

Towards Uncovering How Large Language Models Work: An Interpretability Perspective

Anonymous ACL submission

Abstract

Large language models (LLMs) have led to breakthroughs in language tasks, yet the internal mechanisms that enable their remarkable generalization and reasoning abilities remain opaque. This lack of transparency presents challenges such as hallucinations, toxicity, and misalignment with human values, hindering the safe and beneficial deployment of LLMs. This survey paper aims to uncover the internal working mechanisms underlying LLM functionality through the lens of explainability. First, we review how knowledge is encoded within LLMs via mechanistic interpretability techniques. Then, we summarize how knowledge is embedded in LLM representations by leveraging probing techniques and representation engineering. Additionally, we investigate the training dynamics through a mechanistic perspective to explain phenomena such as grokking and memorization. Lastly, we explore how the insights gained from these explanations can enhance LLM performance through model editing, improve efficiency through pruning, and better align with human values.

1 Introduction

Large language models (LLMs) such as GPT-4 (OpenAI, 2023), LLaMA-2 (Touvron et al., 2023), Claude-3 (AnthropicAI, 2023), and Gemini (Team et al., 2023) have led to tremendous advancements in language understanding and generation, achieving state-of-the-art performance in a wide array of real-world tasks. Despite their superior performance across various tasks, the “how” and “why” behind their generalization and reasoning abilities are still not well understood. This lack of understanding poses several challenges. First, LLMs frequently generate hallucinations and factually incorrect output, which complicates efforts to improve their performance. Second, as LLMs become more powerful, problems surrounding potential toxicity, unfairness, and dishonesty threaten

to undermine user trust. Therefore, there is an urgent need to delve deeper into the inner workings of LLMs to fully address these issues. Gaining insights into how these models operate is a crucial step towards developing robust safeguards and ensuring their responsible deployment.

In this paper, we provide a systematic overview of the existing literature that uncovers the internal working mechanisms of LLMs using explainability techniques (Figure 1). First, we provide a summary of findings on how knowledge is encoded within the architecture of trained LLMs. The explainability technique, mechanistic interpretability, seems promising in providing such explanation. It focuses on the functionality of each model component and interprets how models operate at the level of neurons, circuits, and attention heads. Second, we examine how knowledge is encoded internally in intermediate representations. To this end, representation engineering is adopted to explain specific behavior of the model, such as dishonesty, by analyzing hidden representations. Specifically, representation engineering focuses on identifying patterns for certain behaviors through probing-based methods and employs them to mitigate undesired behaviors (Zou et al., 2023). Third, we inspect the model training process to understand the development of generalization abilities during the training process. Finally, we review how insights from the aforementioned analysis help us improve models in terms of higher performance through model editing, better efficiency through pruning, and better human alignment.

Our work differs from existing survey articles on the explainability of LLMs (Zhao et al., 2023; Wu et al., 2024b; Luo and Specia, 2024; Ferrando et al., 2024), which either summarize explainability techniques or discuss their utilities. In contrast, our goal is to review very recent studies that fundamentally provide insights into trained LLMs and their dynamic training processes. We focus on works

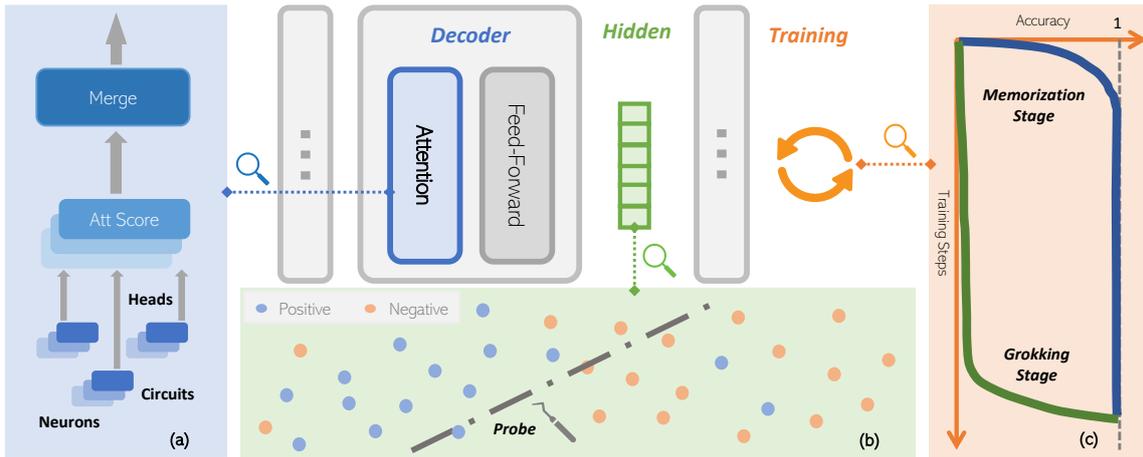


Figure 1: In this work, we review existing progress on how LLMs work, including: (a) how knowledge is encoded within model components; (b) what knowledge is encoded in intermediate representations; and (c) how generalization abilities are achieved during the training process.

083 that uncover how LLMs function and identify the
 084 factors that contribute to their reasoning abilities
 085 via using explainability techniques and monitoring
 086 the training process. We review the state-of-the-
 087 art insights on the inner working mechanisms of
 088 LLMs and explore how these insights can further
 089 enhance model performance and benefit humans.

090 2 How is Knowledge Encoded in Model 091 Architectures?

092 LLMs are built on extensive training datasets and
 093 intricate model architectures, which contribute to
 094 their remarkable emergent abilities (Wei et al.,
 095 2022). However, the exact mechanisms through
 096 which these models acquire and process vast
 097 amounts of knowledge remain unclear. Addition-
 098 ally, the contributions of individual model com-
 099 ponents to the overall function have been largely
 100 unexplored. To fully understand LLMs, recent stud-
 101 ies have shifted to make use of mechanistic inter-
 102 pretability (more details are given in Section A
 103 at the Appendix) to reverse engineer LLMs at a more
 104 granular level such as neurons and attention heads.

105 2.1 Neurons

106 Neurons can activate on knowledge and patterns
 107 within LLMs. They are observed to be *polyse-*
 108 *mantic*, meaning that an individual neuron can
 109 be activated on multiple unrelated terms (Olah
 110 et al., 2020a; Bills et al., 2023). This character-
 111 istic presents a significant challenge in mechanistically
 112 understanding how models operate. Despite the
 113 challenge, recent work has explored the underlying
 114 causes of this polysemantic nature. Two key con-

115 cepts have emerged as instrumental in unraveling
 116 its formation: *Superposition* (Olah et al., 2020a)
 117 and *Monosemanticity* (Bricken et al., 2023).

118 2.1.1 Superposition

119 Superposition describes the phenomenon where
 120 a feature can be spread across multiple neurons,
 121 meanwhile a neuron can also be mixed up with mul-
 122 tiple features. Some researchers believe that this
 123 mechanism originated from an excessive number of
 124 features compared to the number of neurons (Olah
 125 et al., 2020a; Elhage et al., 2022). In exploring
 126 this concept through a toy example, i.e. a ReLU
 127 network, researchers have found that superposition
 128 allows for the representation of additional features.
 129 However, to mitigate interference, a nonlinear fil-
 130 ter needs to be introduced (Elhage et al., 2022).
 131 When features are sparse, superposition effectively
 132 supports the representation of these features and al-
 133 lows computations such as the absolute value func-
 134 tion (Elhage et al., 2022). Neurons within models
 135 can be either monosemantic or polysemantic.

136 Others argue that polysemanticity arises inci-
 137 dentally due to factors encountered during the
 138 training process such as regularization and neu-
 139 ral noise (Lecomte et al., 2024). Mathematical
 140 demonstrations have shown that a constant frac-
 141 tion of feature collisions, introduced through random
 142 initialization, can always result in polysemantic
 143 neurons, even when the number of neurons exceeds
 144 the number of features (Lecomte et al., 2024).

145 Another study investigates polysemanticity
 146 through the lens of the “feature capacity”, denoting
 147 the fraction of embedding dimensions consumed

by a feature in the representation space (Scherlis et al., 2022). By analyzing one-layer and two-layer toy models, this work indicates that features are represented based on their importance in reducing loss. More important features are allocated their own dimensions, while the less critical ones may be overlooked, and the rest will share embedding dimensions (Scherlis et al., 2022). Features only end up sharing dimensions when assigning additional capacity will not result in loss decreasing (Scherlis et al., 2022). Moreover, the relationship between superposition and feature importance has been demonstrated on LLMs (Gurnee et al., 2023). Experiments show that the early layers tend to represent many features in superposition, while the middle layers include dedicated neurons to represent high-level features (Gurnee et al., 2023).

2.1.2 Monosemanticity

Monosemantic neurons, associated with a single concept, are much easier to interpret than polysemantic neurons. Investigating the factors that enhance monosemanticity is meaningful to model interpretation. A research using toy models reveals that changing the loss minimum could improve monosemanticity. Such loss minimum usually coexists with negative biases (Jermyn et al., 2022). However, in reality building a purely monosemantic model is infeasible due to the unmanageable loss (Bricken et al., 2023). Another line of studies seeks to disentangle superposition to reach a monosemantic understanding. The sparse autoencoder emerges as a promising tool serving this purpose, aiming to reconstruct sparsely activated directions that are more interpretable and monosemantic (Cunningham et al., 2023). The method utilizes dictionary learning where features are pre-defined (Sharkey et al., 2022). The effectiveness of this approach largely depends on the comprehensiveness of the pre-defined dictionary. Bricken et al. (2023) utilizes it to interpret a one-layer transformer model with a 512-neuron MLP layer. The sparse autoencoder is trained on MLP activations from 8B data points, with autoencoder sizes ranging from 512 to 13,100 features. Larger autoencoders are able to achieve finer granularity in interpreting features, revealing details that cannot be discovered at the neuron level. These identified features can be used to manipulate the model’s output, offering new ways to control and understand LLMs (Bricken et al., 2023).

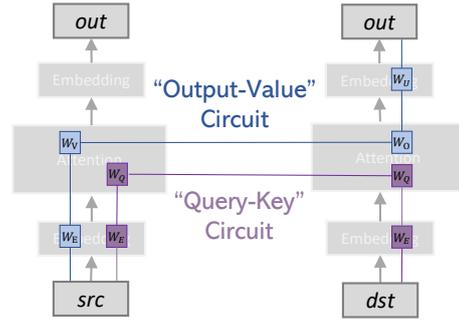


Figure 2: An illustration of a Transformer circuit, which is a key concept in mechanistic interpretability.

2.2 Circuits

Circuit is one of the core concepts in the field of mechanistic interpretability (see Figure 2). It was originally proposed to reverse engineer vision models, in which individual neurons and their connections are viewed as functional units (Olah et al., 2020a). Researchers have found that features in former layers of models act as fundamental units, such as edge detectors. These features are combined through weights to form a circuit unit. This viewpoint is partially evidenced by a few understandable neuron units (or circuits) performing specific functions, such as curve detectors (Cammarata et al., 2020) and high-low frequency detectors (Schubert et al., 2021). Several interesting phenomena have been observed in these circuits. For example, symmetric transformations of basic features, including copying, scaling, flipping, coloring, rotating, can be achieved with basic neurons known as “equivariance” or “motif” (Olah et al., 2020b).

Despite rich insights from vision models, transformer models, with their unique architecture featuring attention blocks, present new challenges. To address these, a mathematical framework specifically for *transformer circuits* has also been proposed (Elhage et al., 2021). This framework simplifies the complex architecture of LLMs by focusing on decoder-only transformer models that have no more than two layers, all made up entirely of attention blocks. Within this toy model, the transformer encompasses input embedding, residual stream, attention layers, and output embeddings. Attention layers read information from the residual stream and then write their output back into it. Consequently, communication is achieved through read and write operations at the layer level.

Each attention head works independently and in parallel, contributing its output to the residual

stream. These heads consist of key, query, output, and value vectors, represented as W_K , W_Q , W_O and W_V . There are two types of circuits: i) “query-key” (QK) circuits; ii) “output-value” (OV) circuits (Elhage et al., 2021), as shown in Figure 2. The QK circuits, formed by $W_Q^T W_K$, play a crucial role in determining which previously learned token to copy information from (Elhage et al., 2021). It is essential for models to recall and retrieve information from earlier context. Conversely, the OV circuits, composed of $W_O W_V$, determine how the current token influences the output logits (Elhage et al., 2021).

The result shows that transformers with no layer can model bigram statistics, predicting the next token from the source token. Adding one layer allows the model to capture both bigram and “skip-trigram” patterns. Interestingly, with two layers, transformer models give rise to a concept termed as “*induction head*” (Section 2.3). These induction heads exist in the second layer and beyond. Usually, they are composed of heads from their previous layer, which is useful in suggesting the next token based on the present ones (Elhage et al., 2021).

A circuit composed of 28 attention heads have been identified to enable indirection object identification tasks as well (Wang et al., 2022). And these heads route information between heads finally to the outputs. Also, Hanna et al. (2024) found a circuit that performs greater-than computation on a set of MLPs in GPT-2.

2.3 Attention heads

A special type of attention head called *induction head* is assumed critical in enabling in-context learning abilities within LLMs (Brown et al., 2020), due to the co-occurrence of induction heads and in-context learning (Olsson et al., 2022). Induction heads also refer to a kind of circuits that complete the pattern by prefix matching and copying previously occurred sequences (Olsson et al., 2022). They are composed of two heads: the first attention head is from the previous layer attending to previous tokens that are followed by the current token, which achieves prefix matching and provides the attend-to token (the token following current token). The second head, i.e. induction head, copies the attend-to token and increases its output logits. More specifically, this rule means that if models have seen similar patterns such as “[A*][B*]” given current token “[A]”, these models are able to pre-

dict “[B]” (Olsson et al., 2022). Despite the single token used in the toy example, long prefix matching such as three consecutive tokens has also been observed in related work (Chan et al., 2022).

As a result, layers with induction heads possess more powerful in-context learning abilities than simple copying. In addition, multiple empirical studies have demonstrated the causal relationships between induction heads and in-context learning abilities by observing the change of in-context learning abilities after manipulating induction heads (Olsson et al., 2022; Chan et al., 2022). Although this theory offers a comprehensive explanation of the mechanisms behind transformer models with only two attention layers, further ablation studies are still needed to validate its effectiveness. It is also important to note that this framework is exclusively based on attention heads, without incorporating MLP layers.

A recent study reveals how factual associations are stored and extracted within LLMs (Geva et al., 2023). A fact consists of a subject, a relation, and an object. During inferring, the subject will be enriched with subject-related attributes at the early MLP sublayers. The relation will be propagated to prediction through attention heads. And the prediction representation will “query” enriched subject to extract object information (Geva et al., 2023). Moreover, Chughtai et al. (2024) believes that there are also mixed heads containing information from both subject and relation. The subject heads, mixed heads, and relation heads are working additively to elicit the outputs.

3 What Knowledge is Encoded in Intermediate Representations?

In the previous section, we summarize existing studies on the architectural composition of knowledge within LLMs, with a focus on their structural components. We highlight how components of LLMs function differently. In this section, we introduce an in-depth review of the knowledge encoded by *representations of LLMs*, including world knowledge and factual knowledge captured within these models. We examine how factors such as the depth of layers and the scale of models influence this encoding process.

3.1 Probing World and Factual Knowledge

To investigate whether the representations of LLMs encode world knowledge and factual knowledge,

335 probing techniques offer crucial insights into the
336 structure and dynamics of these representations.
337 Specifically, probing techniques can identify specific
338 directions within the representation space, and
339 these directions are essential for understanding certain
340 behaviors and the encoding of knowledge (Zou
341 et al., 2023; Liu et al., 2023).

342 Recent studies have demonstrated that LLMs
343 can learn world models and encode them in their
344 representations for certain tasks. One study suc-
345 cessfully uses a set of non-linear probes to uncover
346 world representations within models, specifically
347 in the context of the game of Othello (Li et al.,
348 2022). It demonstrates models’ ability to track the
349 board state, and make predictions without being
350 explicitly to do so (Li et al., 2022). Furthermore,
351 another work finds that linear representation struc-
352 tures can also perform well on predictions, sim-
353 ply by altering the expression of the board state at
354 each timestamp (Nanda et al., 2023b). The linear
355 and non-linear explanations reveal how models per-
356 ceive the world naturally, which might be different
357 from humans. Additionally, by analyzing repre-
358 sentations of spatial datasets, one study reveals the
359 model’s ability to learn linear representations of
360 space and time across multiple levels (Gurnee and
361 Tegmark, 2023).

362 LLMs are also capable of encoding factual
363 knowledge. Marks and Tegmark (2023) craft self-
364 curated true/false datasets to study the geometry
365 of representations of true/false statements derived
366 from a model’s residual stream. By applying princi-
367 pal component analysis (PCA), a clear linear struc-
368 ture emerges. The truth directions are leveraged
369 to mediate the model’s dishonest behaviors locally.
370 Another research avenue explores vectors related
371 to toxicity within MLP blocks through singular
372 value decomposition (SVD). The identified dimen-
373 sions are simply subtracted to efficiently achieve
374 mitigation (Lee et al., 2024).

375 Function vectors have also been discovered
376 within the attention heads of LLMs, which trigger
377 the execution of a certain task across diverse inputs.
378 For example, Todd et al. (2023) found that these
379 function vectors are shown in various in-context
380 learning tasks, and can execute related tasks de-
381 spite zero-shot inputs. Also, causal interventions
382 at the neuron level can help identify the individ-
383 ual neurons encoding spatial coordinates and time
384 information (Gurnee and Tegmark, 2023).

385 Lastly, representations associated with undesir-

386 able behaviors of LLMs, such as dishonesty, toxic-
387 ity, hallucinations, can also be extracted. Typically,
388 a direction in the representation space is identified
389 as contributing to a specific behavior. This direc-
390 tion will then be used to adjust the representations
391 so that models’ behaviors can be controlled (Zou
392 et al., 2023). For example, Li et al. (2024) employs
393 this technique to probe and enhance the truthful-
394 ness of models. Azaria and Mitchell (2023) also
395 successfully distinguishes the truthfulness of state-
396 ments by simply training a classifier on model rep-
397 resentations. A recent work has been developed to
398 identify hallucination tokens from the response by
399 integrating a range of classifiers that are trained on
400 each layer from separate hidden parts: MLPs and
401 attention layers (CH-Wang et al., 2023).

3.2 Role of Layer Depth and Model Scale 402

403 The influence of layer depth and model scale on
404 representations has been an interesting research di-
405 rection. Empirically, research shows that a range of
406 knowledge is well trained until the middle layers.
407 For example, Gurnee and Tegmark (2023) demon-
408 strate that space and time representations reach the
409 best quality up to half of the layers in a range of
410 open-source LLMs. Besides, the function vectors
411 with strong causal effects are also collected from
412 the middle layers of LLMs, while the effects are
413 near zero in the deeper layers (Todd et al., 2023).
414 Furthermore, another study shows that different lev-
415 els of concepts are well learned in different layers,
416 where simpler tasks are learned in the early layers
417 while complex tasks can only be well learned in
418 the deeper layers (Jin et al., 2024; Ju et al., 2024).
419 However, the underlying reason why the middle
420 layers perform so well remains unexplored.

421 It is generally believed that more capabilities are
422 gained as the models scale up (Wei et al., 2022).
423 Some recent studies have also supported this hy-
424 pothesis in certain cases. For example, the space
425 and time representations are more precise as the
426 models scale up (Gurnee and Tegmark, 2023). But
427 the inner mechanism leading to better performance
428 when model scales remain unknown.

4 How is Generalization Ability Achieved During Training? 429

430 In the preceding sections, we analyze LLMs in
431 a post-hoc manner, focusing on neurons, connec-
432 tions, attention heads, and representations to under-
433 stand how knowledge is acquired within models.
434

In this section, we discuss the dynamic training process of models to understand how the generalization ability is achieved during the training process. We will particularly examine two important phenomena observed in relation to generalization: grokking and memorization. Here, grokking indicates the phenomenon where models suddenly improve validation accuracy after overfitting. Investigating grokking can shed light on how generalization emerges during training. Moreover, examining memorization, where models rely on statistical patterns rather than causal relationships, can help disentangle the roles of generalization versus the roles of memorization in model behaviors.

4.1 Understanding Grokking

Grokking is a phenomenon in which models suddenly improve their validation accuracy after severely overfitting on over-parameterized neural networks (Power et al., 2022). The surge in validation accuracy is generally interpreted as a gain of generalization ability.

4.1.1 A Data Perspective

Experiments implemented on a two-layer decoder-only transformer network have shown that grokking is closely related to factors such as data, representations, and regularization. Smaller datasets require more optimization steps for grokking to occur (Power et al., 2022). Conversely, more samples can decrease the number of steps needed for generalization (Zhu et al., 2024). The minimal amount of data needed for grokking also depends on the minimal number of data points required to learn a robust representation (Liu et al., 2022a). Furthermore, it has been found that generalization often coincides with well-structured embeddings. Additionally, regularization measures can accelerate the onset of grokking, with weight decay standing out as particularly effective in strengthening generalization capabilities (Liu et al., 2022a). A recent study proposes that massive datasets in LLMs make grokking less conceivable (Zhu et al., 2024).

4.1.2 Weight Norms

When examining the weight norms of the final layers in models that do not use regularization techniques, a phenomenon, termed as *slingshot mechanism*, has been observed. It describes a cyclic behavior during the terminal phase of training, where there are oscillations between stable and unstable

regimes, i.e., training loss spike. The spike occurs with a phase where weight norms grow, followed by a phase of norm plateau. Thilak et al. (2022) point out that grokking, non-trivial feature adaptation, occurs only at the beginning of slingshots. The appearance of the slingshot effect and grokking can be modulated by adjusting the optimizer parameters, especially when using adaptive optimizers such as Adam (Kingma and Ba, 2014). However, it is unclear whether this observation holds universally across various scenarios.

Additionally, another concept called the *LU mechanism* has also been proposed, describing dynamics between loss and weight norms (Liu et al., 2022b). In algorithmic datasets, an L-shaped training loss and a U-shaped test loss reduction concerning weight norms are identified, implying an optimal range for initializing weight norms. Nevertheless, this finding does not seamlessly transfer to real-world machine learning tasks, where large initialization and small weight decay are often necessary. Lyu et al. (2023); Mohamadi et al. (2023) attribute it to a competition between the early-phase implicit bias favoring kernel predictors induced by large initialization and a late-phase implicit bias favoring min-norm/margin predictors promoted by small weight decay. Similarly, Merrill et al. (2023) conclude that this competition manifests a competition between a dense subnetwork in the initial phase and a sparse one after grokking. However, a recent research on deep neural networks believes that feature ranks/linear probing accuracy could be a better indicator of phase transition than weight norm (Fan et al., 2024).

4.1.3 Test Loss

Double descent captures the pattern where a model’s test accuracy at the log level initially improves, then drops due to overfitting, and finally increases again after gaining generalization abilities (Nakkiran et al., 2021). This pattern is more noticeable in the test loss. A unified framework has been developed to integrate grokking with double descent, treating them as two manifestations of the same underlying process (Davies et al., 2023). The framework attributes the transition of generalization to slower pattern learning, which has been further supported by Kumar et al. (2023). This transition is demonstrated to exist at the level of both epochs and models.

533	4.2 Memorization		
534	<i>Memorization</i> often refers to the phenomenon that	Research has shown that it is possible to edit factual	583
535	models predict with statistical features rather than	knowledge by modifying the weights of specific	584
536	causal relations. The study using slightly corrupted	neurons in MLPs. One study successfully adopts	585
537	algorithmic datasets with two-layer neural models	this approach by altering neural computations re-	586
538	has revealed that memorization can coexist with	lated to recall of factual knowledge (Meng et al.,	587
539	generalization. And memorization can be miti-	2022). Another study expands this method further	588
540	gated by pruning relevant neurons or by regular-	to allow multiple edits at the same time (Meng	589
541	ization (Doshi et al., 2023). Although different	et al., 2023). Although these methods are effective	590
542	regularization methods might not share learning	for targeted edits, their capabilities on updating	591
543	goals, they all contribute to better representations.	relevant knowledge and preventing forgetting still	592
544	And the training process in the study consists of	require further investigation (Cohen et al., 2023).	593
545	two stages: i) the grokking process, ii) the decay of	Interestingly, a recent study indicates that the	594
546	memorization learning (Doshi et al., 2023). How-	paragraphs memorized by a model can be pin-	595
547	ever, the underlying causes behind this process are	pointed using high-gradient weights in attention	596
548	not yet fully understood. Besides, the assumption	heads of the lower layers (Stoehr et al., 2024).	597
549	that regularization is the key to this process is under	This research employs localization techniques to	598
550	debate, especially in light of observing grokking in	identify specific attention heads, which are then	599
551	absence of regularization (Kumar et al., 2023). The	fine-tuned to unlearn the memorized knowledge.	600
552	importance of the rate of feature learning and the	This approach holds promise in enhancing privacy	601
553	number of necessary features is favored in explana-	protection in large language models, although a	602
554	tions, challenging the role of the weight norm (Ku-	comprehensive evaluation is still needed.	603
555	mar et al., 2023).	Besides, facts are also encoded in the representa-	604
556	Interestingly, a study hypothesizes that memo-	tion space, making representation a natural candi-	605
557	rization constitutes a phase of grokking (Nanda	date to edit models' outputs. So far, most studies fo-	606
558	et al., 2023a). The study finds that grokking in-	cus on modifying representations at inference time,	607
559	cludes three distinct stages: memorization, circuit	while the influence of permanent modifications has	608
560	formation, and memorization cleanup (Nanda et al.,	barely been studied. A recent work provides a more	609
561	2023a). The study identifies an algorithm that uti-	precise way to edit model representations to change	610
562	lizes Discrete Fourier Transforms and trigonomet-	their output distributions (Hernandez et al., 2023).	611
563	ric identities to achieve modular addition through	Instead of only adding the derived vectors into out-	612
564	analyzing the model's weights. The circuits en-	put representations, this study directly changes the	613
565	abling this algorithm seem to evolve in a steady	embedding of a related entity so as to trigger tar-	614
566	manner instead of randomly walking. Varma et al.	getted outputs. As a result, the position of the mod-	615
567	(2023) concludes that the efficiency of memori-	ified entity in the embedding space has changed,	616
568	sation and generalization depends on the size of	leading to causal influence on model generations.	617
569	dataset. However, our understanding towards the		
570	relationship between memorization and grokking	5.2 Model Pruning for Better Efficiency	618
571	is still limited.		
572	5 How to Make Use of The Insights?		
573	In the preceding three sections, we have explored	In contrast to deciphering the inner workings of	619
574	how knowledge is encoded within LLMs (Section	models, one study examines the differences be-	620
575	2), and how this knowledge is encoded in their rep-	tween pre-training and fine-tuning phases with	621
576	resentations (Section 3). Building on these insights,	mechanistic interpretability tools. It reveals that	622
577	this section emphasizes on how we can leverage	fine-tuning retains all the capabilities learned in	623
578	our in-depth understanding of LLMs to enhance	the pre-training phase. Transformations between	624
579	their performance through editing, improve their	pre-training and fine-tuning stem from "wrappers"	625
580	efficiency via pruning, and better align them with	in MLPs learned on top of models. Interestingly,	626
581	human values and preferences.	these wrappers can be eliminated by pruning a few	627
		neurons or retraining on an unrelated downstream	628
		task (Jain et al., 2023). This discovery sheds light	629
		on potential safety concerns associated with current	630
		alignment approaches.	631

Different from pruning neurons, the idea of representation engineering, that is directly manipulating representations without the need for optimization or additional labeled data, has also been demonstrated effective in model pruning. Some work attempts to fine-tune models with representation engineering and achieves a comparable and even better performance than state-of-the-art fine-tuning techniques (Wu et al., 2024a,c). One work employs forward passes from two topics and derives their difference vectors, which are used in inference time without additional fine-tuning (Turner et al., 2023). Wu et al. (2024a) also demonstrates the feasibility of fine-tuning models through editing representations. Unlike conventional parameter-efficient fine-tuning (PEFT), representation editing focuses on learning an additional group of trainable parameters to modify representations directly other than models’ parameters. And the trainable parameters have been reduced to a factor of 32 compared to that of LoRA (Hu et al., 2021; Wu et al., 2024a). Another approach utilizes the distributed alignment search of Geiger et al. (2024) to find a set of linear subspace implementing interventions. This method outperforms most PEFT models on a range of tasks (Wu et al., 2024c).

5.3 Model Alignment to Human Values

From the mechanistic perspective, practical applications tend to evaluate model alignments with different tools. Inspired by induction heads, a recent work measures bias scores of attention heads in pre-trained LLMs, focusing on specific stereotypes. It implemented a method to ensure the accuracy of identifying biased heads by comparing the changes of attention score between biased and regular heads. Through masking identified biased heads, the study effectively reduces the gender bias encoded in the model (Yang et al., 2023). Besides, another work localizes attention heads that are responsible to lie with linear probing and activation patching. A set of intentionally designed prompts is used to instruct LLMs to be dishonest. Meanwhile, linear probes are trained to classify true/false activations of heads. Then, the selected activations are patched with those of honest behaviors to observe the changes of outputs. Multiple attention heads across five layers are causally located in Campbell et al. (2023).

Recently, representation engineering has emerged as a promising avenue for detecting biases

within embedding space. A notable study suggests that MLPs operate on token representations to alter the distribution of output vocabulary (Geva et al., 2022). After reverse engineering MLPs, it is believed that the output from each feed-forward layer can be seen as sub-updates to output vocabulary distributions, essentially promoting certain high-level concepts. This insight has been used effectively to mitigate toxicity levels in LLMs (Geva et al., 2022). Another line of work finds multiple representation vectors within MLPs that encourage models’ undesired behaviors. These vectors are decomposed using singular value decomposition, allowing researchers to pinpoint specific dimensions that contribute to toxicity (Lee et al., 2024).

6 Conclusions and Looking Beyond

In this paper, we explore techniques to uncover the inner workings of LLMs through an explainability lens. We provide a systematic overview of how explainability techniques can reveal the architectural composition of knowledge within LLMs and the encoding of knowledge in their internal representations. Furthermore, we inspect training dynamics through a mechanistic perspective to explain phenomena like “grokking” that can explain generalization abilities of LLMs. Lastly, we reviewed how insights from these explainability analyses can enhance LLM performance through model editing, improve efficiency via pruning, and better align models with human preferences.

Although there is some preliminary progress in uncovering the inner workings of LLMs, looking beyond, there exist several critical challenges and opportunities. First, LLMs have encoded a vast amount of real-world knowledge into their architectures and parameters. However, current research has only revealed a small fraction of the encoded knowledge. Future efforts should focus on developing scalable techniques that can effectively analyze and interpret the intricate knowledge structures embedded within LLMs. Second, LLMs have demonstrated remarkable reasoning abilities that exhibit human-like cognitive abilities. However, our current understanding of how these high-level reasoning abilities emerge from the interplay of architectural components and training dynamics is limited. More efforts are needed to reveal the intricate mechanisms that give rise to these advanced reasoning capabilities.

732 Limitations

733 In this paper, we intend to integrate available tech-
734 niques that enable us to learn the inner workings of
735 LLMs. Despite the valuable perspectives provided,
736 our study has several limitations. First, we do not
737 explore the complete landscape of relevant XAI
738 methods for understanding LLMs, due to space con-
739 straints. Other techniques like concept-based expla-
740 nations, example-based explanations, and counter-
741 factual explanations may also provide some useful
742 insights into the inner workings of LLMs. These
743 methods could potentially uncover additional as-
744 pects or offer complementary viewpoints that are
745 not covered by the mechanistic interpretability and
746 representation engineering approaches discussed
747 in this paper. Furthermore, while we try to provide
748 a comprehensive overview of the current state-of-
749 the-art, the field of explainable AI for LLMs is
750 rapidly evolving. New techniques, theories, and
751 findings may emerge that could reshape or extend
752 our understanding of how LLM works. Continuous
753 monitoring and incorporating these developments
754 will be crucial to maintaining a comprehensive and
755 up-to-date perspective on this topic.

756 References

757 AnthropicAI. 2023. [Introducing claude](#).

758 Amos Azaria and Tom Mitchell. 2023. The internal
759 state of an llm knows when its lying. *arXiv preprint*
760 *arXiv:2304.13734*.

761 Steven Bills, Nick Cammarata, Dan Moss-
762 ing, Henk Tillman, Leo Gao, Gabriel Goh,
763 Ilya Sutskever, Jan Leike, Jeff Wu, and
764 William Saunders. 2023. Language mod-
765 els can explain neurons in language models.
766 [https://openai-public.blob.core.windows-](https://openai-public.blob.core.windows.net/neuron-explainer/paper/index.html)
767 [.net/neuron-explainer/paper/index.html](https://openai-public.blob.core.windows.net/neuron-explainer/paper/index.html).

768 Trenton Bricken, Adly Templeton, Joshua Batson,
769 Brian Chen, Adam Jermy, Tom Conerly, Nick
770 Turner, Cem Anil, Carson Denison, Amanda Askell,
771 Robert Lasenby, Yifan Wu, Shauna Kravec, Nicholas
772 Schiefer, Tim Maxwell, Nicholas Joseph, Zac
773 Hatfield-Dodds, Alex Tamkin, Karina Nguyen,
774 Brayden McLean, Josiah E Burke, Tristan Hume,
775 Shan Carter, Tom Henighan, and Christopher
776 Olah. 2023. Towards monosemanticity: Decom-
777 posing language models with dictionary learning.
778 *Transformer Circuits Thread*. [https://transformer-](https://transformer-circuits.pub/2023/monosemantic-features/index.html)
779 [circuits.pub/2023/monosemantic-](https://transformer-circuits.pub/2023/monosemantic-features/index.html)
780 [features/index.html](https://transformer-circuits.pub/2023/monosemantic-features/index.html).

781 Tom B Brown, Benjamin Mann, Nick Ryder, Melanie
782 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind

Neelakantan, Pranav Shyam, Girish Sastry, Amanda
783 Askell, et al. 2020. Language models are few-shot
784 learners. *Advances in Neural Information Processing*
785 *Systems (NeurIPS)*. 786

Nick Cammarata, Gabriel Goh, Shan Carter,
787 Ludwig Schubert, Michael Petrov, and
788 Chris Olah. 2020. [Curve detectors](#). *Distill*.
789 <https://distill.pub/2020/circuits/curve-detectors>. 790

James Campbell, Richard Ren, and Phillip Guo.
791 2023. Localizing lying in llama: Understanding in-
792 structed dishonesty on true-false questions through
793 prompting, probing, and patching. *arXiv preprint*
794 *arXiv:2311.15131*. 795

Sky CH-Wang, Benjamin Van Durme, Jason Eisner,
796 and Chris Kedzie. 2023. Do androids know they’re
797 only dreaming of electric sheep? *arXiv preprint*
798 *arXiv:2312.17249*. 799

Lawrence Chan, Adrià Garriga-Alonso, Nicholas
800 Goldwosky-Dill, Ryan Greenblatt, Jenny Nitishin-
801 skaya, Ansh Radhakrishnan, Buck Shlegeris, and
802 Nate Thomas. 2022. Causal scrubbing, a method
803 for rigorously testing interpretability hypotheses. *AI*
804 *Alignment Forum*. 805

Bilal Chughtai, Lawrence Chan, and Neel Nanda. 2023.
806 Neural networks learn representation theory: Reverse
807 engineering how networks perform group operations.
808 In *ICLR 2023 Workshop on Physics for Machine*
809 *Learning*. 810

Bilal Chughtai, Alan Cooney, and Neel Nanda. 2024.
811 Summing up the facts: Additive mechanisms
812 behind factual recall in llms. *arXiv preprint*
813 *arXiv:2402.07321*. 814

Roi Cohen, Eden Biran, Ori Yoran, Amir Globerson,
815 and Mor Geva. 2023. Evaluating the ripple effects
816 of knowledge editing in language models. *arXiv*
817 *preprint arXiv:2307.12976*. 818

Arthur Conmy, Augustine N Mavor-Parker, Aengus
819 Lynch, Stefan Heimersheim, and Adrià Garriga-
820 Alonso. 2023. Towards automated circuit discov-
821 ery for mechanistic interpretability. *arXiv preprint*
822 *arXiv:2304.14997*. 823

Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert
824 Huben, and Lee Sharkey. 2023. Sparse autoencoders
825 find highly interpretable features in language models.
826 *arXiv preprint arXiv:2309.08600*. 827

Xander Davies, Lauro Langosco, and David Krueger.
828 2023. Unifying grokking and double descent. *arXiv*
829 *preprint arXiv:2303.06173*. 830

Darshil Doshi, Aritra Das, Tianyu He, and Andrey Gro-
831 mov. 2023. To grok or not to grok: Disentangling
832 generalization and memorization on corrupted algo-
833 rithmic datasets. *arXiv preprint arXiv:2310.13061*. 834

835	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. 2022. Toy models of superposition. <i>Transformer Circuits Thread</i> . https://transformer-circuits.pub/2022/toy_model/index.html .	889
836		890
837		891
838		
839		892
840		893
841		894
842		895
		896
843	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. 2021. A mathematical framework for transformer circuits. <i>Transformer Circuits Thread</i> . https://transformer-circuits.pub/2021/framework/index.html .	
844		
845		897
846		898
847		899
848		900
849		901
850		
851		902
852		903
853		904
854	Simin Fan, Razvan Pascanu, and Martin Jaggi. 2024. Deep grokking: Would deep neural networks generalize better? <i>arXiv preprint arXiv:2405.19454</i> .	905
855		906
856		907
		908
857	Javier Ferrando, Gabriele Sarti, Arianna Bisazza, and Marta R Costa-jussà. 2024. A primer on the inner workings of transformer-based language models. <i>arXiv preprint arXiv:2405.00208</i> .	909
858		
859		
860		
861	Dan Friedman, Alexander Wettig, and Danqi Chen. 2023. Learning transformer programs. <i>arXiv preprint arXiv:2306.01128</i> .	910
862		911
863		912
		913
864	Atticus Geiger, Zhengxuan Wu, Christopher Potts, Thomas Icard, and Noah Goodman. 2024. Finding alignments between interpretable causal variables and distributed neural representations. In <i>Causal Learning and Reasoning</i> , pages 160–187. PMLR.	914
865		915
866		
867		
868		
869	Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. 2023. Dissecting recall of factual associations in auto-regressive language models. <i>arXiv preprint arXiv:2304.14767</i> .	916
870		917
871		918
872		
873	Mor Geva, Avi Caciularu, Kevin Ro Wang, and Yoav Goldberg. 2022. Transformer feed-forward layers build predictions by promoting concepts in the vocabulary space. <i>arXiv preprint arXiv:2203.14680</i> .	919
874		920
875		921
876		922
		923
877	Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2021. Transformer feed-forward layers are key-value memories. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 5484–5495.	924
878		
879		
880		
881		
882	Andrey Gromov. 2023. Grokking modular arithmetic. <i>arXiv preprint arXiv:2301.02679</i> .	925
883		926
		927
884	Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris Bertsimas. 2023. Finding neurons in a haystack: Case studies with sparse probing. <i>arXiv preprint arXiv:2305.01610</i> .	928
885		929
886		
887		
888		
	Wes Gurnee and Max Tegmark. 2023. Language models represent space and time. <i>arXiv preprint arXiv:2310.02207</i> .	930
		931
	Michael Hanna, Ollie Liu, and Alexandre Variengien. 2024. How does gpt-2 compute greater-than?: Interpreting mathematical abilities in a pre-trained language model. <i>Advances in Neural Information Processing Systems</i> , 36.	932
		933
		934
		935
		936
	Peter Hase, Mohit Bansal, Been Kim, and Asma Ghandeharioun. 2023. Does localization inform editing? surprising differences in causality-based localization vs. knowledge editing in language models. <i>arXiv preprint arXiv:2301.04213</i> .	937
		938
		939
		940
	Evan Hernandez, Belinda Z. Li, and Jacob Andreas. 2023. Inspecting and editing knowledge representations in language models.	930
		903
		904
	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. <i>arXiv preprint arXiv:2106.09685</i> .	905
		906
		907
		908
		909
	Samyak Jain, Robert Kirk, Ekdeep Singh Lubana, Robert P Dick, Hidenori Tanaka, Edward Grefenstette, Tim Rocktäschel, and David Scott Krueger. 2023. Mechanistically analyzing the effects of fine-tuning on procedurally defined tasks. <i>arXiv preprint arXiv:2311.12786</i> .	910
		911
		912
		913
		914
		915
	Adam S Jermyn, Nicholas Schiefer, and Evan Hubinger. 2022. Engineering monosemanticity in toy models. <i>arXiv preprint arXiv:2211.09169</i> .	916
		917
		918
	Mingyu Jin, Qinkai Yu, Jingyuan Huang, Qingcheng Zeng, Zhenting Wang, Wenyue Hua, Haiyan Zhao, Kai Mei, Yanda Meng, Kaize Ding, Fan Yang, Mengnan Du, and Yongfeng Zhang. 2024. Exploring concept depth: How large language models acquire knowledge at different layers?	919
		920
		921
		922
		923
		924
	Tianjie Ju, Weiwei Sun, Wei Du, Xinwei Yuan, Zhaochun Ren, and Gongshen Liu. 2024. How large language models encode context knowledge? a layer-wise probing study. <i>arXiv preprint arXiv:2402.16061</i> .	925
		926
		927
		928
		929
	Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. <i>arXiv preprint arXiv:1412.6980</i> .	930
		931
		932
	Pang Wei Koh and Percy Liang. 2017. Understanding black-box predictions via influence functions. In <i>International conference on machine learning</i> , pages 1885–1894. PMLR.	933
		934
		935
		936
	Tanishq Kumar, Blake Bordelon, Samuel J Gershman, and Cengiz Pehlevan. 2023. Grokking as the transition from lazy to rich training dynamics. <i>arXiv preprint arXiv:2310.06110</i> .	937
		938
		939
		940

941	Victor Lecomte, Kushal Thaman, Rylan Schaeffer, Naomi Bashkansky, Trevor Chow, and Sanmi Koyejo. 2024. What causes polysemanticity? an alternative origin story of mixed selectivity from incidental causes. In <i>ICLR 2024 Workshop on Representational Alignment</i> .	Samuel Marks and Max Tegmark. 2023. The geometry of truth: Emergent linear structure in large language model representations of true/false datasets. <i>arXiv preprint arXiv:2310.06824</i> .	995 996 997 998
947	Andrew Lee, Xiaoyan Bai, Itamar Pres, Martin Wattenberg, Jonathan K Kummerfeld, and Rada Mihalcea. 2024. A mechanistic understanding of alignment algorithms: A case study on dpo and toxicity. <i>arXiv preprint arXiv:2401.01967</i> .	Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in gpt. <i>Advances in Neural Information Processing Systems</i> , 35:17359–17372.	999 1000 1001 1002
952	Noam Levi, Alon Beck, and Yohai Bar-Sinai. 2023. Grokking in linear estimators—a solvable model that groks without understanding. <i>arXiv preprint arXiv:2310.16441</i> .	Kevin Meng, Arnab Sen Sharma, Alex J Andonian, Yonatan Belinkov, and David Bau. 2023. Mass-editing memory in a transformer. In <i>The Eleventh International Conference on Learning Representations</i> .	1003 1004 1005 1006 1007
956	Kenneth Li, Aspen K Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2022. Emergent world representations: Exploring a sequence model trained on a synthetic task. <i>arXiv preprint arXiv:2210.13382</i> .	William Merrill, Nikolaos Tsilivis, and Aman Shukla. 2023. A tale of two circuits: Grokking as competition of sparse and dense subnetworks. <i>arXiv preprint arXiv:2303.11873</i> .	1008 1009 1010 1011
961	Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2024. Inference-time intervention: Eliciting truthful answers from a language model. <i>Advances in Neural Information Processing Systems</i> , 36.	Mohamad Amin Mohamadi, Zhiyuan Li, Lei Wu, and Danica Sutherland. 2023. Grokking modular arithmetic can be explained by margin maximization. In <i>NeurIPS 2023 Workshop on Mathematics of Modern Machine Learning</i> .	1012 1013 1014 1015 1016
966	Tom Lieberum, Matthew Rahtz, János Kramár, Geoffrey Irving, Rohin Shah, and Vladimir Mikulik. 2023. Does circuit analysis interpretability scale? evidence from multiple choice capabilities in chinchilla. <i>arXiv preprint arXiv:2307.09458</i> .	Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever. 2021. Deep double descent: Where bigger models and more data hurt. <i>Journal of Statistical Mechanics: Theory and Experiment</i> , 2021(12):124003.	1017 1018 1019 1020 1021
971	Wenhao Liu, Xiaohua Wang, Muling Wu, Tianlong Li, Changze Lv, Zixuan Ling, Jianhao Zhu, Cenyuan Zhang, Xiaoqing Zheng, and Xuanjing Huang. 2023. Aligning large language models with human preferences through representation engineering. <i>arXiv preprint arXiv:2312.15997</i> .	Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. 2023a. Progress measures for grokking via mechanistic interpretability. <i>arXiv preprint arXiv:2301.05217</i> .	1022 1023 1024 1025
977	Ziming Liu, Ouail Kitouni, Niklas S Nolte, Eric Michaud, Max Tegmark, and Mike Williams. 2022a. Towards understanding grokking: An effective theory of representation learning. <i>Advances in Neural Information Processing Systems</i> , 35:34651–34663.	Neel Nanda, Andrew Lee, and Martin Wattenberg. 2023b. Emergent linear representations in world models of self-supervised sequence models. <i>arXiv preprint arXiv:2309.00941</i> .	1026 1027 1028 1029
982	Ziming Liu, Eric J Michaud, and Max Tegmark. 2022b. Omnigrok: Grokking beyond algorithmic data. <i>arXiv preprint arXiv:2210.01117</i> .	Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter. 2020a. <i>Zoom in: An introduction to circuits</i> . <i>Distill</i> . https://distill.pub/2020/circuits/zoom-in .	1030 1031 1032 1033
985	Scott M Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. <i>Advances in neural information processing systems</i> , 30.	Chris Olah, Nick Cammarata, Chelsea Voss, Ludwig Schubert, and Gabriel Goh. 2020b. <i>Naturally occurring equivariance in neural networks</i> . <i>Distill</i> . https://distill.pub/2020/circuits/equivariance .	1034 1035 1036 1037
988	Haoyan Luo and Lucia Specia. 2024. From understanding to utilization: A survey on explainability for large language models. <i>arXiv preprint arXiv:2401.12874</i> .	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. 2022. In-context learning and induction heads. <i>Transformer Circuits Thread</i> . https://transformer-circuits.pub/2022/in-context-learning-and-induction-heads/index.html .	1038 1039 1040 1041 1042 1043 1044 1045 1046 1047 1048
991	Kaifeng Lyu, Jikai Jin, Zhiyuan Li, Simon S Du, Jason D Lee, and Wei Hu. 2023. Dichotomy of early and late phase implicit biases can provably induce grokking. <i>arXiv preprint arXiv:2311.18817</i> .	OpenAI. 2023. <i>Gpt-4 technical report</i> .	1049

1050	Vardan Papyan, XY Han, and David L Donoho. 2020.	Vikrant Varma, Rohin Shah, Zachary Kenton, János	1105
1051	Prevalence of neural collapse during the terminal	Kramár, and Ramana Kumar. 2023. Explaining	1106
1052	phase of deep learning training. <i>Proceedings of</i>	grokking through circuit efficiency. <i>arXiv preprint</i>	1107
1053	<i>the National Academy of Sciences</i> , 117(40):24652–	<i>arXiv:2309.02390</i> .	1108
1054	24663.		
1055	Alethea Power, Yuri Burda, Harri Edwards, Igor	Kevin Wang, Alexandre Variengien, Arthur Conmy,	1109
1056	Babuschkin, and Vedant Misra. 2022. Grokking:	Buck Shlegeris, and Jacob Steinhardt. 2022. In-	1110
1057	Generalization beyond overfitting on small algorithmic	terpretability in the wild: a circuit for indirect ob-	1111
1058	datasets. <i>arXiv preprint arXiv:2201.02177</i> .	ject identification in gpt-2 small. <i>arXiv preprint</i>	1112
		<i>arXiv:2211.00593</i> .	1113
1059	Marco Tulio Ribeiro, Sameer Singh, and Carlos	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel,	1114
1060	Guestrin. 2016. "why should i trust you?" explaining	Barret Zoph, Sebastian Borgeaud, Dani Yogatama,	1115
1061	the predictions of any classifier. In <i>Proceedings of</i>	Maarten Bosma, Denny Zhou, Donald Metzler, et al.	1116
1062	<i>the 22nd ACM SIGKDD international conference on</i>	2022. Emergent abilities of large language models.	1117
1063	<i>knowledge discovery and data mining</i> , pages 1135–	<i>Transactions on Machine Learning Research</i> .	1118
1064	1144.		
1065	Adam Scherlis, Kshitij Sachan, Adam S Jermyn, Joe	Muling Wu, Wenhao Liu, Xiaohua Wang, Tianlong Li,	1119
1066	Benton, and Buck Shlegeris. 2022. Polysemantic-	Changze Lv, Zixuan Ling, Jianhao Zhu, Cenyuan	1120
1067	ity and capacity in neural networks. <i>arXiv preprint</i>	Zhang, Xiaoqing Zheng, and Xuanjing Huang.	1121
1068	<i>arXiv:2210.01892</i> .	2024a. Advancing parameter efficiency in fine-	1122
1069	Ludwig Schubert, Chelsea Voss, Nick Cam-	tuning via representation editing. <i>arXiv preprint</i>	1123
1070	marata, Gabriel Goh, and Chris Olah. 2021.	<i>arXiv:2402.15179</i> .	1124
1071	High-low frequency detectors . <i>Distill</i> .	Xuansheng Wu, Haiyan Zhao, Yaochen Zhu, Yucheng	1125
1072	https://distill.pub/2020/circuits/frequency-edges .	Shi, Fan Yang, Tianming Liu, Xiaoming Zhai, Wenlin	1126
1073	Lee Sharkey, Dan Braun, and beren. 2022. [Interim	Yao, Jundong Li, Mengnan Du, et al. 2024b. Usable	1127
1074	research report] Taking features out of superposition	xai: 10 strategies towards exploiting explainability in	1128
1075	with sparse autoencoders . Accessed 2024-01-23.	the llm era. <i>arXiv preprint arXiv:2403.08946</i> .	1129
1076	Niklas Stoehr, Mitchell Gordon, Chiyuan Zhang, and	Zhengxuan Wu, Aryaman Arora, Zheng Wang, Atti-	1130
1077	Owen Lewis. 2024. Localizing paragraph mem-	cus Geiger, Dan Jurafsky, Christopher D. Manning,	1131
1078	orization in language models. <i>arXiv preprint</i>	and Christopher Potts. 2024c. Reft: Representation	1132
1079	<i>arXiv:2403.19851</i> .	finetuning for language models .	1133
1080	Gemini Team, Rohan Anil, Sebastian Borgeaud,	Yi Yang, Hanyu Duan, Ahmed Abbasi, John P Lalor,	1134
1081	Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu,	and Kar Yan Tam. 2023. Bias a-head? analyzing bias	1135
1082	Radu Soricut, Johan Schalkwyk, Andrew M Dai,	in transformer-based language model attention heads.	1136
1083	Anja Hauth, et al. 2023. Gemini: a family of	<i>arXiv preprint arXiv:2311.10395</i> .	1137
1084	highly capable multimodal models. <i>arXiv preprint</i>	Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu,	1138
1085	<i>arXiv:2312.11805</i> .	Huiqi Deng, Hengyi Cai, Shuaiqiang Wang, Dawei	1139
1086	Vimal Thilak, Etai Littwin, Shuangfei Zhai, Omid	Yin, and Mengnan Du. 2023. Explainability for large	1140
1087	Saremi, Roni Paiss, and Joshua M. Susskind. 2022.	language models: A survey. <i>ACM Transactions on</i>	1141
1088	The slingshot mechanism: An empirical study of	<i>Intelligent Systems and Technology (TIST)</i> .	1142
1089	adaptive optimizers and the Grokking Phenomenon .	Ziqian Zhong, Ziming Liu, Max Tegmark, and Jacob	1143
1090	In <i>Has it Trained Yet? NeurIPS 2022 Workshop</i> .	Andreas. 2023. The clock and the pizza: Two stories	1144
1091	Eric Todd, Millicent L Li, Arnab Sen Sharma, Aaron	in mechanistic explanation of neural networks. <i>arXiv</i>	1145
1092	Mueller, Byron C Wallace, and David Bau. 2023.	<i>preprint arXiv:2306.17844</i> .	1146
1093	Function vectors in large language models. <i>arXiv</i>	Xuekai Zhu, Yao Fu, Bowen Zhou, and Zhouhan	1147
1094	<i>preprint arXiv:2310.15213</i> .	Lin. 2024. Critical data size of language mod-	1148
1095	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	els from a grokking perspective. <i>arXiv preprint</i>	1149
1096	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	<i>arXiv:2401.10463</i> .	1150
1097	Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti	Andy Zou, Long Phan, Sarah Chen, James Campbell,	1151
1098	Bhosale, et al. 2023. Llama 2: Open founda-	Phillip Guo, Richard Ren, Alexander Pan, Xuwang	1152
1099	tion and fine-tuned chat models. <i>arXiv preprint</i>	Yin, Mantas Mazeika, Ann-Kathrin Dombrowski,	1153
1100	<i>arXiv:2307.09288</i> .	et al. 2023. Representation engineering: A top-	1154
1101	Alex Turner, Lisa Thiergart, David Udell, Gavin Leech,	down approach to ai transparency. <i>arXiv preprint</i>	1155
1102	Ulisse Mini, and Monte MacDiarmid. 2023. Acti-	<i>arXiv:2310.01405</i> .	1156
1103	vation addition: Steering language models without		
1104	optimization. <i>arXiv preprint arXiv:2308.10248</i> .		

A Mechanistic Interpretability

Mechanistic interpretability refers to the process of zooming into neural networks to understand the underlying components and mechanisms that drive their behaviors, also known as reverse engineering (Olah et al., 2020a). Just as the microscope revealed the world of cells, looking inside neural networks provides a glimpse into rich inner structures of models. This approach diverges from conventional interpretability methods that aim to explain the overall behaviors through features, neural activations, data instances etc. Instead, it draws inspiration from other fields, such as neuroscience and biology, to investigate individual neurons and their connections. By tracking each neuron and weight, an intricate picture emerges on how neural networks operate through interconnected “circuits” that implement meaningful algorithms. On this delicate scale, neural networks become approachable systems rather than black boxes. Neurons play an understandable role and their circuits of connections implement factual relationships about the world. We can thus observe the step-by-step construction of high-level concepts, such as circle detectors, animal faces, cars, and logical operations (Olah et al., 2020a). In essence, zooming into the micro-level mechanics of LLMs enables deeper comprehension of their macro-level behaviors. Such mechanistic perspective represents a paradigm shift in interpretability towards unpacking the causal factors that drive model outputs.

A.1 Role in the General XAI Field

Mechanistic interpretability in XAI represents a paradigm shift towards a deeper and more fundamental understanding of deep neural network (DNN) models (Zhao et al., 2023).

- **Global versus Local Interpretation:** Mechanistic interpretability diverges from the traditional local focus of XAI, which concentrates on explaining specific predictions made by deep learning models, e.g., feature attribution techniques. Instead, it adopts a global approach, aiming to comprehend DNN models as a whole through the lens of high-level concepts and circuits.
- **Post-hoc Analysis versus Intrinsic Design:** Mechanistic interpretability aims to decipher the complexities inherent in pre-trained DNN models in a post-hoc way. This contrasts with efforts to create models that are mechanistically interpretable by design (Friedman et al., 2023).

- **Model-Specific versus Model-Agnostic:** Unlike some XAI methods such as LIME (Ribeiro et al., 2016) and SHAP (Lundberg and Lee, 2017), which are model-agnostic, mechanistic interpretability is a model-specific explanation. It requires tailor-made designs for each distinct LLM, analyzing their unique characteristics.
- **White-box versus Black-box:** Mechanistic interpretability aligns with white-box analysis, requiring direct access to a model’s internal parameters and activations. This is in contrast to black-box XAI tools such as LIME and SHAP, which operate solely based on the model’s inputs and outputs.

In summary, mechanistic interpretability in XAI is a critical approach to gain a profound understanding of DNN models. It emphasizes a **global** and **post-hoc** perspective, focusing on **model-specific**, **white-box** analysis to decipher the inner workings and intrinsic logic of complex AI systems. This approach is pivotal to advance transparency and build trust for LLMs, especially in high-stake scenarios where grasping “why” behind AI systems is as crucial as the decisions themselves.

A.2 Why Mechanistic Interpretability?

The question naturally arises: *Why has XAI research on LLMs moved towards the more specialized domain in mechanistic interpretability?* Exploring this shift can shed light on the evolving needs and challenges in this field. In this section, we attempt to look through several factors that we believe have played a major role in steering the shift.

Alignment Requirement. In the age of LLMs, the standards for model performance have become more rigorous, not just in terms of accuracy but also in addressing crucial social concerns like dishonesty and fairness. Under this circumstance, the challenge of aligning LLMs with our values and expectations has become a pressing concern, one that demands a deep understanding and effective control of these models. To tackle these challenges, mechanistic interpretability stands out as a promising approach, offering a way to understand the underlying workings of these models.

Understanding Reasoning Capability. The field of XAI in machine learning has made significant progress with techniques designed to provide valuable insights to end users, such as feature attributions (Ribeiro et al., 2016) and example-based

1257	explanations (Koh and Liang, 2017). These techniques have been proven to be quite effective in computer vision tasks, where the demands for complex alignment were less strict. However, as LLMs become more sophisticated, their reasoning capability has transformed from mere pattern recognition to a form of complex, human-like cognition. This advancement in LLMs’ reasoning abilities renders traditional XAI methods obsolete and less competent in interpreting their behaviors.	1307
1258		
1259		
1260		
1261		
1262		
1263		
1264		
1265		
1266		
1267	Understanding Inner Working of LLMs. Moreover, alongside the strong reasoning abilities of LLMs, their notorious deep and intricate architectures are raising new concerns. Since the inner workings of these models are multifaceted and intricate, new challenges in explaining models at the structure level have emerged. Conventional global interpretability techniques, which are adept at uncovering the high-level knowledge acquired in different components of models, fall short when providing sights into the functions and the evolution of knowledge within these models. This issue is further confounded as LLMs scale aggressively, making neuron-level and layer-level insights increasingly insufficient. This complexity highlights the urgent need for innovative approaches that enable us to zoom in models and provide more in-depth, mechanistic understandings at various levels.	
1268		
1269		
1270		
1271		
1272		
1273		
1274		
1275		
1276		
1277		
1278		
1279		
1280		
1281		
1282		
1283		
1284		
1285	Alternatively, mechanistic interpretability aims to unravel the inner workings of LLMs, providing insights into the “how” and “why” behind their decision-making processes. Specifically, mechanistic interpretability focuses on the causal relationships and underlying mechanisms within models. This not only is more suited to the advanced nature of LLMs, but is also crucial to ensure transparency, trust, and reliability in their applications.	
1286		
1287		
1288		
1289		
1290		
1291		
1292		
1293		
1294	A.3 Mechanistic Interpretability Theories	
1295	Most of the current work on mechanistic interpretability is based on vision models, and some recent work has begun to investigate Transformer models. In this section, we introduce some core concepts and pivotal phenomena in the field of mechanistic interpretability. Since LLMs are too complicated to analyze locally, simple yet artificial models are purposely designed to investigate their characteristics and internal mechanisms. We will introduce the main assumptions and observations made under this setting, including <i>circuits</i> , <i>induction heads</i> , <i>superposition</i> , <i>polysemanticity</i> ,	
1296		
1297		
1298		
1299		
1300		
1301		
1302		
1303		
1304		
1305		
1306		
	and <i>monosemanticity</i> .	1307
	B Mechanistic Interpretability v.s. Representation Engineering	1308
		1309
	In this section, we provide further discussion on different explanation scales of two techniques. Further, we provide our understanding towards their	1310
		1311
		1312
	Explanability Scale. These two techniques explain LLMs at opposite scales.	1313
		1314
	• Micro-scale: Mechanistic interpretability focuses on dissecting the intricate inner workings of LLMs at the neuron and circuit levels. It aims at illustrating how models function and process certain tasks with subnetworks.	1315
		1316
		1317
		1318
		1319
	• Macro-scale: Representation engineering places representations, rather than neurons or circuits, as the central unit of analysis. The goal is to understand and control cognitive behaviors by studying their manifestations in learned representation spaces.	1320
		1321
		1322
		1323
		1324
		1325
	Roles in XAI. Two techniques are providing multifaceted perspectives in the field of XAI. Representation engineering embodies how well embeddings capture the essence of data. Good representations are crucial to making accurate predictions. The visualization of representation can also implicitly demonstrate the quality of learning. On the other hand, through the lens of mechanistic interpretability, we can investigate relations between models’ abilities like generalization and training dynamics. Examining the evolution of models from initialization to generalization, we can reveal characteristics of generalization, such as sparsity. These characteristics could serve as benchmarks for what constitutes “good learning”. Apart from that, mechanistic interpretability is known to explain individual functional components and potentially improve model performance in the future.	1326
		1327
		1328
		1329
		1330
		1331
		1332
		1333
		1334
		1335
		1336
		1337
		1338
		1339
		1340
		1341
		1342
		1343
	Potential to Alignment. At the current stage, both techniques have witnessed preliminary applications in LLM alignment. Mechanistic interpretability plays a crucial role in locating knowledge or biases at the level of attention heads, while representation engineering is primarily employed in targeting undesired behaviors at the level of layers. Despite the distinct focus of each approach within models, both have proven effective in identifying biases and highlighting practical steps for improvement. However, they are still incompetent in uncovering rudimentary causes behind these biases.	1344
		1345
		1346
		1347
		1348
		1349
		1350
		1351
		1352
		1353
		1354
		1355

C Research Challenges

In this section, we outline the research challenges that deserve future efforts from the community.

C.1 The Validity of Existing Theories

While theories that attempt to explain the mechanisms behind the capabilities of transformer models are promising, their empirical support is not definitive. For example, understanding induction heads is key to explain transformer models because they are recognized as foundations for in-context learning abilities. However, as highlighted by [Olsson et al. \(2022\)](#), defining what exactly an induction head is remains somewhat elusive. Similarly, the proposition of a mathematical framework to explain circuits inside a simplified network opens up an interesting avenue of research. Although [Lieberum et al. \(2023\)](#) conclude that circuit analysis is feasible on LLMs, this theoretical framework has not been thoroughly tested with empirical studies. Besides, these theoretical models rely on idealized assumptions such as superposition and often lack ground truth. This further complicates the task of validating these theories.

C.2 The Curse of Dimensionality

Another challenge is that the parameters we can explain are much less than a third of all parameters in LLMs. These explanations focus on components of attention heads, and although dictionary learning helps to partially understand polysemantic neurons, there is still a vast territory that remains unexplored. The rest majority of these model parameters are tied to MLP layers, which are notoriously difficult to fully comprehend ([Olsson et al., 2022](#)). Their compositions are more complicated than those of attention layers, making the analysis process considerably more arduous and perplexing. For instance, [Geva et al. \(2021\)](#) believes that the output of MLPs is a composition of memories including textual patterns and output distributions. [Meng et al. \(2022\)](#) attempt to modify MLPs to edit factual knowledge in LLMs. However, the effectiveness of editing has been put into doubt by another work ([Hase et al., 2023](#)).

C.3 Evaluation of Concepts and Circuits

A key challenge in mechanistic interpretability is validating and ensuring the accuracy of proposed conceptual explanations and functional circuits. Unlike straightforward metrics in machine learn-

ing to assess predictions, interpretation evaluation lacks clear ground truth. As noted in [Chan et al. \(2022\)](#), we are short of tools to measure the degree to which explanations interpret the relevant phenomenon. Existing ad-hoc ablation methods, i.e. standard zero and mean ablations, are neither universal nor scalable. Exploring measurements from various angles, such as causal scrubbing, which involves randomly sampling inputs to patch activations without disturbing the input distribution, could enrich our evaluation dimensions. Moreover, manual inspections are challenging in identifying circuits within LLMs. Our understanding of automatically discovering these circuits is still developing ([Conmy et al., 2023](#)). Heterogeneous mechanistic explanations can be generated in networks trained on simple tasks such as modular additions ([Zhong et al., 2023](#)). This suggests that even in seemingly simple scenarios, the outcomes of circuit analysis can be uncertain. Additionally, different models learned on similar tasks might learn same family of circuits, but the precise circuits learned by individual networks are not the same ([Chughtai et al., 2023](#)).

C.4 Conflicted Explanations

There are other observations in understanding observations, such as neural collapse ([Papayan et al., 2020](#)), yet there is a notable gap in understanding how these observations are interconnected. The root causes of these observations often lead to conflicting viewpoints. For example, [Gromov \(2023\)](#) suggests that grokking might be triggered by the learning of a new feature. Unfortunately, the leap in generalization could be too subtle to notice without a hierarchical model ([Gromov, 2023](#)). On the other hand, there is some debate around linking grokking with generalization ([Levi et al., 2023](#)). Moreover, a significant limitation of these studies is their focus on arithmetic datasets instead of real-world datasets, which casts doubt on how broadly these findings can be applied. To fully understand the generalization of models and reconcile these conflicting views, a holistic examination of how these observations relate to each other and their impact on training dynamics across models is essential.