is often to develop a general understanding of the LM component that is independent of any specific task, making it essential for its interpretation to remain consistent across different tasks (Section 5.3). The process of interpreting an LM component overlaps with the “component interpretation” step of interpreting an LM behavior; however, to obtain a comprehensive interpretation of an LM component, one often needs to iterate the interpretation process across multiple tasks, which goes beyond the scope of a particular behavior.

Considering the overlap in the procedures for the two sub-categories of circuit study, in what follows, we will mainly focus on presenting the workflow for interpreting an LM behavior.

5.2.1 General Workflow for Circuit Study

A circuit study typically involves the following steps: (1) **Choose an LM behavior and a dataset:** Select the LM behavior for which we are searching the circuit. (2) **Define the computational graph:** Describe the nodes and edges of the computational graph. (3) **Localization:** Identify all the important nodes and edges connecting them. (4) **Interpretation:** Explain the role of all nodes and edges in implementing the LM behavior. (5) **Evaluation:** Evaluate the faithfulness of the discovered circuit.

Choose an LM Behavior and a Dataset The first step in circuit discovery involves selecting a specific LM behavior and curating a dataset that showcases the behavior. The LM should have high task performance on the dataset, as this suggests the presence of internal mechanisms or a circuit that supports the behavior. In contrast, low task performance indicates that the mechanism may be poorly defined or absent, preventing further investigation.

Define the Computational Graph In most existing literature (Conmy et al., 2023; Syed et al., 2023; Hanna et al., 2024a;b), a circuit is defined as a directed acyclic graph (DAG) connecting individual computational components in the LM. For example, Hanna et al. (2024b) defined the GPT-2 model as DAG, where 144 attention heads (12 heads/layer × 12 layers) and the 12 FF sub-layers (1 per layer) in a GPT-2 model were considered nodes in its DAG and discovered circuits as sub-graphs from it. Following Elhage

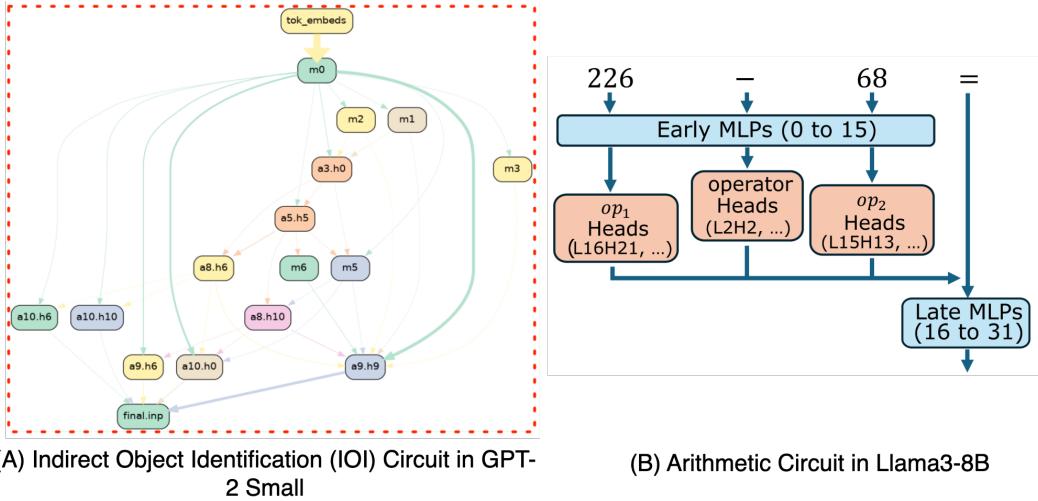


Figure 12: (A) The IOI circuit, discovered by Conmy et al. (2023) as a reproduction of Wang et al. (2022a), where the computational graph of the model consists of nodes that are not directly connected across adjacent layers but still connected in effect by writing to and reading from the RS. (B) The arithmetic circuit, discovered by Nikankin et al. (2024), is represented as a computational graph where multiple nodes can be assigned to the same LM component, one per input position, when positional information suggests distinct functional roles for the same component.

et al. (2021), many works also considered the connections between LM components in *non-adjacent* layers as valid edges in the LM’s DAG, due to the additive nature of the RS (Eq 1). That is, even if two components (e.g., the FF sub-layers in Layer 2 and Layer 5) are not connected through direct computations across adjacent layers, they are still connected in effect by writing to and reading from the RS, as shown in Figure 12. In practice, one can define a circuit at various levels of abstraction and granularity depending on the scope of interpretation. For example, instead of considering each LM component as a single node, one can define multiple nodes of the same LM component, one for each input position, when the position information may imply different roles of the same component (Wang et al., 2022a; Nikankin et al., 2024), as shown in Figure 12.

Localization Once the computational graph for the LM is defined, the next step is to identify all the important nodes and edges responsible for implementing the behavior. This process is termed “localization”. Localization can be broken down into *Localization of Nodes* and *Localization of Edges*. **(1) Localization of nodes:** Intervention methods, such as zero ablation (Olsson et al., 2022) and mean ablation (Wang et al., 2022a), are commonly used to measure the importance of individual nodes in the model. These methods involve intervening in a node’s output and observing its impact on the model’s behavior, with important nodes being those whose intervention results in a deviation above a certain threshold in model performance. **(2) Localization of edges:** Path patching (Wang et al., 2022a) is a commonly used intervention method to localize edges. Other intervention-based techniques, such as zero ablation, mean ablation, and resampling ablation, can also be applied for localization. However, these methods can be computationally-intensive as we need to measure the importance of each edge individually. To address this issue, automated localization techniques such as ACDC (Conmy et al., 2023) and EAP (Syed et al., 2023), have been proposed, as discussed in Section 4.2.3.

Interpretation Once all the important LM components and edges have been identified during the localization step, we interpret the functional role of each component. This interpretation step can be further divided into two sub-steps – **(a) Generating a hypothesis:** The first step involves generating a hypothesis about the function of each LM component. This process is often guided by analyzing the behavior of the component under relevant inputs using various techniques, such as visualization (e.g., attention visualization) and vocabulary projection methods. A human can then generate a plausible hypothesis based on the

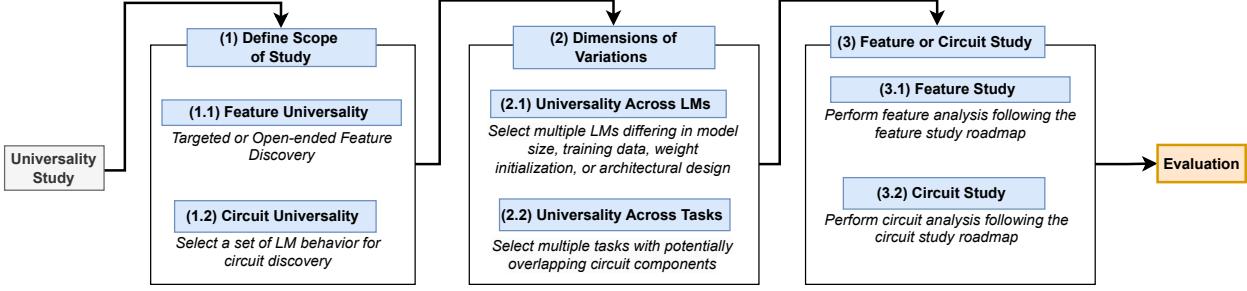


Figure 13: A task-centric beginner’s roadmap for the study of universality.

analysis (e.g., attention head H was observed to decrease the logit of the token it attends to and thus is hypothesized to function as a copy suppression head). **(b) Validating the hypothesis:** The hypothesis can be validated using various techniques, including both vocabulary projection and intervention-based methods. For instance, if the hypothesis suggests that attention head H is a copy suppression head, then knocking out (or noising) attention head H should remove the suppression effect. If the hypothesis is validated, the model component is considered interpreted; otherwise, further analysis and new hypotheses are needed.

Evaluation The discovered circuits are evaluated for three criteria – faithfulness, minimality, and completeness. *Faithfulness* is evaluated by comparing the full model vs. the partial one with the localized circuit alone (Olsson et al., 2022; Wang et al., 2022a; Marks et al., 2024b); the circuit is considered faithful when the performance gap between the full and the partial LM is small. On the other hand, *minimality* measures whether all LM components in the circuits are necessary, often by randomly ablating components in the circuit and computing the change in the LM behavior (Wang et al., 2022a). A significant change in behavior would indicate that all parts of the circuit are necessary and the explanation is minimal. Minimality is also referred to as *sparsity* (Bhaskar et al., 2024). Finally, *completeness* checks if the circuit includes all the LM components used by the LM for exhibiting the behavior of interest (Wang et al., 2022a; Marks et al., 2024b). Specifically, Wang et al. (2022a) measures the completeness by comparing the full model vs identified circuit behavior under random ablations of the circuit components. If the circuit is complete, LM behavior should remain similar to the whole model even under random ablations.

5.3 Study of Universality

The study of universality (Figure 13) investigates whether similar features and circuits can be found across different LMs or tasks. To this end, the study of universality involves investigating the internal mechanisms of several LMs and tasks with a primary research question as – “*Do the same features and circuits exist across LMs and tasks?*”

5.3.1 General Workflow for Universality Study

The general workflow for universality consists of the following steps.

Scope of Universality The study of universality can be mainly studied under two dimensions – (1) *Universality of features*: the degree of similarity of the features contained within representations across different LMs; (2) *Universality of circuits*: whether a particular LM behavior is implemented using similar circuits across LMs, and whether individual model components perform the same specialized functions (e.g., induction heads [Olsson et al. 2022], successor heads [Gould et al. 2023]) across multiple LMs and tasks. The universality of features can be studied in both the targeted and open-ended feature study settings. On the other hand, the study of circuit universality only makes sense when a particular LM behavior is consistently observed across multiple LMs and/or tasks, as the goal is to examine whether various LMs or a single LM in various tasks implement the same behavior using similar circuits.

Dimension of Variations To measure **universality across LMs**, it is essential to analyze multiple LMs that differ in aspects such as model size, training data, weight initialization, or architectural design. For instance, Gurnee et al. (2024) trained five GPT2-Small and GPT2-Medium models from different random seeds to study the universality of neurons. In the case of studying **universality across tasks**, the dimension of variation can be the task choice and the specific task setting (e.g., data distribution). For instance, Merullo et al. (2023) studied two different tasks, an indirect object identification task and a colored objects task, to determine whether certain LM components are reused across tasks while maintaining the same function.

Feature and Circuit Study The feature and circuit study can then be performed across LMs and/or tasks, as discussed in Sections 5.1.1 and 5.2.1.

Evaluation Finally, we compare the similarities and differences of the discovered circuits and features across different LMs and/or tasks. Specifically, the features and circuits that are found across various LMs are considered to be universal, while those unique to specific LMs are considered model-specific. Likewise, one can draw conclusions about task-agnostic vs. task-specific circuits. It is important to note that the question of “do similar features and circuits exist across LMs and tasks” may not be binary; the experiments in practice instead reflect the degree to which universality holds. For instance, Gurnee et al. (2023) found that only 1-5% of the neurons were universal across five different GPT2-Small and GPT2-Medium models.

6 Case Studies of Beginner’s Roadmap

We present case studies for research work in MI, detailing how the high-level workflows in our roadmap were employed in papers that conducted these analyses. These case studies serve as *concrete examples* to help readers new to the field effectively map research questions to specific research activities outlined in our roadmap.

6.1 Case Study for Feature Study

6.1.1 Targeted Feature Study with Probing

Gurnee et al. (2023) proposed and employed a probing technique, *sparse probing*, to conduct targeted feature discovery. **(1) Hypothesis Generation:** Gurnee et al. (2023) sought to discover 100 pre-defined features across seven LMs of varying sizes. Furthermore, they also aimed to localize features to a specific neuron or sets of neurons within the activations. For instance, to determine whether an LM encodes the “is_french_text” feature in the FF activations in the first layer (f_i^1 ’s) of the LM, the following steps were taken for **(2) Verifying the Hypothesis:** **(2.1) Dataset preparation:** A labeled dataset consisting of a train and test set was created for training and evaluating the probe. This dataset included examples with the “is_french_text” feature (i.e., French text) and examples without this feature. **(2.2) Activation Extraction:** The corresponding f_i^1 activations for each example in the dataset were obtained when they were used as input text of the LM. Whenever the examples consisted of multiple tokens, elementwise aggregation (e.g., mean or max) of the activations over the token span were considered. **(2.3) Training a Sparse Probe:** A sparse probe was trained using the extracted activation a_i^1 and their corresponding labels in the train set. To localize the feature to a specific set of neurons, the activations were constrained to have at most k -non-zero coefficients. To determine which neurons should have non-zero coefficients, a technique called *adaptive thresholding* was employed which involves training a series of classifiers that iteratively decrease the value of k . At each step, the probe is retrained to use only the top- k neurons with the highest coefficient magnitudes from the previous step. **(3) Evaluation:** Finally, the trained sparse probe was evaluated using F1 score, where if the probe achieves a high F1 score in the test set, it suggests that the neurons with non-zero coefficients have “is_french_text” feature encoded within them. Furthermore, they also measure the precision and recall to get more insight into the granularity of features encoded by the set of neurons. For instance, high precision and low recall of the probing classifier may indicate that the identified subset of neurons may encode a more specific feature than the feature being probed for (e.g. “is_french_noun” instead of “is_french_text”). In their work, Gurnee et al. (2023) repeated Steps 2-3 multiple times to de-

determine whether activations at other layers also encode the “is_french_text” feature, and applied the same to multiple LMs to study the universality of this feature.

6.1.2 Open-ended Feature Study with SAEs

Bricken et al. (2023) conducted open-ended feature discovery on a one-layer toy LM using SAE, visualization, and human and automatic interpretation techniques. The following steps were taken to discover features encoded in the LM activations: **(1) Observe:** SAEs were first trained on the activations obtained by feeding the LM sentences from its training set as input. After training, a higher-dimensional, sparse activation was obtained for each LM activation. The next step was to accurately label each element of the sparse activations, referred to as *features* by Bricken et al. (2023), with a correct description. For instance, a feature might activate in response to phrases related to music. A user interface was then developed for each element of the sparse activations to provide useful information for accurate labeling. This included text examples of when each sparse feature activates, the effects they have on the logits when active, examples of how they influence token probabilities if the feature is ablated, and other relevant details. All of these observations provide the intermediate explanations of each activation. **(2-3) Explain and Evaluation:** Human evaluators were tasked with labeling the features using all the information provided in the interface, guided by a scoring rubric that included instructions such as “On a scale of 0–3, rate your confidence in this interpretation” for evaluating the interpretability of the feature discovery. Additionally, automatic explanations (Bills et al., 2023) were also used.

6.2 Case Study for Circuit Study

Wang et al. (2022a) conducted one of the first circuit analyses based on GPT-2 Small, following these steps. **(1) Choose an LM Behavior and a Dataset:** The study examined a circuit in GPT-2 Small responsible for solving the indirect object identification (IOI) task: Given a sentence such as “When Mary and John went to the store, John gave a drink to”, complete the sentence with “Mary” (IO) and not “John” (S). To investigate the behavior, a dataset was synthesized to showcase the IOI task. **(2) Define the Computational Graph:** The study defined the computational graph with attention heads as nodes and connections between them as edges. Notably, components such as FFs, LayerNorms, and embedding matrices were excluded from the analysis, allowing the investigation to focus solely on understanding the attention heads. In addition, the study considered the same attention heads at different token positions as distinct components; as a result, the computational graph to be discovered is *position dependent*. **(3) Localization:** Using the synthesized dataset, path patching was applied to identify the important nodes and edges in the circuit. Specifically, to determine whether an attention head a_i^l is important, path patching for the pathway $a_i^l \rightarrow \text{logits}$ was applied and the effect on the logit difference between the direct (John) and indirect (Mary) object was measured, where a high logit difference indicates that the a_i^l plays a crucial role in the IOI task. After performing path patching across all attention heads in GPT-2 Small, 26 attention heads were identified as important nodes for the IOI circuit. **(4) Interpretation:** This step involves investigating the role of each of the 26 important attention heads identified in the previous step. For instance, to understand the role of a_9^9 , the following two substeps were performed. **(4.1) Generating a hypothesis:** To generate a plausible hypothesis on the role of a_9^9 , attention visualization was analyzed, which showed that a_9^9 strongly attends to the indirect object token (e.g., “Mary”). Based on this observation, a hypothesis was formed: a_9^9 is a *name mover head* that attends to the correct name and copies whatever it attends to. **(4.2) Validating the hypothesis:** They validated the hypothesis that attention head a_9^9 is a *name mover head* by demonstrating that a_9^9 (1) attends to the correct name and (2) copies whatever name it attends to. First, they analyzed attention patterns and found that these heads strongly focused on the indirect object (IO) token, with higher attention correlating ($\rho > 0.81$) with increased logit output in the direction of the correct name token (IO token). To further confirm their copying behavior, they studied what values were written via the OV matrix of a_9^9 by simulating the perfect attention score of a_9^9 to the IO token and measuring the resulting logit probabilities. Specifically, they applied the logit lens to $h_{\text{name}}^1 W_V^{9,9} W_O^9$ (i.e., h_{name}^1 multiplied by the OV matrix of the head a_9^9) to examine what values these heads wrote into the RS at the position of each name token. The “copy score” that quantified how often the input name appeared in the top 5 logits, exceeded 95% for Name Mover Heads, compared to less than 20% for average heads. These results

confirm that Name Mover Heads play a crucial role in copying and transferring name information through the model’s computation. **(5) Evaluation:** The study assessed the discovered circuit using three evaluation metrics—faithfulness, minimality, and completeness—as discussed in Section 5.2.1.

7 Findings and Applications

7.1 Feature Study

Monosemantic vs. Polysemantic Neurons Earlier work studied neurons as a natural candidate for encoding features, where a high value at a particular position would indicate the presence of that feature. This led to the discovery of several interesting neurons such as *sentiment neurons* (Tigges et al., 2023), *Base64 neurons* (Elhage et al., 2022a), *French neurons* (Gurnee et al., 2023), *knowledge neurons* (Dai et al., 2022), *skill neurons* (Wang et al., 2022b), *positional neurons* (Voita et al., 2023), etc. The question of whether neurons encode specific features has been explored extensively, not only in the context of transformer-based LMs but also in earlier work on word embeddings (Mikolov et al., 2013), recurrent neural networks (RNNs; Radford et al., 2017), and vision models (Cammarata et al., 2020), predating the field of MI itself. However, many of these studies found that neurons within LMs are not *monosemantic*, i.e. they do not activate only for a single feature. In contrast, they are *polysemantic*, i.e. they activate in response to multiple unrelated features (Elhage et al., 2022a; Gurnee et al., 2023; Elhage et al., 2022b). For instance, a single neuron could activate for both French texts and texts encoded in Base64. This polysemantic nature of neurons shows the initial hypothesis that each neuron encodes a single feature is not always true.

Superposition The discovery of polysemantic neurons has led to the hypothesis of *superposition*, i.e., a model can represent a greater number of independent features than the number of available neurons (Elhage et al., 2022b). In this framework, features are encoded as directions formed by a linear combination of multiple neurons rather than being tied to a single neuron. Although superposition enables the encoding of a larger number of features, it also introduces potential interference among overlapping representations, where activating one feature may unintentionally activate others that share the same representation space. Empirical findings suggest that the natural sparsity of feature activations helps mitigate this issue of interference (Elhage et al., 2022b). Since only a small subset of features is typically active at any given time, the risk of unintended activations is reduced, allowing the neural model to maintain distinguishable representations despite the constraints imposed by superposition. Furthermore, models appear to prioritize encoding the most important features in more monosemantic fashion, while less critical features are either distributed across multiple neurons or omitted altogether, thereby minimizing performance degradation caused by interference (Scherlis et al., 2022). Beyond its role in feature representation, superposition has also been observed to aid in data memorization for smaller datasets and in learning more generalized features for larger ones, with a notable transition between these regimes marked by the phenomenon of double descent (Henighan et al., 2023).

Tackling Superposition with SAEs Recently, SAEs have gained popularity as a method for interpreting representations that encode features in superposition (Bricken et al., 2023; Sharkey et al., 2022; Riggs, 2023; Cunningham et al., 2023; Lieberum et al., 2024; Marks et al., 2024a; He et al., 2024). Various SAE variants with different architectures have been proposed, as discussed in Section 4.3. Early studies on SAEs have shown promising results, with the extracted features from SAEs appearing more interpretable than those from neurons, according to both human analysis and automatic explanation scores (Bricken et al., 2023; Makelov, 2024). Despite these promising findings, several challenges persist. For instance, SAEs do not perfectly reconstruct the activation which is often attributed to the trade-off between sparsity and reconstruction loss, raising concerns regarding the faithfulness of the extracted representations (Gurnee, 2024; Muhammed et al., 2024; Chanin et al., 2024; Anders & Bloom, 2024). For instance, Templeton et al. (2024) speculated that SAEs may have only captured a fraction of features encoded in model activations, with many rare or highly specific concepts remaining undetected due to their infrequent activation. Similarly, recent studies (Braun et al., 2024; Makelov et al., 2024) have also raised concerns regarding the functional usefulness of the extracted features, i.e., whether these features are actively used by the model for inference or merely reflect patterns present in the training data used to train the SAE. Although SAEs with improved

architectures have been proposed to mitigate some of these limitations (Muhamed et al., 2024; Braun et al., 2024), further advancements in SAEs are needed to ensure that the extracted features are both representative of the model’s underlying computation and practically useful for interpretability.

Linear Representation Hypothesis Most MI studies rely on the *linear representation hypothesis* (Mikolov et al., 2013; Pennington et al., 2014) for performing feature studies. The hypothesis states that *features are encoded linearly in the representation space of the model*. For example, SAEs adopt this hypothesis by assuming that activations can be expressed as sparse linear combinations of feature vectors, while probing assumes that a linear model can detect the presence or absence of a feature in activations. Specifically, the linear representation hypothesis posits that neural networks have two properties: *linearity*, i.e. the network’s activation space consists of meaningful (linear) vectors, each representing a feature, and *decomposability*, i.e., network activations can be decomposed and described in terms of these independent features. Notably, this hypothesis predates MI and has been extensively studied to understand how high-level concepts are stored in word embeddings (Mikolov et al., 2013), latent representations of variational autoencoders (Wang et al., 2024d), and LM representations (Dai et al., 2022; Radford et al., 2017). A well-known example of the hypothesis is the presence of semantic relationships such as $\overrightarrow{\text{man}} - \overrightarrow{\text{woman}} \approx \overrightarrow{\text{king}} - \overrightarrow{\text{queen}}$ within word embeddings, providing evidence that high-level concepts (e.g., gender) are represented as specific directions in the activation space. An increasing body of work has found evidence supporting the hypothesis by discovering features encoded as directions within LM representations such as sentiment (Tigges et al., 2023), space and time (Gurnee & Tegmark, 2023), language (Bricken et al., 2023), truth (Marks & Tegmark, 2023), and refusal (Arditi et al., 2024). However, some studies have also found the existence of non-linear features that contradict the linear representation hypothesis. For instance, Engels et al. (2024) discovered multi-dimensional features, such as circular representations for days of the week and months of the year, demonstrating that some features are inherently irreducible to a single dimension. This finding contradicts the previous notion of linear representation hypothesis where each feature is represented by one-dimensional directions. Subsequently, Olah (2023) pointed out that since multi-dimensional features remain mathematically linear, existing techniques should be capable of detecting them with a slight change. However, Csordás et al. (2024) found that gated recurrent neural networks (RNNs) encode sequential information using non-linear “onion representations”, where tokens are represented by magnitude rather than direction, leading to layered features that do not reside in distinct linear subspaces. Given these contradictory findings, further research is crucial to either validate or refute the linear representation hypothesis, as it forms the foundation of much feature analysis.

7.2 Circuit Study

7.2.1 Interpreting LM Behaviors

Circuit studies have identified a range of behaviors across LMs of different sizes, as listed in Table 3. These circuits are defined at varying levels of granularity, where nodes in the computation graph are defined as the outputs of MHA and FF sublayers (Olsson et al., 2022; Wang et al., 2022a) or as SAE features (Marks et al., 2024b; Cunningham et al., 2023; Kissane et al., 2024). In addition, Merullo et al. (2023) and Quirke et al. (2024) showed that the same components (e.g., induction heads) are reused by different circuits (e.g., IOI and induction circuits) to implement different tasks, demonstrating the generalizability of interpreted components. Beyond circuit discovery, some studies (Heimersheim & Janiak, 2023; Marks et al., 2024b) have shown that the identified circuits can be enhanced for better performance, highlighting the potential for practical downstream applications. For instance, Marks et al. (2024b) proposed a technique, Spurious Human-interpretable Feature Trimming (SHIFT), which improves the generalization of a classifier by removing the spurious features in the discovered circuit. This approach paves the way for new possibilities in human-AI collaboration.

However, these studies have also revealed new challenges in circuit discovery. For instance, circuit discovery methods were found to generalize poorly to adversarial examples, indicating that existing approaches are sensitive to the datasets used in circuit identification and evaluation (uit de Bos & Garriga-Alonso, 2024). In addition, some studies (Zhong et al., 2024; Nanda et al., 2023a) have also shown that LMs implement different algorithms in tandem to solve the same task and the interactions between these parallel algorithms

| Circuits | Circuit Descriptions | Circuit Nodes | LMs |
|---|--|--|-------------------------------|
| Induction (Elhage et al., 2021) | Completes sentences like “Mr D urs ley was thin and bold. Mr D” with “urs” | 2 attention heads | 2-layer toy LM, GPT2-Small |
| Indirect Object Identification (Wang et al., 2022a) | Completes sentences like “When John and Mary went to the store, John gave a drink to” with “Mary” as opposed to “John” | 26 attention heads | GPT2-Small |
| Docstring (Heimersheim & Janiak, 2023) | Predicts argument names in the docstring: <pre>def port(self, load, size, files): ... :param load: crime population :param size: unit dark :param</pre> is completed by “files” | 8 attention heads | 4-layer toy LM |
| Greater-Than Operation (Hanna et al., 2024a) | Completes sentences such as “The war lasted from the year 1732 to the year 17” predict valid two-digit end years greater than 32 | 8 attention heads and 4 FF sub-layers | GPT2-Small |
| Gender Bias (Chintam et al., 2023) | Completes sentences such as “The {profession} said that” with gendered pronouns reflecting common stereotypes | 21 attention heads and FF sublayers | GPT2-Small |
| Subject-verb Agreement (Marks et al., 2024b) | Completes sentences like “The keys in the cabinet” with “are” as opposed to “is” | < 100 SAE features | Pythia-70M |
| Arithmetic Calculation (Nikankin et al., 2024) | $1 + 4 \rightarrow 5; 8 - 4 \rightarrow 4; 2 * 4 \rightarrow 8; 8/4 \rightarrow 2$ | 6, 6, 20, and 6 attention heads, along with all FF sub-layers, correspond to Addition, Subtraction, Multiplication, and Division, respectively | Llama3-8B |

Table 3: List of circuits identified in various MI studies. Note that the list does not encompass all circuits discovered by the MI community.

are highly sensitive to training and architectural hyperparameters. Finally, scalability remains a concern for larger LMs. To demonstrate the scalability of current techniques for identifying circuits in LLMs, Lieberum et al. (2023) identified the circuit used for the multiple-choice question-answering task on the 70B Chinchilla LLM (Hoffmann et al., 2022).

7.2.2 Interpreting LM components

Interpreting Transformer Components The study of circuits has also yielded insights into the functionalities of transformer components.

The **RS** can be viewed as a one-way communication channel that transfers information from earlier to later layers. Furthermore, Elhage et al. (2021) hypothesized that MHA and FF in different layers write their output in different subspaces of the RS, which prevents interference of information. In addition, nostalgebraist (2020) proposed to view the RS as an LM’s current “guess” for the output, which is iteratively refined layer-by-layer.

The **MHA sublayers** are responsible for moving information between tokens, which enables information from other tokens (i.e., context) to be incorporated into each token’s representation. Elhage et al. (2021) showed that each attention head in a layer operates independently and can be interpreted independently. Several MI studies have shown that these attention heads seem to have specialized roles. For instance, “negative heads” discovered in GPT2-small by McDougall et al. (2023) are responsible for reducing the logit values of the tokens that have already appeared in the context. Other notably identified attention heads include previous token heads (Wang et al., 2022a), duplicate token heads (Wang et al., 2022a), copying heads (Elhage et al., 2021), induction heads (Olsson et al., 2022), and successor heads (Gould et al., 2023).

The **FF sublayers** are attributed for the majority of feature extraction (Gurnee et al., 2023), storing and recalling pre-trained knowledge (Meng et al., 2022) and arithmetic computation (Stolfo et al., 2023a). More generally, Geva et al. (2021) viewed FF sublayers as key-value stores where the outputs of the first layer (W_k^l) of the FF sublayer serves as keys that activate values (stored knowledge) within the weight matrices

of the second layer (W_v^l). Furthermore, they demonstrated that earlier FF layers typically process shallow (syntactic or grammatical) input patterns, while later layers focus more on semantic patterns (e.g., text related to TV shows).

The **LayerNorms** (Ba, 2016) are primarily employed to stabilize and accelerate the training of LMs. However, their role in model computation during inference remains less understood. Many MI studies assume that LayerNorm layers do not play a significant role during inference and thus often ignore their computation. To investigate this assumption, Heimersheim (2024) demonstrated that the LayerNorm in a pre-trained GPT2-small can be removed by fine-tuning the model on a small fraction of the training data without the LayerNorm, suggesting that it may not be essential for the core computational process. In contrast, other studies have revealed that LayerNorm can have additional functionalities; for instance, Winsor (2020) showed that LayerNorm can serve as a general-purpose non-linear function capable of performing various classification tasks, while Stolfo et al. (2024) proposed that it is instrumental in implementing confidence regularization within LMs. These findings highlight that while LayerNorm is critical during training, its exact contributions during inference merit further exploration.

7.3 Universality

We summarize the findings on the universality of features and circuits below.

7.3.1 Universality of Features

The extent of feature universality remains largely unexplored, with no definitive conclusions. Gurnee et al. (2024) found that only 1-5% of neurons in randomly initialized GPT-2 models exhibit universality, though the polysemantic nature of neurons may have obscured the identification of universal features. To address the problem, recent studies (Lan et al., 2024; Wang et al., 2024b) have used SAEs to analyze feature universality across LMs, revealing a high degree of similarity in SAE features across models of varying architectures, sizes, and training regimes. Investigating the degree of feature universality and its dependence on factors such as initialization, model size, and loss function remains a critical open challenge.

7.3.2 Universality of Circuits

Similar to feature universality, studies on circuit universality have yielded mixed results. Early circuit analyses identified components such as induction heads (Olsson et al., 2022), successor heads (Gould et al., 2023), and duplication heads (Wang et al., 2022a), across multiple LMs. Similarly, Merullo et al. (2023) found that different circuits implementing different tasks (IOI and colored objects tasks) reuse the same components (e.g., induction heads), suggesting a degree of universality of circuits across tasks studied. Additionally, Tigges et al. (2024) found that specialized attention heads responsible for tasks like indirect object identification (IOI) (Wang et al., 2022a), and greater-than task (Hanna et al., 2024a) consistently emerge while training models of different sizes after a similar number of tokens are processed during training. However, Zhong et al. (2024) discovered that two LMs trained with different initializations can develop qualitatively different circuits for the modular addition task. Similarly, Chughtai et al. (2023) found that LMs trained to perform group composition on finite groups with different random weight initializations on the same task do not develop similar representations and circuits.

7.4 Findings on Model Capabilities

In-Context Learning (ICL) ICL is an emergent ability of LLMs that enables them to adapt to new tasks based solely on instructions or a few demonstrations at inference time (Wei et al., 2022a). Elhage et al. (2021) studied a simplified case of ICL and discovered an *induction circuit* composed of attention heads with specialized roles (e.g., *induction heads*), which were then found to be crucial even for general cases of ICL (Olsson et al., 2022; Bansal et al., 2023). Following this, Bietti et al. (2024) examined how a transformer model balances memorizing general knowledge (*global bigrams*) with adapting to sequence-specific information (*in-context bigrams*). They observed that global bigrams are learned first, followed by the gradual development of an induction head mechanism for handling in-context bigrams. Similarly, Reddy

(2023) showed that the distributional properties of training data (e.g., items in natural data tend to occur in clusters rather than being uniformly distributed over time [Chan et al. 2022b]) are crucial for the development of induction heads. Further investigations in various synthetic settings have uncovered additional variants of induction heads. For instance, Edelman et al. (2024) investigated ICL in bigram models trained on samples from a Markov chain and demonstrated that transformers learn *statistical induction heads* when trained on such data. Ren et al. (2024) further investigated few-shot ICL and identified *semantic induction heads*, which, unlike prior induction heads, model the semantic relationship between the input and the output token (e.g., “*I have a nice pen for writing. The pen is nice to*” → “*write*”). Besides the study of induction heads, Hendel et al. (2023) demonstrated that an LM maps the set of demonstrations \mathcal{S} in ICL to a *task vector* that essentially represents the mapping/rule described in \mathcal{S} . This task vector is then used by the model to generate task-relevant outputs.

Reasoning Recently, sufficiently large LMs have been shown to exhibit reasoning capabilities (Wei et al., 2022a; Huang & Chang, 2022). MI studies have investigated LMs on various reasoning tasks to understand how and to what extent LMs perform reasoning. For instance, Stolfo et al. (2023b) studied arithmetic reasoning and found that attention heads are responsible for transferring information from operand and operator tokens to the RS of the answer or output token, with FF modules subsequently calculating the answer token. Rai & Yao (2024) employed a neuron activation analysis to study how Chain-of-Thought (CoT) prompting (Wei et al., 2022b) elicits the arithmetic reasoning capability of LMs. Their study reveals that specific FF neurons encode arithmetic reasoning concepts and that their activation is necessary but not sufficient for arithmetic reasoning. Besides arithmetic reasoning, Dutta et al. (2024) studied CoT multi-step reasoning over fictional ontologies and found that LMs seem to deploy multiple pathways in parallel to compute the final answer. Additionally, Brinkmann et al. (2024) discovered an interpretable algorithm in LM for the task of path-finding in trees. On the other hand, Saparov et al. (2024) showed that transformers learn to perform graph search/multi-hop reasoning using a parallel “path-merging” algorithm by extracting the circuit for this algorithm; so the embeddings of each fact will include information about facts provable in k hops, where k increases exponentially with the number of layers. Similarly, Biran et al. (2024) also investigated multi-hop reasoning on factual information by employing Patchscopes and showed that while the early layers effectively handle the first-hop reasoning, later layers often fail because they may no longer contain the necessary mechanisms to perform the second-hop reasoning required to extract the correct answer. Finally, Men et al. (2024) identified that the look-ahead planning mechanism in LMs is primarily facilitated by the multi-head attention in middle layers at the last token.

Knowledge Mechanisms LMs acquire a vast amount of knowledge (e.g., facts, grammar, commonsense, concepts, etc.) during pretraining; however, our understanding of how they store, recall, and utilize the knowledge remains incomplete (Wang et al., 2024c). Geva et al. (2021) posit that *knowledge is stored in the FF sub-layers*, where the first transformation layer of FF generates keys that activate corresponding values stored in the second FF transformation layer, functioning as a key-value memory system (more detail in Section 7.2.2). In addition, they also employed vocabulary projection methods to discover various kinds of semantic and syntactic knowledge stored in the FF sub-layers. Similarly, recent studies have also shown that factual (Meng et al., 2022) and commonsense knowledge (Gupta et al., 2023) are stored in the middle-layer FF sub-layers by employing intervention-based techniques. At a more granular level, specific FF neurons such as arithmetic neurons (Rai & Yao, 2024) and entropy neurons (Stolfo et al., 2024) have also been identified to store specific knowledge. In addition to FF sub-layers, recent studies indicate that *knowledge is also stored in the MHA sub-layers*, where each attention head seems to encode a specific type of knowledge (Gould et al., 2023; Geva et al., 2023). Besides localizing knowledge storage within LMs, Geva et al. (2023) further studied the circuit for the recall of facts or knowledge, i.e., given a subject (e.g., “Beats music”) and relation (e.g., “is owned by”), predicting the fact or attribute (e.g., “Apple”). Their study revealed a three-step mechanism: subject enrichment via FF sub-layers, relation propagation at the end token position via MHA sub-layers, and attribute extraction by later MHA sub-layers. Moreover, Chughtai et al. (2024) found that the recall of facts in LMs leverages several distinct, independent, and qualitatively different mechanisms or circuits. Although each mechanism alone might not suffice, their additive combination creates constructive interference that ultimately leads to the correct answer.

Others As listed in Section 7.2, prior work has also studied LM capabilities in tasks such as IOI (Wang et al., 2022a), modular addition (Nanda et al., 2023b), and greater-than operations (Hanna et al., 2024a), leading to the discovery of circuits that implement these tasks. Compared with the interpretation of ICL and reasoning, these studies not only justified the rationale of a capability but also revealed its underlying algorithm through circuits.

7.5 Findings on Learning Dynamics

Phase Changes during LM Training Prior studies have observed sudden shifts in LMs’ capabilities, called “phase changes” (Olsson et al., 2022; Power et al., 2022; Wei et al., 2022a). These changes are considered key steps during LM training. MI has been applied to examine the relationship between the emergence of features and circuits and these phase changes. For example, Olsson et al. (2022) found correlations between phase changes and the formation of induction circuits, suggesting that the development of these circuits underlies the phase change. In the task of symbol manipulation, Nanda et al. (2023a); Varma et al. (2023) discovered a similar correlation contributing to LM grokking (Power et al., 2022), a phenomenon of LM generalizing that can occur when they are trained beyond overfitting. Chen et al. (2024) found that sudden drops in the loss during training correspond to the acquisition of attention heads that recognize specific syntactic relations. Similarly, Wang et al. (2024a) examined the grokking phenomenon in implicit reasoning tasks and found that LMs can perform implicit reasoning through grokking, where the gradual formation of a generalizing circuit drives the process. Finally, Huang et al. (2024b) provided a unified explanation for grokking, double descent (Nakkiran et al., 2021), and emergent abilities (Wei et al., 2022a) as a competition between memorization and generalization circuits.

Effects of Post-training MI studies have investigated how fine-tuning enhances LM capabilities by analyzing underlying mechanistic changes in task-relevant circuits, LM representations, etc. These studies suggest that fine-tuning does not fundamentally change the mechanisms but enhances existing ones. For instance, Jain et al. (2023) observed that these changes are often localized in a few model weights and that simply pruning them can restore lost pre-training capabilities. Similarly, Jain et al. (2024) examined the effects of safety fine-tuning and found that it minimally alters the model, with the fine-tuning primarily adjusting FF weights to push unsafe inputs into a “null space”, thereby making the model significantly less sensitive to unsafe inputs. Jiang et al. (2024) examined catastrophic forgetting in fine-tuning, suggesting that it occurs not by erasing pre-existing skills but by introducing specialized reasoning patterns that overshadow the original capabilities. On the other hand, Lee et al. (2024) investigated how direct preference optimization (DPO) (Rafailov et al., 2023) alters the internal mechanisms of LMs to reduce toxicity. They found that DPO does not eliminate the LM’s ability to generate toxic outputs but instead reduces toxicity by learning to bypass regions in the model associated with toxic behavior. These findings imply that models aligned with DPO remain vulnerable to being “jailbroken” or reverted to an unaligned state, as the underlying toxic vectors persist and can be reactivated.

7.6 Applications of MI

7.6.1 Model Enhancement

Knowledge Editing LMs are known to store factual knowledge encountered during pre-training (Petroni et al., 2019; Cohen et al., 2023). For instance, when an LM is prompted with “The space needle is in the city of”, it may retrieve the stored facts and correctly predict “Seattle”. However, these stored facts may be incorrect or outdated over time, leading to factually incorrect generation (Cohen et al., 2024). MI has been found to be a helpful tool for addressing the problem, including understanding where and how facts are stored within LMs, how they are recalled during inference time, and providing approaches for knowledge editing (Meng et al., 2022; 2023; Geva et al., 2023; Sharma et al., 2024). For instance, Meng et al. (2022) used activation patching to localize components that are responsible for storing factual knowledge, and then edited the fact (e.g., replacing “Seattle” with “Paris”) by only updating the parameters of those components.

LM Generation Steering LM generation steering involves controlling an LM’s output by manipulating its activations at inference time. For instance, Geva et al. (2022) proposed a method to suppress toxic

language generation by identifying and manually activating neurons in FF layers responsible for promoting non-toxic or safe words. Similarly, Templeton et al. (2024) identified safety-related features (e.g., unsafe code, gender bias) and manipulated their activations to steer the LM towards (un)desired behaviors (e.g., safe code generation, unbiased text generation). Additionally, Nanda et al. (2023b) demonstrated that an LM’s output can be altered (e.g., flipping a player turn in the game of Othello from YOURS to MINE) by pushing its activation in the direction of a linear vector representing the desired behavior, which was identified using a linear probe.

7.6.2 AI Safety

AI safety is an important concern that MI aims to address. Although LMs undergo multiple rounds of safety fine-tuning, they frequently exhibit low reliability and susceptibility to harmful prompts (Casper et al., 2023b; MacDiarmid et al., 2024). Within MI, “enumerative safety” aims to address AI safety by enumerating all features in LMs and inspecting those related to dangerous capabilities or intentions (Elhage et al., 2022b; Olah & Jermyn, 2023). To this end, Templeton et al. (2024) identified several safety-relevant features that not only activate when the LM exhibits specific behaviors but also causally influence the LM’s output; however, the specific circuits that use these features to implement the behavior have not yet been identified. Similarly, Geva et al. (2022) discovered several non-toxic neurons and promoted them to steer the LM’s generation of non-toxic tokens. Recently, a new line of research in MI has centered on *latent space monitoring* to detect and prevent harmful model behaviors. Specifically, rather than inspecting the input and output texts for harmful model behaviors, this approach analyzes the latent representation in the model. To this end, these approaches have identified several safety-related linear directions within model latent representations that are useful for monitoring and limited control of AI safety-related behaviors such as truthfulness (Bürger et al., 2024), refusal (Arditi et al., 2024), and jailbreaking (Ball et al., 2024; Li et al., 2024). Finally, circuit-level insights have proven useful for detecting prompt injection attacks (Belrose et al., 2023) and identifying adversarial examples in tasks such as IOI (Wang et al., 2022a).

7.6.3 Others

Insights from MI have also been used for other downstream tasks. For example, Marks et al. (2024b) improved the generalization of classifiers by identifying and ablating spurious features that humans consider to be task-irrelevant. Geva et al. (2022) proposed self-supervised early exit prediction for efficient inference, drawing insights from their investigations of the FF sublayers’ role in token prediction.

8 Discussion and Future Work

8.1 Advancement in MI Techniques

Despite several promising recent advancements in MI techniques, the current set of techniques still faces a number of critical limitations that hinder their effectiveness and general applicability. One of the primary challenges is **scalability**, as many existing techniques (e.g., SAE, intervention-based techniques) require significant computational resources for interpreting larger LMs and increasingly complex model behaviors (e.g., reasoning). Furthermore, the *heavy reliance on human interpretation* of current techniques (e.g., interpreting SAE results) exacerbates scalability concerns while also introducing subjectivity, raising issues about the reproducibility and reliability of the findings. To address these challenges, developing more *automated techniques* for feature and circuit analysis is a promising direction, as it could reduce human effort while enhancing both the rigor and scalability of interpretability research. Additionally, exploring techniques that are *less computationally demanding*, such as attribution patching (Nanda et al., 2022), is another promising research direction. Beyond the scalability concern, the field has often focused on individual techniques in isolation rather than **combining complementary techniques** to improve their applicability and mitigate limitations of individual techniques. For instance, Patchscopes (Ghandeharioun et al., 2024) combines intervention-based and vocabulary-based techniques to obtain more expressive techniques, as discussed in Section 4.2. Finally, **addressing challenges of individual techniques**, as outlined under *Technical Advancements* in Section 4, is crucial for these techniques to reliably and effectively serve MI studies. For

instance, while intervention-based methods frequently suffer from out-of-distribution issues due to input corruption methods, probing techniques are inherently correlational and lack causal explanatory power. Exploring new techniques that are more efficient, automated, and integrated, while also addressing the limitations of existing techniques is essential for advancing the field and ensuring the reliability and applicability of MI insights.

8.2 Practical Utility of MI Studies

While the importance of deriving actionable insights from MI studies for downstream applications is widely recognized (Doshi-Velez & Kim, 2017), most MI studies often showcase these applications only as extrinsic evaluations on toy tasks, lacking rigorous comparisons against alternative methods or baselines. Furthermore, most MI studies, especially circuit studies, involve investigating simple LM behavior without downstream applications, often criticized as “streetlight interpretability” (Bereska & Gavves, 2024; Casper, 2023b; Wang, 2022). Similarly, although some studies have been conducted on “production-level” LMs (Lieberum et al., 2023; Templeton et al., 2024), the majority of studies still investigate smaller LMs, raising concerns about generalizability given the mixed findings on universality. Despite these challenges, a few MI studies have demonstrated promising practical downstream applications. For instance, recent studies have shown the potential of MI techniques in enhancing AI safety, notably through red-teaming efforts to identify model vulnerabilities (Arditi et al., 2024; Casper et al., 2023a). Additionally, several studies have also leveraged MI to develop new techniques (Zou et al., 2024; Ashuach et al., 2024; Guo et al., 2024; Pochinkov & Schoots, 2024) and perform evaluation (Deeb & Roger, 2024; Hong et al., 2024; Lynch et al., 2024) on tasks such as machine unlearning (Liu et al., 2025) and knowledge editing (Wang et al., 2024c). Finally, Kitouni et al. (2024) explored whether mechanistic approaches can uncover scientific knowledge embedded in models trained on prediction tasks. In summary, the advancement of MI research will benefit from a stronger emphasis on practical utility, which is crucial for reinforcing the field’s credibility and demonstrating its real-world impact.

8.3 Standardized Benchmarks and Metrics

Evaluating interpretability results is inherently challenging due to the lack of ground truth (Zhou et al., 2022), and current MI studies also employ various *ad hoc* evaluation approaches, potentially leading to inconsistent comparisons (Zhang & Nanda, 2023). On the other hand, there have been proposals for standardized evaluation benchmarks. RAVEL (Huang et al., 2024a) evaluates techniques (e.g., SAEs) that disentangle polysemantic neurons into monosemantic features. Tracr (Lindner et al., 2024) converts human-readable RASP programs (Weiss et al., 2021) into transformer models, enabling the construction of models with known ground-truth features and algorithms. These synthetic models provide a controlled setting to validate feature and circuit discovery algorithms, allowing researchers to assess whether these algorithms can accurately identify the known ground-truth features and mechanisms. However, these proposed benchmarks are still insufficient (Räuker et al., 2023). For instance, features and algorithms identified on synthetic transformers in the Tracr benchmark may not generalize to naturally-trained transformers. Thus, more effort is needed to develop standard evaluation techniques and to ensure their wide adoption.

8.4 Automated Hypothesis Generation in MI Practices

MI research often requires researchers to propose a hypothesis about the underlying mechanisms of LM behavior. However, this can be a laborious and non-scalable process. For example, to formulate plausible hypotheses while interpreting LM components in a circuit study, researchers need to manually analyze the activation patterns or intervention results (Section 5.2.1). This reliance on human intuition poses a potential scalability bottleneck, particularly when studying large-scale LLMs and complex behaviors. To mitigate this challenge, developing automated hypothesis-generation methods, potentially with human oversight, will be crucial for improving efficiency and scalability in MI research.

8.5 Generalizing Beyond MI for More Intuitive Applications

While MI was initially defined around the three objects as we described in Section 3.1, recent research has generalized beyond this scope for more intuitive applications. For instance, while existing feature studies have focused on decoding low-level input properties as features from the activations of LMs, recent works by Feng & Steinhardt (2023); Chalmers (2025) have called for decoding higher-level *proposition* from the activations. Consider finding safety-relevant features for enumeration safety (Elhage et al., 2022b) as an example. While being able to discover features encoding properties such as “kill” is helpful, this discovery alone cannot describe the safety propensity of the model; rather, higher-level propositions such as “I support killing humans” or “I oppose killing humans” can more directly allow one to measure the safety concern of an LM. A notable example of work in this direction is Feng & Steinhardt (2023), which investigates how individual features are combined to form higher-level concepts and how these concepts are represented within model activations. Specifically, it identified *binding vectors*—a mechanism that enables LMs to encode associations such as `lives(Alice, Paris)` and `lives(Bob, Bangkok)` when given input like “Alice lives in the capital city of France. Bob lives in the capital city of Thailand.” Understanding how models represent features at a higher level of abstraction seems critical for a more comprehensive and intuitive interpretation of LM behavior, highlighting the necessity for further research in this area.

8.6 A New Paradigm of Human-AI Collaboration Driven by MI

Several studies have explored the utility of post-hoc feature-attribution explanations (Ribeiro et al., 2016; Lundberg, 2017; Sundararajan et al., 2017) for improving human-AI team performance. However, most of these studies (Bansal et al., 2021; Chen et al., 2023; Carton et al., 2020; Rai et al., 2024) have found that these feature-attribution explanation fails to improve human-AI team performance and instead lead to a decline in performance compared to stand-alone AI. This decline in performance, however, may be attributed to the limitations of feature-based explanations, which simply highlight important input tokens as the model explanation without offering deeper insights into the model’s decision-making process. In contrast, MI provides a more comprehensive understanding by analyzing the model’s internal mechanisms, uncovering what features are extracted from the input tokens and how they influence the final output. Besides post-hoc explanation methods, textual explanations such as chain-of-thought (CoT) have also been proposed to be leveraged for model explanation; however, recent research suggests that CoT explanations can be unfaithful and misleading (Turpin et al., 2023). In light of these findings, MI not only has the potential to provide the most faithful explanations as it mostly relies on causal techniques but also presents a new type of “mechanistic explanation”, opening up new potential for human-AI collaboration. However, limited work has been done to investigate whether the mechanistic explanation generated through MI analysis can be used to improve human-AI team performance. One notable effort in employing the mechanistic understanding of LM behavior to improve human-AI team collaboration is provided by (Marks et al., 2024b). In this work, the authors proposed Spurious Human-interpretable Feature Trimming (SHIFT), a technique that enhances classifier generalization by removing the spurious features with the help of human evaluators. Their experiments demonstrate that SHIFT not only eliminates unintended biases but also improves the classifier’s overall performance. However, the practical application of MI findings such as SHIFT have only been considered as an extrinsic evaluation to validate MI studies, and as a result, they have been tested in only limited domains and benchmarks.

Acknowledgements

DR and ZY were sponsored by the National Science Foundation (#2311468/#2423813) and the Department of Computer Science at the College of Computing and Engineering at George Mason University. We appreciate comments from readers of the earlier version of this paper.

References

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

preprint arXiv:2303.08774, 2023.

- Abien Fred Agarap. Deep learning using rectified linear units (relu). *arXiv preprint arXiv:1803.08375*, 2018.
- Evan Anders and Joseph Bloom. Examining language model performance with reconstructed activations using sparse autoencoders. <https://www.lesswrong.com/posts/8QRH8wKcnKGhpAu2o/examining-language-model-performance-with-reconstructed>, 2024.
- Omer Antverg and Yonatan Belinkov. On the pitfalls of analyzing individual neurons in language models. In *International Conference on Learning Representations*, 2021.
- Andy Arditi, Oscar Obeso, Aaquib Syed, Daniel Paleka, Nina Panickssery, Wes Gurnee, and Neel Nanda. Refusal in language models is mediated by a single direction. *arXiv preprint arXiv:2406.11717*, 2024.
- Tomer Ashuach, Martin Tutek, and Yonatan Belinkov. Revs: Unlearning sensitive information in language models via rank editing in the vocabulary space. *arXiv preprint arXiv:2406.09325*, 2024.
- Jimmy Lei Ba. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- Sarah Ball, Frauke Kreuter, and Nina Panickssery. Understanding jailbreak success: A study of latent space dynamics in large language models, 2024.
- Gagan Bansal, Tongshuang Wu, Joyce Zhou, Raymond Fok, Besmira Nushi, Ece Kamar, Marco Tulio Ribeiro, and Daniel Weld. Does the whole exceed its parts? the effect of ai explanations on complementary team performance, 2021.
- Hritik Bansal, Karthik Gopalakrishnan, Saket Dingliwal, Sravan Bodapati, Katrin Kirchhoff, and Dan Roth. Rethinking the role of scale for in-context learning: An interpretability-based case study at 66 billion scale. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 11833–11856, 2023.
- Yonatan Belinkov. Probing classifiers: Promises, shortcomings, and advances. *Computational Linguistics*, 48(1):207–219, 2022.
- Nora Belrose, Zach Furman, Logan Smith, Danny Halawi, Igor Ostrovsky, Lev McKinney, Stella Biderman, and Jacob Steinhardt. Eliciting latent predictions from transformers with the tuned lens. *arXiv preprint arXiv:2303.08112*, 2023.
- Yoshua Bengio, Geoffrey Hinton, Andrew Yao, Dawn Song, Pieter Abbeel, Trevor Darrell, Yuval Noah Harari, Ya-Qin Zhang, Lan Xue, Shai Shalev-Shwartz, et al. Managing extreme ai risks amid rapid progress. *Science*, 384(6698):842–845, 2024.
- Leonard Bereska and Efstratios Gavves. Mechanistic interpretability for ai safety—a review. *arXiv preprint arXiv:2404.14082*, 2024.
- Adithya Bhaskar, Alexander Wettig, Dan Friedman, and Danqi Chen. Finding transformer circuits with edge pruning. *arXiv preprint arXiv:2406.16778*, 2024.
- Alberto Bietti, Vivien Cabannes, Diane Bouchacourt, Herve Jegou, and Leon Bottou. Birth of a transformer: A memory viewpoint. *Advances in Neural Information Processing Systems*, 36, 2024.
- Steven Bills, Nick Cammarata, Dan Mossing, Henk Tillman, Leo Gao, Gabriel Goh, Ilya Sutskever, Jan Leike, Jeff Wu, and William Saunders. Language models can explain neurons in language models. <https://openai.com/index/language-models-can-explain-neurons-in-language-models/>, 2023.
- Eden Biran, Daniela Gottesman, Sohee Yang, Mor Geva, and Amir Globerson. Hopping too late: Exploring the limitations of large language models on multi-hop queries. *arXiv preprint arXiv:2406.12775*, 2024.
- Dan Braun, Jordan Taylor, Nicholas Goldowsky-Dill, and Lee Sharkey. Identifying functionally important features with end-to-end sparse dictionary learning. *arXiv preprint arXiv:2405.12241*, 2024.

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. Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*, 2023. <https://transformer-circuits.pub/2023/monosemantic-features/index.html>.

Jannik Brinkmann, Abhay Sheshadri, Victor Levoso, Paul Swoboda, and Christian Bartelt. A mechanistic analysis of a transformer trained on a symbolic multi-step reasoning task. *arXiv preprint arXiv:2402.11917*, 2024.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.

Lennart Bürger, Fred A Hamprecht, and Boaz Nadler. Truth is universal: Robust detection of lies in llms. *arXiv preprint arXiv:2407.12831*, 2024.

Nick Cammarata, Gabriel Goh, Shan Carter, Ludwig Schubert, Michael Petrov, and Chris Olah. Curve detectors. *Distill*, 5(6):e00024–003, 2020.

Nicola Cancedda. Spectral filters, dark signals, and attention sinks. *arXiv preprint arXiv:2402.09221*, 2024.

Samuel Carton, Qiaozhu Mei, and Paul Resnick. Feature-based explanations don’t help people detect misclassifications of online toxicity, 2020.

Stephen Casper. The engineer’s interpretability sequence. <https://www.lesswrong.com/s/a6ne2ve5uturEEQK7>, 2023a.

Stephen Casper. Is mechanistic interpretability about to be practically useful? <https://www.lesswrong.com/s/a6ne2ve5uturEEQK7/p/Es2qzCvhJ8QYsckaA>, 2023b.

Stephen Casper, Tong Bu, Yuxiao Li, Jiawei Li, Kevin Zhang, Kaivalya Hariharan, and Dylan Hadfield-Menell. Red teaming deep neural networks with feature synthesis tools. *Advances in Neural Information Processing Systems*, 36:80470–80516, 2023a.

Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémie Scheurer, Javier Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, et al. Open problems and fundamental limitations of reinforcement learning from human feedback. *arXiv preprint arXiv:2307.15217*, 2023b.

David J Chalmers. Propositional interpretability in artificial intelligence. *arXiv preprint arXiv:2501.15740*, 2025.

Lawrence Chan, Adrià Garriga-Alonso, Nicholas Goldwosky-Dill, Ryan Greenblatt, Jenny Nitishinskaya, Ansh Radhakrishnan, Buck Shlegeris, and Nate Thomas. Causal scrubbing, a method for rigorously testing interpretability hypotheses. *AI Alignment Forum*, 2022a. <https://www.alignmentforum.org/posts/JvZhhzyCHu2Yd57RN/causal-scrubbing-a-method-for-rigorously-testing>.

Stephanie Chan, Adam Santoro, Andrew Lampinen, Jane Wang, Aaditya Singh, Pierre Richemond, James McClelland, and Felix Hill. Data distributional properties drive emergent in-context learning in transformers. *Advances in Neural Information Processing Systems*, 35:18878–18891, 2022b.

Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. A survey on evaluation of large language models. *ACM Transactions on Intelligent Systems and Technology*, 15(3):1–45, 2024.

-
- David Chanin, James Wilken-Smith, Tomáš Dulka, Hardik Bhatnagar, and Joseph Bloom. A is for absorption: Studying feature splitting and absorption in sparse autoencoders. *arXiv preprint arXiv:2409.14507*, 2024.
- Angelica Chen, Ravid Schwartz-Ziv, Kyunghyun Cho, Matthew L. Leavitt, and Naomi Saphra. Sudden drops in the loss: Syntax acquisition, phase transitions, and simplicity bias in mlms. In *The Twelfth International Conference on Learning Representations*, 2024.
- Valerie Chen, Q Vera Liao, Jennifer Wortman Vaughan, and Gagan Bansal. Understanding the role of human intuition on reliance in human-ai decision-making with explanations, 2023.
- Abhijith Chintam, Rahel Beloch, Willem Zuidema, Michael Hanna, and Oskar van der Wal. Identifying and adapting transformer-components responsible for gender bias in an English language model. In Yonatan Belinkov, Sophie Hao, Jaap Jumelet, Najoung Kim, Arya McCarthy, and Hosein Mohedbi (eds.), *Proceedings of the 6th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP*, pp. 379–394, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.blackboxnlp-1.29. URL <https://aclanthology.org/2023.blackboxnlp-1.29/>.
- Bilal Chughtai, Lawrence Chan, and Neel Nanda. A toy model of universality: Reverse engineering how networks learn group operations. In *International Conference on Machine Learning*, pp. 6243–6267. PMLR, 2023.
- Bilal Chughtai, Alan Cooney, and Neel Nanda. Summing up the facts: Additive mechanisms behind factual recall in llms. *arXiv preprint arXiv:2402.07321*, 2024.
- Roi Cohen, Mor Geva, Jonathan Berant, and Amir Globerson. Crawling the internal knowledge-base of language models. In *Findings of the Association for Computational Linguistics: EACL 2023*, pp. 1856–1869, 2023.
- Roi Cohen, Eden Biran, Ori Yoran, Amir Globerson, and Mor Geva. Evaluating the ripple effects of knowledge editing in language models. *Transactions of the Association for Computational Linguistics*, 12: 283–298, 2024.
- Arthur Conmy, Augustine Mavor-Parker, Aengus Lynch, Stefan Heimersheim, and Adrià Garriga-Alonso. Towards automated circuit discovery for mechanistic interpretability. *Advances in Neural Information Processing Systems*, 36:16318–16352, 2023.
- Alexis Conneau, Germán Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. What you can cram into a single vector: Probing sentence embeddings for linguistic properties. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2126–2136, 2018.
- Alan Cooney and Neel Nanda. Circuitsvis. <https://github.com/TransformerLensOrg/CircuitsVis>, 2023.
- Róbert Csordás, Christopher Potts, Christopher D Manning, and Atticus Geiger. Recurrent neural networks learn to store and generate sequences using non-linear representations. *arXiv preprint arXiv:2408.10920*, 2024.
- Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoencoders find highly interpretable features in language models. *arXiv preprint arXiv:2309.08600*, 2023.
- Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. Knowledge neurons in pretrained transformers. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 8493–8502, 2022.
- Guy Dar, Mor Geva, Ankit Gupta, and Jonathan Berant. Analyzing transformers in embedding space. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 16124–16170, 2023.

Aghyad Deeb and Fabien Roger. Do unlearning methods remove information from language model weights? *arXiv preprint arXiv:2410.08827*, 2024.

Alexander Yom Din, Taelin Karidi, Leshem Choshen, and Mor Geva. Jump to conclusions: Short-cutting transformers with linear transformations. *arXiv preprint arXiv:2303.09435*, 2023.

Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608*, 2017.

Subhabrata Dutta, Joykirat Singh, Soumen Chakrabarti, and Tanmoy Chakraborty. How to think step-by-step: A mechanistic understanding of chain-of-thought reasoning. *arXiv preprint arXiv:2402.18312*, 2024.

Benjamin L Edelman, Ezra Edelman, Surbhi Goel, Eran Malach, and Nikolaos Tsilivis. The evolution of statistical induction heads: In-context learning markov chains. *arXiv preprint arXiv:2402.11004*, 2024.

Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021. <https://transformer-circuits.pub/2021/framework/index.html>.

Nelson Elhage, Tristan Hume, Catherine Olsson, Neel Nanda, Tom Henighan, Scott Johnston, Sheer ElShowk, Nicholas Joseph, Nova DasSarma, Ben Mann, Danny Hernandez, Amanda Askell, Kamal Ndousse, Andy Jones, Dawn Drain, Anna Chen, Yuntao Bai, Deep Ganguli, Liane Lovitt, Zac Hatfield-Dodds, Jackson Kernion, Tom Conerly, Shauna Kravec, Stanislav Fort, Saurav Kadavath, Josh Jacobson, Eli Tran-Johnson, Jared Kaplan, Jack Clark, Tom Brown, Sam McCandlish, Dario Amodei, and Christopher Olah. Softmax linear units. *Transformer Circuits Thread*, 2022a. <https://transformer-circuits.pub/2022/solu/index.html>.

Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. Toy models of superposition. *Transformer Circuits Thread*, 2022b. https://transformer-circuits.pub/2022/toy_model/index.html.

Joshua Engels, Eric J Michaud, Isaac Liao, Wes Gurnee, and Max Tegmark. Not all language model features are linear. *arXiv preprint arXiv:2405.14860*, 2024.

N Benjamin Erichson, Zhewei Yao, and Michael W Mahoney. Jumprelu: A retrofit defense strategy for adversarial attacks. *arXiv preprint arXiv:1904.03750*, 2019.

Jiahai Feng and Jacob Steinhardt. How do language models bind entities in context? *arXiv preprint arXiv:2310.17191*, 2023.

Javier Ferrando, Gabriele Sarti, Arianna Bisazza, and Marta R Costa-jussà. A primer on the inner workings of transformer-based language models. *arXiv preprint arXiv:2405.00208*, 2024.

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

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

Atticus Geiger, Zhengxuan Wu, Christopher Potts, Thomas Icard, and Noah Goodman. Finding alignments between interpretable causal variables and distributed neural representations. In *Causal Learning and Reasoning*, pp. 160–187. PMLR, 2024.

-
- Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers are key-value memories. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 5484–5495, 2021.
- Mor Geva, Avi Caciularu, Kevin Wang, and Yoav Goldberg. 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*, pp. 30–45, 2022.
- Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. Dissecting recall of factual associations in auto-regressive language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 12216–12235, 2023.
- Asma Ghandeharioun, Avi Caciularu, Adam Pearce, Lucas Dixon, and Mor Geva. Patchscopes: a unifying framework for inspecting hidden representations of language models. In *Proceedings of the 41st International Conference on Machine Learning*, pp. 15466–15490, 2024.
- Nicholas Goldowsky-Dill, Chris MacLeod, Lucas Sato, and Aryaman Arora. Localizing model behavior with path patching. *arXiv preprint arXiv:2304.05969*, 2023.
- Rhys Gould, Euan Ong, George Ogden, and Arthur Conmy. Successor heads: Recurring, interpretable attention heads in the wild. In *The Twelfth International Conference on Learning Representations*, 2023.
- Phillip Guo, Aaquib Syed, Abhay Sheshadri, Aidan Ewart, and Gintare Karolina Dziugaite. Mechanistic unlearning: Robust knowledge unlearning and editing via mechanistic localization. *arXiv preprint arXiv:2410.12949*, 2024.
- Anshita Gupta, Debanjan Mondal, Akshay Krishna Sheshadri, Wenlong Zhao, Xiang Lorraine Li, Sarah Wiegreffe, and Niket Tandon. Editing common sense in transformers. *arXiv preprint arXiv:2305.14956*, 2023.
- Wes Gurnee. Sae reconstruction errors are (empirically) pathological. In *AI Alignment Forum*, pp. 16, 2024.
- Wes Gurnee and Max Tegmark. Language models represent space and time. *arXiv preprint arXiv:2310.02207*, 2023.
- Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris Bertsimas. Finding neurons in a haystack: Case studies with sparse probing. *Transactions on Machine Learning Research*, 2023.
- Wes Gurnee, Theo Horsley, Zifan Carl Guo, Tara Rezaei Kheirkhah, Qinyi Sun, Will Hathaway, Neel Nanda, and Dimitris Bertsimas. Universal neurons in gpt2 language models. *arXiv preprint arXiv:2401.12181*, 2024.
- Michael Hanna, Ollie Liu, and Alexandre Variengien. How does gpt-2 compute greater-than?: Interpreting mathematical abilities in a pre-trained language model. *Advances in Neural Information Processing Systems*, 36, 2024a.
- Michael Hanna, Sandro Pezzelle, and Yonatan Belinkov. Have faith in faithfulness: Going beyond circuit overlap when finding model mechanisms. *arXiv preprint arXiv:2403.17806*, 2024b.
- Zhengfu He, Wentao Shu, Xuyang Ge, Lingjie Chen, Junxuan Wang, Yunhua Zhou, Frances Liu, Qipeng Guo, Xuanjing Huang, Zuxuan Wu, et al. Llama scope: Extracting millions of features from llama-3.1-8b with sparse autoencoders. *arXiv preprint arXiv:2410.20526*, 2024.
- Stefan Heimersheim. You can remove gpt2’s layernorm by fine-tuning. *arXiv preprint arXiv:2409.13710*, 2024.
- Stefan Heimersheim and Jett Janiak. A circuit for python docstrings in a 4-layer attention-only transformer. <https://www.alignmentforum.org/posts/u6KXXmKFbXfWzoAXn/acircuit-for-python-docstrings-in-a-4-layer-attention-only>, 2023.

Stefan Heimersheim and Neel Nanda. How to use and interpret activation patching. *arXiv preprint arXiv:2404.15255*, 2024.

Roei Hendel, Mor Geva, and Amir Globerson. In-context learning creates task vectors. *arXiv preprint arXiv:2310.15916*, 2023.

Tom Henighan, Shan Carter, Tristan Hume, Nelson Elhage, Robert Lasenby, Stanislav Fort, Nicholas Schiefer, and Christopher Olah. Superposition, memorization, and double descent. *Transformer Circuits Thread*, 6:24, 2023.

Evan Hernandez, Arnab Sen Sharma, Tal Haklay, Kevin Meng, Martin Wattenberg, Jacob Andreas, Yonatan Belinkov, and David Bau. Linearity of relation decoding in transformer language models. *arXiv preprint arXiv:2308.09124*, 2023.

Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.

Yihuai Hong, Lei Yu, Haiqin Yang, Shauli Ravfogel, and Mor Geva. Intrinsic evaluation of unlearning using parametric knowledge traces. *arXiv preprint arXiv:2406.11614*, 2024.

Jie Huang and Kevin Chen-Chuan Chang. Towards reasoning in large language models: A survey. *arXiv preprint arXiv:2212.10403*, 2022.

Jing Huang, Zhengxuan Wu, Christopher Potts, Mor Geva, and Atticus Geiger. Ravel: Evaluating interpretability methods on disentangling language model representations. *arXiv preprint arXiv:2402.17700*, 2024a.

Yufei Huang, Shengding Hu, Xu Han, Zhiyuan Liu, and Maosong Sun. Unified view of grokking, double descent and emergent abilities: A perspective from circuits competition. *arXiv preprint arXiv:2402.15175*, 2024b.

Samyak Jain, Robert Kirk, Ekdeep Singh Lubana, Robert P Dick, Hidenori Tanaka, Edward Grefenstette, Tim Rocktäschel, and David Scott Krueger. Mechanistically analyzing the effects of fine-tuning on procedurally defined tasks. *arXiv preprint arXiv:2311.12786*, 2023.

Samyak Jain, Ekdeep Singh Lubana, Kemal Oksuz, Tom Joy, Philip HS Torr, Amartya Sanyal, and Puneet K Dokania. What makes and breaks safety fine-tuning? a mechanistic study. *arXiv preprint arXiv:2407.10264*, 2024.

Stanisław Jastrzębski, Devansh Arpit, Nicolas Ballas, Vikas Verma, Tong Che, and Yoshua Bengio. Residual connections encourage iterative inference. *arXiv preprint arXiv:1710.04773*, 2017.

Gangwei Jiang, Caigao Jiang, Zhaoyi Li, Siqiao Xue, Jun Zhou, Linqi Song, Defu Lian, and Ying Wei. Interpretable catastrophic forgetting of large language model fine-tuning via instruction vector. *arXiv preprint arXiv:2406.12227*, 2024.

Shahar Katz, Yonatan Belinkov, Mor Geva, and Lior Wolf. Backward lens: Projecting language model gradients into the vocabulary space. *arXiv preprint arXiv:2402.12865*, 2024.

Connor Kissane, Robert Krzyzanowski, Joseph Isaac Bloom, Arthur Conmy, and Neel Nanda. Interpreting attention layer outputs with sparse autoencoders. *arXiv preprint arXiv:2406.17759*, 2024.

Oual Kitouni, Niklas Nolte, Víctor Samuel Pérez-Díaz, Sokratis Trifinopoulos, and Mike Williams. From neurons to neutrons: A case study in interpretability. *arXiv preprint arXiv:2405.17425*, 2024.

Robert Krzyzanowski, Connor Kissane, Arthur Conmy, and Neel Nanda. We inspected every head in gpt-2 small using saes so you don't have to. In *Alignment Forum*, 2024.

Michael Lan, Philip Torr, Austin Meek, Ashkan Khakzar, David Krueger, and Fazl Barez. Sparse autoencoders reveal universal feature spaces across large language models. *arXiv preprint arXiv:2410.06981*, 2024.

Andrew Lee, Xiaoyan Bai, Itamar Pres, Martin Wattenberg, Jonathan K Kummerfeld, and Rada Mihalcea. A mechanistic understanding of alignment algorithms: A case study on dpo and toxicity. In *International Conference on Machine Learning*, pp. 26361–26378. PMLR, 2024.

Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-time intervention: Eliciting truthful answers from a language model, 2024.

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

Tom Lieberum, Senthooran Rajamanoharan, Arthur Conmy, Lewis Smith, Nicolas Sonnerat, Vikrant Varma, János Kramár, Anca Dragan, Rohin Shah, and Neel Nanda. Gemma scope: Open sparse autoencoders everywhere all at once on gemma 2. *arXiv preprint arXiv:2408.05147*, 2024.

Johnny Lin. Neuronpedia: Interactive reference and tooling for analyzing neural networks, 2023. URL <https://www.neuronpedia.org>. Software available from neuronpedia.org.

David Lindner, János Kramár, Sebastian Farquhar, Matthew Rahtz, Tom McGrath, and Vladimir Mikulik. Tracr: Compiled transformers as a laboratory for interpretability. *Advances in Neural Information Processing Systems*, 36, 2024.

Sijia Liu, Yuanshun Yao, Jinghan Jia, Stephen Casper, Nathalie Baracaldo, Peter Hase, Yuguang Yao, Chris Yuhan Liu, Xiaojun Xu, Hang Li, et al. Rethinking machine unlearning for large language models. *Nature Machine Intelligence*, pp. 1–14, 2025.

Scott Lundberg. A unified approach to interpreting model predictions, 2017.

Aengus Lynch, Phillip Guo, Aidan Ewart, Stephen Casper, and Dylan Hadfield-Menell. Eight methods to evaluate robust unlearning in llms. *arXiv preprint arXiv:2402.16835*, 2024.

Monte MacDiarmid, Timothy Maxwell, Nicholas Schiefer, Jesse Mu, Jared Kaplan, David Duvenaud, Sam Bowman, Alex Tamkin, Ethan Perez, Mrinank Sharma, Carson Denison, and Evan Hubinger. Simple probes can catch sleeper agents, 2024. URL <https://www.anthropic.com/news/probes-catch-sleeper-agents>.

Aleksandar Makelov. Sparse autoencoders match supervised features for model steering on the ioi task. In *ICML 2024 Workshop on Mechanistic Interpretability*, 2024.

Aleksandar Makelov, George Lange, and Neel Nanda. Towards principled evaluations of sparse autoencoders for interpretability and control. *arXiv preprint arXiv:2405.08366*, 2024.

Luke Marks, Alisdair Paren, David Krueger, and Fazl Barez. Enhancing neural network interpretability with feature-aligned sparse autoencoders. *arXiv preprint arXiv:2411.01220*, 2024a.

Samuel Marks and Max Tegmark. The geometry of truth: Emergent linear structure in large language model representations of true/false datasets. *arXiv preprint arXiv:2310.06824*, 2023.

Samuel Marks, Can Rager, Eric J Michaud, Yonatan Belinkov, David Bau, and Aaron Mueller. Sparse feature circuits: Discovering and editing interpretable causal graphs in language models. *arXiv preprint arXiv:2403.19647*, 2024b.

Callum McDougall, Arthur Conmy, Cody Rushing, Thomas McGrath, and Neel Nanda. Copy suppression: Comprehensively understanding an attention head. *arXiv preprint arXiv:2310.04625*, 2023.

-
- Tianyi Men, Pengfei Cao, Zhuoran Jin, Yubo Chen, Kang Liu, and Jun Zhao. Unlocking the future: Exploring look-ahead planning mechanistic interpretability in large language models. *arXiv preprint arXiv:2406.16033*, 2024.
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372, 2022.
- Kevin Meng, Arnab Sen Sharma, Alex Andonian, Yonatan Belinkov, and David Bau. Mass-editing memory in a transformer. In *The Eleventh International Conference on Learning Representations*, 2023.
- Jack Merullo, Carsten Eickhoff, and Ellie Pavlick. Circuit component reuse across tasks in transformer language models. In *The Twelfth International Conference on Learning Representations*, 2023.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems*, 26, 2013.
- Aashiq Muhammed, Mona Diab, and Virginia Smith. Decoding dark matter: Specialized sparse autoencoders for interpreting rare concepts in foundation models. *arXiv preprint arXiv:2411.00743*, 2024.
- Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever. Deep double descent: Where bigger models and more data hurt. *Journal of Statistical Mechanics: Theory and Experiment*, 2021(12):124003, 2021.
- Neel Nanda, Chris Olah, Catherine Olsson, Nelson Elhage, and Hume Tristan. Attribution patching: Activation patching at industrial scale, 2023. <https://www.neelnanda.io/mechanistic-interpretability-attribution-patching>, 2022.
- Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. Progress measures for grokking via mechanistic interpretability. *arXiv preprint arXiv:2301.05217*, 2023a.
- Neel Nanda, Andrew Lee, and Martin Wattenberg. Emergent linear representations in world models of self-supervised sequence models. *arXiv preprint arXiv:2309.00941*, 2023b.
- Yaniv Nikankin, Anja Reusch, Aaron Mueller, and Yonatan Belinkov. Arithmetic without algorithms: Language models solve math with a bag of heuristics. *arXiv preprint arXiv:2410.21272*, 2024.
- nostalgia. Interpreting gpt: the logit lens. *AI Alignment Forum*, 2020. <https://www.lesswrong.com/posts/AcKRB8wDpdaN6v6ru/interpreting-gpt-the-logit-lens>.
- Chris Olah and Adam Jermyn. What would be the most safety-relevant features in language models? (circuits updates - july 2023). *Transformer Circuits Thread*, 2023. <https://transformer-circuits.pub/2023/july-update/index.html#safety-features>.
- Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter. Zoom in: An introduction to circuits. *Distill*, 2020. <https://distill.pub/2020/circuits/zoom-in>.
- Christopher Olah. What is a linear representation? what is a multidimensional feature?, 2023. URL <https://transformer-circuits.pub/2024/july-update/#linear-representations>.
- Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. In-context learning and induction heads. *Transformer Circuits Thread*, 2022. <https://transformer-circuits.pub/2022/in-context-learning-and-induction-heads/index.html>.
- Koyena Pal, Jiuding Sun, Andrew Yuan, Byron C Wallace, and David Bau. Future lens: Anticipating subsequent tokens from a single hidden state. *arXiv preprint arXiv:2311.04897*, 2023.

-
- Judea Pearl. *Causality*. Cambridge university press, 2009.
- Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pp. 1532–1543, 2014.
- Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. Language models as knowledge bases? 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)*, pp. 2463–2473, 2019.
- Nicholas Pochinkov and Nandi Schoots. Dissecting language models: Machine unlearning via selective pruning. *arXiv preprint arXiv:2403.01267*, 2024.
- Alethea Power, Yuri Burda, Harri Edwards, Igor Babuschkin, and Vedant Misra. Grokking: Generalization beyond overfitting on small algorithmic datasets. *arXiv preprint arXiv:2201.02177*, 2022.
- Philip Quirke, Clement Neo, and Fazl Barez. Increasing trust in language models through the reuse of verified circuits. *arXiv preprint arXiv:2402.02619*, 2024.
- Alec Radford, Rafal Jozefowicz, and Ilya Sutskever. Learning to generate reviews and discovering sentiment. *arXiv preprint arXiv:1704.01444*, 2017.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36:53728–53741, 2023.
- Daking Rai and Ziyu Yao. An investigation of neuron activation as a unified lens to explain chain-of-thought eliciting arithmetic reasoning of llms. *arXiv preprint arXiv:2406.12288*, 2024.
- Daking Rai, Rydia R Weiland, Kayla Margaret Gabriella Herrera, Tyler H Shaw, and Ziyu Yao. Understanding the effect of algorithm transparency of model explanations in text-to-sql semantic parsing. *arXiv preprint arXiv:2410.16283*, 2024.
- Senthooran Rajamanoharan, Arthur Conmy, Lewis Smith, Tom Lieberum, Vikrant Varma, János Kramár, Rohin Shah, and Neel Nanda. Improving dictionary learning with gated sparse autoencoders. *arXiv preprint arXiv:2404.16014*, 2024a.
- Senthooran Rajamanoharan, Tom Lieberum, Nicolas Sonnerat, Arthur Conmy, Vikrant Varma, János Kramár, and Neel Nanda. Jumping ahead: Improving reconstruction fidelity with jumprelu sparse autoencoders. *arXiv preprint arXiv:2407.14435*, 2024b.
- Tilman Räuker, Anson Ho, Stephen Casper, and Dylan Hadfield-Menell. Toward transparent ai: A survey on interpreting the inner structures of deep neural networks. In *2023 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)*, pp. 464–483. IEEE, 2023.
- Gautam Reddy. The mechanistic basis of data dependence and abrupt learning in an in-context classification task. In *The Twelfth International Conference on Learning Representations*, 2023.
- Jie Ren, Qipeng Guo, Hang Yan, Dongrui Liu, Xipeng Qiu, and Dahua Lin. Identifying semantic induction heads to understand in-context learning. *arXiv preprint arXiv:2402.13055*, 2024.
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. " 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*, pp. 1135–1144, 2016.

Logan Riggs. Found 600+ monosemantic features in a small lm using sparse autoencoders. <https://www.lesswrong.com/posts/wqRqb7h6ZC48iDgfK/tentatively-found-600-monosemantic-features-in-a-small-lm>, 2023.

Mansi Sakravadia, Arham Khan, Aswathy Ajith, Daniel Grzenda, Nathaniel Hudson, André Bauer, Kyle Chard, and Ian Foster. Attention lens: A tool for mechanistically interpreting the attention head information retrieval mechanism. *arXiv preprint arXiv:2310.16270*, 2023.

Abulhair Saparov, Srushti Pawar, Shreyas Pimpalgaonkar, Nitish Joshi, Richard Yuanzhe Pang, Vishakh Padmakumar, Seyed Mehran Kazemi, Najoung Kim, and He He. Transformers struggle to learn to search. *arXiv preprint arXiv:2412.04703*, 2024.

Naomi Saphra and Sarah Wiegreffe. Mechanistic? *arXiv preprint arXiv:2410.09087*, 2024.

Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, et al. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*, 2022.

Adam Scherlis, Kshitij Sachan, Adam S Jermyn, Joe Benton, and Buck Shlegeris. Polysemanticity and capacity in neural networks. *arXiv preprint arXiv:2210.01892*, 2022.

Lee Sharkey, Dan Braun, and Beren Millidge. Taking features out of superposition with sparse autoencoders, 2022. <https://www.alignmentforum.org/posts/z6QQJbtpkEAX3Aojj/interim-research-report-taking-features-out-of-superposition>, 2022.

Arnab Sen Sharma, David Atkinson, and David Bau. Locating and editing factual associations in mamba. *arXiv preprint arXiv:2404.03646*, 2024.

Alessandro Stolfo, Yonatan Belinkov, and Mrinmaya Sachan. A mechanistic interpretation of arithmetic reasoning in language models using causal mediation analysis. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 7035–7052, 2023a.

Alessandro Stolfo, Yonatan Belinkov, and Mrinmaya Sachan. Understanding arithmetic reasoning in language models using causal mediation analysis. *arXiv preprint arXiv:2305.15054*, 2023b.

Alessandro Stolfo, Ben Wu, Wes Gurnee, Yonatan Belinkov, Xingyi Song, Mrinmaya Sachan, and Neel Nanda. Confidence regulation neurons in language models. *arXiv preprint arXiv:2406.16254*, 2024.

Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *International conference on machine learning*, pp. 3319–3328. PMLR, 2017.

Aaquib Syed, Can Rager, and Arthur Conmy. Attribution patching outperforms automated circuit discovery. *arXiv preprint arXiv:2310.10348*, 2023.

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, Calum McDougall, Monte MacDiarmid, C. Daniel Freeman, Theodore R. Sumers, Edward Rees, Joshua Batson, Adam Jermyn, Shan Carter, Chris Olah, and Tom Henighan. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. *Transformer Circuits Thread*, 2024. <https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html>.

Ian Tenney, Dipanjan Das, and Ellie Pavlick. Bert rediscovers the classical nlp pipeline. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4593–4601, 2019a.

Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R Thomas McCoy, Najoung Kim, Benjamin Van Durme, Samuel R Bowman, Dipanjan Das, et al. What do you learn from context? probing for sentence structure in contextualized word representations. *arXiv preprint arXiv:1905.06316*, 2019b.

Curt Tigges, Oskar John Hollinsworth, Atticus Geiger, and Neel Nanda. Linear representations of sentiment in large language models. *arXiv preprint arXiv:2310.15154*, 2023.

-
- Curt Tigges, Michael Hanna, Qinan Yu, and Stella Biderman. Llm circuit analyses are consistent across training and scale. *arXiv preprint arXiv:2407.10827*, 2024.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Miles Turpin, Julian Michael, Ethan Perez, and Samuel Bowman. Language models don't always say what they think: Unfaithful explanations in chain-of-thought prompting. *Advances in Neural Information Processing Systems*, 36:74952–74965, 2023.
- Niels uit de Bos and Adrià Garriga-Alonso. Adversarial circuit evaluation. In *ICML 2024 Workshop on Mechanistic Interpretability*, 2024.
- Vikrant Varma, Rohin Shah, Zachary Kenton, János Kramár, and Ramana Kumar. Explaining grokking through circuit efficiency. *arXiv preprint arXiv:2309.02390*, 2023.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Jesse Vig. Visualizing attention in transformer-based language representation models. *arXiv preprint arXiv:1904.02679*, 2019.
- Elena Voita, Javier Ferrando, and Christoforos Nalmpantis. Neurons in large language models: Dead, n-gram, positional. *arXiv preprint arXiv:2309.04827*, 2023.
- Boshi Wang, Xiang Yue, Yu Su, and Huan Sun. Grokked transformers are implicit reasoners: A mechanistic journey to the edge of generalization. *arXiv preprint arXiv:2405.15071*, 2024a.
- Junxuan Wang, Xuyang Ge, Wentao Shu, Qiong Tang, Yunhua Zhou, Zhengfu He, and Xipeng Qiu. Towards universality: Studying mechanistic similarity across language model architectures. *arXiv preprint arXiv:2410.06672*, 2024b.
- Kevin Wang. Gears-level mental models of transformer interpretability. <https://www.lesswrong.com/posts/X26ksz4p3wSyyCKN/gears-level-mental-models-of-transformer-interpretability>, 2022.
- Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. Interpretability in the wild: a circuit for indirect object identification in gpt-2 small. In *The Eleventh International Conference on Learning Representations*, 2022a.
- Mengru Wang, Yunzhi Yao, Ziwen Xu, Shuofei Qiao, Shumin Deng, Peng Wang, Xiang Chen, Jia-Chen Gu, Yong Jiang, Pengjun Xie, et al. Knowledge mechanisms in large language models: A survey and perspective. *arXiv preprint arXiv:2407.15017*, 2024c.
- Xiaozhi Wang, Kaiyue Wen, Zhengyan Zhang, Lei Hou, Zhiyuan Liu, and Juanzi Li. Finding skill neurons in pre-trained transformer-based language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 11132–11152, 2022b.
- Zihao Wang, Lin Gui, Jeffrey Negrea, and Victor Veitch. Concept algebra for (score-based) text-controlled generative models. *Advances in Neural Information Processing Systems*, 36, 2024d.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. *Transactions on Machine Learning Research*, 2022a.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022b.

-
- Laura Weidinger, Jonathan Uesato, Maribeth Rauh, Conor Griffin, Po-Sen Huang, John Mellor, Amelia Glaese, Myra Cheng, Borja Balle, Atoosa Kasirzadeh, et al. Taxonomy of risks posed by language models. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, pp. 214–229, 2022.
- Gail Weiss, Yoav Goldberg, and Eran Yahav. Thinking like transformers. In *International Conference on Machine Learning*, pp. 11080–11090. PMLR, 2021.
- Eric Winsor. Re-examining layernorm. *AI Alignment Forum*, 2020. <https://www.alignmentforum.org/posts/jfG6vdJZCwTQmG7kb/re-examining-layernorm>.
- Benjamin Wright and Lee Sharkey. Addressing feature suppression in saes. In *AI Alignment Forum*, pp. 16, 2024.
- Yifan Yao, Jinhao Duan, Kaidi Xu, Yuanfang Cai, Zhibo Sun, and Yue Zhang. A survey on large language model (llm) security and privacy: The good, the bad, and the ugly. *High-Confidence Computing*, pp. 100211, 2024.
- Fred Zhang and Neel Nanda. Towards best practices of activation patching in language models: Metrics and methods. In *The Twelfth International Conference on Learning Representations*, 2023.
- Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu, Huiqi Deng, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, and Mengnan Du. Explainability for large language models: A survey. *ACM Transactions on Intelligent Systems and Technology*, 15(2):1–38, 2024.
- Ziqian Zhong, Ziming Liu, Max Tegmark, and Jacob Andreas. The clock and the pizza: Two stories in mechanistic explanation of neural networks. *Advances in Neural Information Processing Systems*, 36, 2024.
- Yilun Zhou, Serena Booth, Marco Tulio Ribeiro, and Julie Shah. Do feature attribution methods correctly attribute features? In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 9623–9633, 2022.
- Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, et al. Representation engineering: A top-down approach to ai transparency. *arXiv preprint arXiv:2310.01405*, 2023.
- Andy Zou, Long Phan, Justin Wang, Derek Duenas, Maxwell Lin, Maksym Andriushchenko, J Zico Kolter, Matt Fredrikson, and Dan Hendrycks. Improving alignment and robustness with circuit breakers. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.