

Towards Uncovering How Large Language Models Work: An Interpretability Perspective

Haiyan Zhao¹, Fan Yang², Bo Shen¹, Ali Payani³, Himabindu Lakkaraju⁴, Mengnan Du¹

¹New Jersey Institute of Technology ²Wake Forest University

³Cisco Research ⁴Harvard University

{hz54, bo.shen, mengnan.du}@njit.edu, yangfan@wfu.edu, apayani@cisco.com, hlakkaraju@hbs.edu

Abstract

Large language models (LLMs) have shown remarkable performance in tackling natural language tasks, yet the internal mechanisms that enable their impressive generalization and reasoning abilities remain opaque. This lack of transparency presents significant challenges in fundamentally eliminating undesirable behaviors such as hallucinations and toxicity, 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 what knowledge is embedded in LLM representations by leveraging probing techniques and representation engineering. Additionally, we investigate the training dynamics to explore models' generalization abilities through grokking and memorization. Finally, we explore how the insights gained from these explanations can further enhance LLM performance through model editing, improve efficiency through pruning, and better align with human values.

1 Introduction

Large language models (LLMs) have led to tremendous advancements in natural 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 spread misinformation, promote biased or harmful views, and even compromise the safety and wellbeing of society. Therefore, there is an urgent need to fully understand the inner workings of LLMs to address these issues. Gaining insights into how these models operate is a crucial first step towards developing robust safeguards and ensuring

their responsible deployment, although our understanding is still in the very early stages.

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 (MI), is promising in explaining models' internals at the level of neurons, circuits and attention heads with activations and weights [Saphra and Wiegrefe, 2024]. Second, we examine what knowledge is encoded internally in intermediate representations. To this end, representation engineering (RE) is adopted to explain specific behaviors of models, such as dishonesty, by analyzing hidden/intermediate representations. Specifically, RE focuses on identifying patterns for certain behaviors through probing-based methods and employs them to mitigate undesired behaviors by steering models at inference time [Zou *et al.*, 2023]. Third, we inspect the model training process to understand the development of generalization abilities. 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 alignment with human values.

Our work differs from existing survey articles on the explainability of LLMs [Zhao *et al.*, 2024; Ferrando *et al.*, 2024], which either summarize explainability techniques or discuss their utilities. In contrast, our goal is to review recent studies that provide insights into trained LLMs and their dynamic training processes. By focusing on model components, representation space, and training processes, we aim to synthesize a line of work particularly focused on uncovering how LLMs function and identify the factors that contribute to their reasoning abilities. Beyond reviewing insights on the inner working mechanisms of LLMs, we further explore how these insights are employed to enhance model performance and benefit humanity.

2 Transformer-based LLMs

Before delving into the techniques and insights about how LLMs work, it is essential to establish a foundational understanding of the LLM architecture. This section introduces key components of the decoder-only LLMs, which have been the

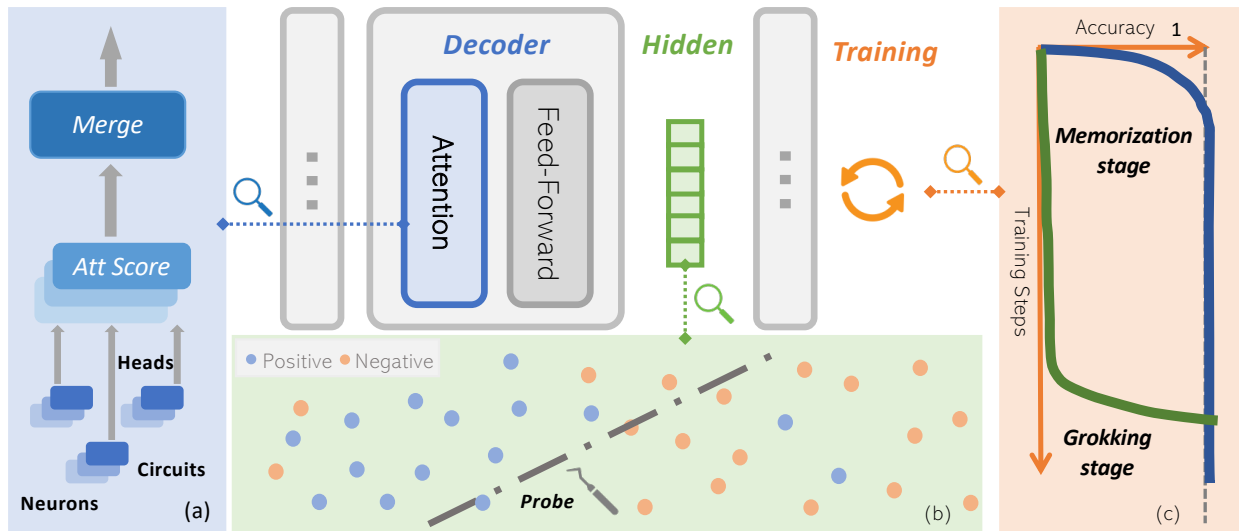


Figure 1: In this work, we review existing progress on how LLMs work, including: (a) how knowledge is encoded within model components; (b) what knowledge is encoded in intermediate representations; and (c) how generalization abilities are achieved during the training process.

84 dominant architecture for LLMs in recent years [Vaswani *et*
85 *al.*, 2017].

86 2.1 Decoder-only LLM Architecture

87 Decoder-only (autoregressive) architecture has been widely
88 used in popular LLMs such as GPT-3.5, GPT-4, Gemma-2,
89 and Llama-3.1. A decoder-only Transformer model typically
90 consists of the following main components (see Figure 2 (a)):

- 91 • *Input Embedding Layer*: This layer converts input tokens
92 (words or subwords) into embeddings/representations.
93 Since Transformers don't inherently understand sequence
94 order, positional encodings are added to the input embed-
95 dings to provide position information for each token.
- 96 • *Multiple Transformer Layers*: The core of the model con-
97 sists of a stack of identical Transformer layers. Each layer
98 contains: a) Multi-head Self-attention: This allows the
99 model to attend to different parts of the input sequence
100 when processing each token. b) Feed-forward Neural Net-
101 work: A simple fully connected network. c) Layer Nor-
102 malization: Applied after each sub-layer to stabilize the
103 learning process. d) Residual Connections: These skip-
104 connections help in training deeper networks by allowing
105 gradients to flow more easily through the model.
- 106 • *Output Layer*: The final layer typically projects the rep-
107 resentations onto the vocabulary space, producing logits for
108 each token in the vocabulary.

109 In autoregressive models, input sequences are processed
110 token by token, with each token attending to all previous to-
111 kens in the sequence via the self-attention mechanism. Dur-
112 ing training, the model learns to predict the next token given
113 the previous tokens. At inference time, the model generates
114 each next token by repeatedly sampling from its predicted
115 probability distribution and feeding this predicted token back
116 as input for the next step.

117 2.2 Intermediate Representations

118 Intermediate representations, also known as hidden represen-
119 tations, are a crucial component in understanding the inner
120 workings of LLMs. The granularity of these representations
121 can vary from the attention head level, where individual heads
122 may specialize in specific patterns, to the layer level, where
123 each layer's output represents a higher-level abstraction of the
124 input. In this work, we refer to intermediate representations
125 as the outputs from "Transformer block" layers (as illustrated
126 in Figure 2 (a)). These layers include both outputs from at-
127 tention blocks capturing contextual information through self-
128 attention mechanisms, and outputs from feed-forward layers
129 which further transform and process this information.

130 Understanding these representations is vital as they offer
131 insights into how models transform information at different
132 layers and reveal what types of knowledge or patterns models
133 have learned with probing tools [Belinkov, 2022]. They have
134 also been extensively used in representation engineering to
135 steer model behaviors. In the following sections, particularly
136 in Section 4, we will explore how these representations are
137 employed to interpret world knowledge, factual information,
138 and even unintended biases learned in models.

139 3 How is Knowledge Encoded in Model Architectures?

140 LLMs' emergent abilities are remarkable, but the underly-
141 ing mechanisms enabling models to learn vast amounts of
142 knowledge remain unclear. To fully understand LLMs, recent
143 studies have been focused on utilizing MI to reverse engineer
144 LLMs at a more subtle level, such as neurons and attention
145 heads. MI emphasizes understanding the causal mechanisms
146 within models, namely the relationship between model com-
147 ponents and their behaviors [Saphra and Wiegrefe, 2024].
148 In this section, we review studies analyzing functionalities of
149

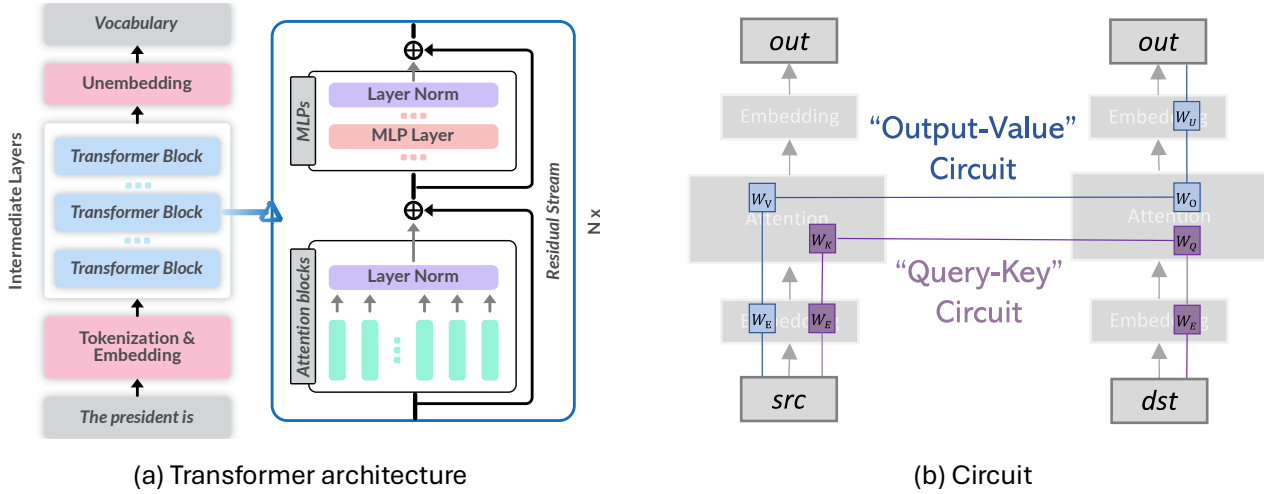


Figure 2: (a) The Decoder-only LLM Architecture. (b) An illustration of a Transformer circuit. *src* and *dst* denote previous tokens and the next predicted token respectively. In the toy model, input tokens are first sent to the embedding layer, then pass through the attention block. The intermediate output is further passed into unembedding layer to decode the destination token.

150 neurons, circuits, and attention heads with toy models and
 151 summarize how LLMs process information and make decisions.
 152 Notably, foundational discoveries in these simplified
 153 models have successfully generalized to explain real-world
 154 LLM behaviors, validating concepts like induction heads and
 155 driving the development of Sparse Autoencoders (SAEs).

156 3.1 Neurons

157 Neurons in LLMs are the fundamental units of computation.
 158 They receive data from inputs and produce output through ac-
 159 tivation functions. When a neuron’s activation value is high
 160 after processing its inputs through the activation function, we
 161 consider the neuron to be “activated”. Neurons can be acti-
 162 vating on multiple unrelated terms (polysemantic) or a single
 163 term (monosemantic) [Olah *et al.*, 2020].

164 Polysemanticity

165 Polysemanticity presents a challenge in mechanistically un-
 166 derstanding how models operate. *Superposition* appears to
 167 be a key cause of polysemantic nature [Olah *et al.*, 2020].
 168 Superposition describes the phenomenon where a feature can
 169 be spread across multiple neurons, while a single neuron can
 170 simultaneously be involved with multiple features.

171 Some researchers suggest that superposition originates
 172 from an excessive number of features compared to neu-
 173 rons [Olah *et al.*, 2020]. In ReLU networks, Elhage *et al.*
 174 [2022] find that superposition enables the representa-
 175 tion of additional features by tolerating some interference.
 176 In certain cases, it even facilitates computations such as
 177 the absolute value function and logical AND [Hänni *et al.*,
 178 2024]. Others argue that polysemanticity can arise inciden-
 179 tally due to multiple training factors including regulariza-
 180 tion and neural noise [Lecomte *et al.*, 2024]. Experiments
 181 show that feature collisions, introduced through random ini-
 182 tialization, consistently result in polysemantic neurons, even
 183 when the number of neurons exceeds the number of fea-
 184 tures [Lecomte *et al.*, 2024]. Another study examines poly-

185 semanticity through the lens of the “feature capacity”, which
 186 denotes the fraction of embedding dimensions consumed by
 187 a feature [Scherlis *et al.*, 2022]. Through the analysis of
 188 toy models, this work indicates that features are represented
 189 based on their importance in reducing loss. More impor-
 190 tant features are allocated their own dimensions, while less
 191 critical ones tend to be polysemantic [Scherlis *et al.*, 2022;
 192 Gurnee *et al.*, 2023]. Gurnee *et al.* [2023] show that early
 193 layers tend to represent many features in superposition, while
 194 middle layers include dedicated neurons to represent high-
 195 level features.

196 Monosemanticity

197 Monosemantic neurons are much easier to interpret. Investi-
 198 gating the factors that enhance monosemanticity is essen-
 199 tial for model interpretation. Jermyn *et al.* [2022] use toy
 200 models to reveal that changing the loss minima could im-
 201 prove monosemanticity. Such loss minimum usually coexists
 202 with negative biases. However, in practice building a purely
 203 monosemantic model is infeasible due to the unmanageable
 204 loss. Another line of studies seeks to disentangle superpo-
 205 sition to reach a monosemantic understanding using sparse
 206 autoencoder (SAE). Instead of focusing on a single neuron,
 207 SAE considers layerwise neurons [Cunningham *et al.*, 2024].
 208 It learns sparse activation directions representing monose-
 209 mantic feature through dictionary learning. However, the ef-
 210 fectiveness of this approach remains unclear. Templeton *et al.*
 211 [2024] find that SAEs produce features representing both
 212 low-level and high-level concepts, such as code errors and
 213 hateful bias-related features. Further, Ferrando *et al.* [2025]
 214 show that SAEs identify features for entity recognition.

215 3.2 Circuits

216 Circuit analysis was originally proposed to reverse engineer
 217 vision models [Olah *et al.*, 2020]. Taking car classification as
 218 an example, researchers suggest that some fundamental fea-
 219 tures such as edge detectors are learned in early layers. These

220 features are then combined through weights to form circuit
221 units like wheel detectors in later layers. This viewpoint is
222 supported by evidence from a few interpretable circuits per-
223 forming specific functions such as curve detection and sym-
224 metric transformations of basic features, including copying,
225 coloring, and rotation.

226 However, Transformer models present new challenges with
227 their unique architecture. To address these challenges, a
228 mathematical framework specifically for *transformer circuits*
229 has been proposed [Elhage *et al.*, 2021]. It simplifies the
230 complex architecture of LLMs by focusing on toy models
231 composed of attention blocks with no more than two layers.
232 Typically, two-layer toy models ensure preserving all neces-
233 sary components of transformer models including input em-
234 bedding, residual stream, attention layers, and output embed-
235 dings, while minimizing architectural complexity.

236 Attention layers communicate by reading information from
237 the residual stream and then writing their output back. Each
238 attention head operates independently and in parallel, with
239 input embedding matrix W_E and output unembedding matrix
240 W_U . These heads consist of key, query, output, and value
241 weights, represented as W_K , W_Q , W_O and W_V . There are
242 two types of circuits: i) “query-key” (QK) circuits formed
243 by $W_Q^T W_K$; ii) “output-value” (OV) circuits composed of
244 $W_U W_{OV}^h W_E$, where h denotes attention head, as shown in
245 Figure 2 (b). The QK circuits are essential for models to
246 recall and retrieve information from earlier context, deter-
247 mining which previous token to copy information from. The
248 computed attention score $W_E^T W_{QK}^h W_E$ indicates how much
249 a destination token attends to a source token. The OV cir-
250 cuits, in turn, determine how the source token influences the
251 output logits [Elhage *et al.*, 2021].

252 By contextualizing the weights in circuits, research shows
253 that circuits are essentially linear or bilinear functions on to-
254 kens. Specifically, Transformers with no layers can model
255 bigram statistics, predicting the next token from the source
256 token. Adding one layer allows the model to capture both
257 bigram and “skip-trigram” patterns. Interestingly, with two
258 layers, Transformer models give rise to “*induction head*” in
259 the second layer and beyond (Section 3.3). These heads are
260 typically composed of heads from their previous layer and are
261 useful in predicting the next token [Elhage *et al.*, 2021].

262 Beyond the above-mentioned theoretical analysis on toy
263 models, circuits implementing specific functions have been
264 identified in real-world LLMs. For example, a circuit com-
265 posed of 26 attention heads in GPT-2 small has been found
266 to enable indirect object identification tasks by transmitting
267 information from name tokens to the final outputs. Other re-
268 search has identified circuits in GPT-2 that perform greater-
269 than computations using a set of MLPs.

270 3.3 Attention heads

271 A special type of attention head called *induction head* is
272 considered critical in enabling in-context learning abilities
273 within LLMs [Brown *et al.*, 2020], due to their co-occurrence
274 and causal relations [Olsson *et al.*, 2022; Chan *et al.*,
275 2022]. Induction heads are circuits that complete patterns
276 through prefix matching and copying previously occurred se-
277 quences [Olsson *et al.*, 2022]. They comprise two heads: the

278 first attention head from the previous layer attends to previous
279 tokens that are followed by the current token, achieving prefix
280 matching and providing the attend-to token (the token follow-
281 ing the current token). The second head copies the attend-to
282 token and increases its output logits. Specifically, this means
283 that if models have encountered patterns such as “[A*][B*]”
284 given the current token “[A]”, they can predict “[B]”. Be-
285 yond the simple example, long prefix matching involving
286 three consecutive tokens has also been observed [Chan *et al.*,
287 2022]. Consequently, layers with induction heads possess
288 more sophisticated in-context learning abilities than simple
289 copying. However, this theory requires further research on
290 real-world LLMs.

291 Other functional heads have also been found. For exam-
292 ple, Gould *et al.* [Gould *et al.*, 2024] discovered successor
293 heads that increment tokens like numbers in a natural order.
294 Another study reveals how factual associations are stored and
295 extracted within LLMs [Geva *et al.*, 2023]. A fact association
296 consists of a subject and an attribute. The subject is enriched
297 with subject-related attributes at the early MLP sublayers,
298 which is then propagated to the prediction. The prediction
299 representation “queries” the enriched subject to extract at-
300 tributes via attention heads. Moreover, Chughtai *et al.* [2024]
301 suggest that there are also mixed heads containing informa-
302 tion from both subject and relation. The subject heads, mixed
303 heads, and relation heads work additively to elicit the outputs.

304 4 What Knowledge is Encoded in 305 Intermediate Representations?

306 In this section, we present an in-depth review of the knowl-
307 edge encoded by *LLM representations*, including world
308 knowledge and factual knowledge captured by models. We
309 examine how factors such as layer depth and model scale im-
310 pact the encoding process.

311 4.1 Probing World and Factual Knowledge

312 Probing techniques play a crucial role in revealing insights
313 into world knowledge and factual knowledge encoded within
314 models. Specifically, they identify key directions within
315 the representation space that are essential for understanding
316 model behaviors and learned knowledge.

317 Recent studies have demonstrated that LLMs can learn
318 world models and encode them in their representations for
319 specific tasks. One study successfully uses non-linear probes
320 to uncover Othello board state representations within mod-
321 els [Li *et al.*, 2023]. It demonstrates models’ ability to track
322 board states and make predictions without explicit instruc-
323 tion. Furthermore, Nanda *et al.* [2023b] find that linear repre-
324 sentation structures can also perform well on predictions by
325 probing the board state at each timestamp through “my color”
326 and “opponent’s color”. These findings reveal how models
327 naturally perceive the world, which may differ from human
328 perception. Additionally, by analyzing representations of
329 spatial datasets, Gurnee *et al.* [2024] demonstrate the model’s
330 ability to learn linear representations of space and time across
331 multiple levels. LLMs are also capable of encoding factual
332 knowledge. Marks *et al.* [2024] craft self-curated true/false
333 datasets to study the geometry of representations of true/false

334 statements derived from models' residual stream. The data
335 are linearly separable under principal component analysis
336 (PCA). The identified truth directions are used to mediate
337 models' dishonest behaviors locally. Another research direc-
338 tion explores toxicity-related vectors within MLPs through
339 singular value decomposition (SVD). The identified dimen-
340 sions can be subtracted to achieve efficient mitigation [Lee *et*
341 *al.*, 2024].

342 Function vectors have also been discovered within LLMs'
343 attention heads of LLMs, triggering specific task execution
344 across diverse inputs. For example, Todd *et al.* [2024] find
345 that these function vectors appear in various in-context learn-
346 ing tasks and can execute related tasks even with zero-shot
347 inputs. Additionally, causal interventions at the neuron level
348 help identify individual neurons encoding spatial coordinates
349 and temporal information [Gurnee and Tegmark, 2024].

350 Lastly, representations has been found to be associated
351 with undesirable LLM behaviors, including dishonesty, tox-
352 icity, hallucination, safety concerns. Research has identified
353 specific directions in the representation space that contribute
354 to these behaviors. These directional vectors can be added to
355 the representation space during inference time to guide model
356 behavior without modifying the model's parameters [Zou *et*
357 *al.*, 2023; Li *et al.*, 2024; Azaria and Mitchell, 2023; Her-
358 nandez *et al.*, 2024]. This technique of modifying model be-
359 havior by manipulating the representation space is known as
360 representation steering or representation intervention, which
361 provides a lightweight and flexible approach to controlling
362 LLM outputs without the need for model retraining or pa-
363 rameter updates.

364 4.2 Role of Layer Depth and Model Scale

365 The influence of layer depth on representation analysis has
366 emerged as a compelling research direction. A particularly
367 intriguing finding is the concentration of well-learned knowl-
368 edge in the middle layers of LLMs. For example, Gurnee
369 *et al.* [2024] demonstrate that space and time representations
370 achieve optimal quality within the first half of layers across
371 LLMs. Additionally, function vectors with strong causal ef-
372 fects are predominantly found in the middle layers of LLMs,
373 while their effects become negligible in deeper layers [Todd
374 *et al.*, 2024]. Furthermore, research indicates that simpler
375 tasks are mastered in earlier layers, while complex tasks re-
376 quire deeper layers for effective learning [Jin *et al.*, 2025;
377 Ju *et al.*, 2024]. However, the underlying mechanisms of this
378 phenomenon remain unclear.

379 As predicted by scaling laws, model capabilities increase
380 as models grow larger. Gurnee *et al.* [2024] show that space
381 and time representations become more precise with model
382 scaling. However, internal mechanisms driving this improved
383 performance during scaling remain poorly understood.

384 5 How is Generalization Ability Achieved 385 During Training?

386 In this section, we examine how models develop generaliza-
387 tion ability during the training process. We focus on two im-
388 portant phenomena associated with generalization: grokking
389 and memorization.

5.1 Understanding Grokking

390 *Grokking* describes a phenomenon where models sud- 391
392 denly demonstrate improved validation accuracy after a pe- 393
394 riod of severe overfitting in over-parameterized neural net- 394
395 works [Power *et al.*, 2022]. This sudden surge in validation 395
396 accuracy is generally interpreted as the acquisition of gener- 396
397 alization ability.

A Data Perspective

398 Experiments on a two-layer decoder-only Transformer net- 398
399 work have demonstrated that grokking is influenced by mul- 399
400 tiple factors, including data characteristics, representations, 400
401 and regularization. Research shows that smaller datasets re- 401
402 quire more optimization steps for grokking to occur [Power 402
403 *et al.*, 2022], while increasing the number of samples reduces 403
404 the steps needed for generalization [Zhu *et al.*, 2024]. The 404
405 minimal data requirement for grokking is tied to the thresh- 405
406 old number of data points needed to learn robust represen- 406
407 tations. However, the massive datasets used in LLMs make 407
408 grokking phenomena less observable [Zhu *et al.*, 2024]. Ad- 408
409 ditionally, regularization techniques can accelerate the onset 409
410 of grokking, with weight decay proving particularly effective 410
411 in enhancing generalization capabilities. Wang *et al.* [2024] 411
412 have shown that Transformers can develop implicit reasoning 412
413 abilities exclusively through grokking. Moreover, their re- 413
414 search reveals that data distribution plays a more crucial role 414
415 in achieving grokking than data size alone.

Weight Norms

416 When analyzing models' final layer weight norms during 416
417 late-stage training, researchers have observed a phenomenon 417
418 called the *slingshot mechanism*, which occurs only in the ab- 418
419 sence of regularization. This mechanism is characterized by 419
420 the simultaneous occurrence of norm growth and training loss 420
421 spikes, with grokking observed at the onset of these sling- 421
422 shots [Thilak *et al.*, 2022]. The manifestation of both the 422
423 slingshot effect and grokking can be controlled by modifying 423
424 optimizer parameters, particularly when using certain adap- 424
425 tive optimizers. However, the universality of this observation 425
426 across different scenarios remains uncertain.

427 Another significant finding is the *LU mechanism*, which 427
428 describes the dynamics between loss and weight norms [Liu 428
429 *et al.*, 2023]. In algorithmic datasets, researchers have iden- 429
430 tified an L-shaped training loss curve and a U-shaped test 430
431 loss reduction relative to weight norms, suggesting an optimal 431
432 range for weight norm initialization. However, this pattern 432
433 doesn't readily apply to real-world machine learning tasks, 433
434 which typically require large initialization values and mini- 434
435 mal weight decay. However, Fan *et al.* [2024] argue that fea- 435
436 ture ranks and linear probing accuracy may serve as more re- 436
437 liable indicators of phase transition compared to weight norm 437
438 measurements.

Test Loss

440 *Double descent* describes a pattern in which a model's test 440
441 accuracy at the log level follows a distinctive trajectory: ini- 441
442 tial improvement, followed by a decline due to overfitting, 442
443 and finally a secondary increase as generalization abilities de- 443
444 velop [Nakkiran *et al.*, 2020]. This pattern becomes more 444
445

Leveraging LLM Understanding for Improvement

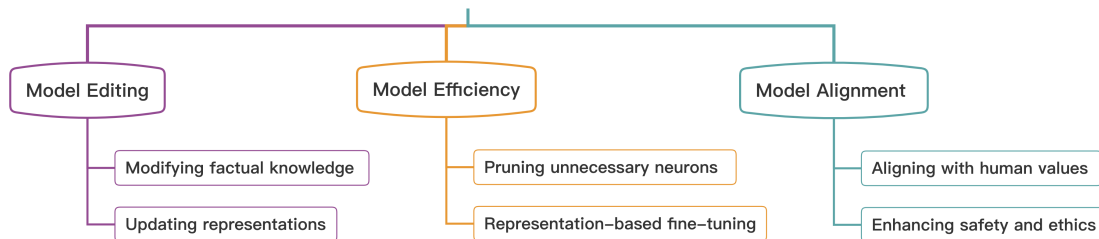


Figure 3: Leveraging understanding of LLMs for targeted improvements. We illustrate three key areas: Model Editing (modifying factual knowledge and updating representations), Model Pruning (improving inference efficiency and reducing model size), and Model Alignment (aligning with human values and enhancing safety and ethics).

446 pronounced when examining test loss. Researchers have de- 488
447 veloped a unified framework that integrates grokking and 489
448 double descent, viewing them as different manifestations of
449 the same underlying process [Davies *et al.*, 2022]. The frame-
450 work suggests that the transition to generalization stems from
451 slower pattern learning, a perspective supported by subse-
452 quent research [Kumar *et al.*, 2024]. This transitional phe-
453 nomenon has been observed to manifest at both the epoch
454 and model levels.

455 5.2 Memorization

456 *Memorization* refers to the phenomenon where models rely 497
457 on statistical features rather than causal relationships for pre- 498
458 dictions. In LLMs, generalization abilities take precedence 499
459 over memorization, particularly when evaluating a model’s 500
460 reasoning capabilities. Understanding memorization is there- 501
461 fore crucial for improving model performance. 502

462 Studies of two-layer neural models have shown that mem- 503
463 orization can exist alongside generalization. Memorization 504
464 can be reduced through either neuron pruning or regulariza- 505
465 tion techniques. While different regularization methods may 506
466 have distinct learning objectives, they all contribute to the 507
467 development of better representations. The terminal phase 508
468 of training encompasses two distinct stages: i) the grokking 509
469 process and ii) the decay of memorization learning [Doshi 510
470 *et al.*, 2024]. However, the underlying mechanisms driving 511
471 this process remain unclear. Furthermore, the hypothesis that 512
472 regularization is central to this process has been challenged, 513
473 particularly given observations of grokking occurring without 514
474 regularization [Kumar *et al.*, 2024]. 515

475 Interestingly, recent research hypothesizes that memoriza- 516
476 tion forms an integral phase of grokking [Nanda *et al.*, 517
477 2023a]. This study identifies three distinct stages in the 518
478 grokking process: memorization, circuit formation, and 519
479 memorization cleanup. Through analysis of model weights, 520
480 the researchers identified an algorithm using Discrete Fourier 521
481 Transforms and trigonometric identities to achieve modular 522
482 addition. Xie *et al.* [2024] illustrate the coexistence of 523
483 heavy memorization and improved generalization after fine- 524
484 tuning. Wang *et al.* [2025] observe that as model size in- 525
485 creases, factual question answering exhibits increased memo- 526
486 rization, while machine translation and reasoning tasks show 527
487 enhanced generalization. Moreover, Menta *et al.* [2025] argue

that memorization primarily occurs in deeper attention heads, 488
while earlier layers play a crucial role in generalizations. 489

490 6 How to Make Use of The Insights?

491 Building on the insights from previous sections, we examine 492
492 how our in-depth understanding of LLMs can be leveraged 493
493 to enhance their performance through editing, improve effi- 494
494 ciency via pruning, and better align them with human values 495
495 and preferences (Figure 3).

496 6.1 Model Editing for Better Performance

497 Recent research has made significant strides in model edit- 498
498 ing. Meng *et al.* [2022] demonstrate the ability to edit fac- 499
499 tual knowledge by modifying specific MLP neuron weights, 500
500 successfully altering neural computations related to factual 501
501 recall. Stoehr *et al.* [2024] further reveal that memorized 502
502 paragraphs can be identified through high-gradient weights 503
503 in attention heads of lower layers. This research direction fo- 504
504 cuses on localizing specific attention heads and fine-tuning 505
505 them to unlearn memorized knowledge, showing promise for 506
506 enhancing privacy protection in LLMs. However, signifi- 507
507 cant challenges remain. Hu *et al.* [2024] demonstrate that 508
508 lifelong knowledge editing often fails due to superposition, 509
509 potentially affecting unrelated knowledge. Similarly, Gu *et al.* 510
510 [2024] identify potential damage to models’ general capa- 511
511 bilities from editing interventions.

512 The representation space offers another avenue for model 513
513 editing, as facts are encoded within these representations. 514
514 While most current research focuses on modifying represen- 515
515 tations during inference time, the impact of permanent mod- 516
516 ifications remains largely unexplored. A breakthrough study 517
517 proposes a more precise method for editing model represen- 518
518 tations to alter output distributions [Hernandez *et al.*, 2024]. 519
519 Rather than merely adding steering vectors during inference, 520
520 this approach directly modifies entity embeddings to trigger 521
521 targeted outputs, resulting in shifted entity positions in the 522
522 embedding space that causally influence model generations.

523 6.2 Improving Model Efficiency

524 In contrast to deciphering models’ inner workings, one 525
525 study examines the differences between pre-training and fine- 526
526 tuning phases using MI tools. It reveals that fine-tuning pre- 527
527 serves all capabilities learned during pre-training. The trans-

528 formations between pre-training and fine-tuning arise from
529 “wrappers” in MLPs learned on top of models. These wrap-
530 pers can be eliminated by pruning a few neurons or retraining
531 on an unrelated downstream task [Jain *et al.*, 2024]. This
532 discovery raises potential safety concerns regarding current
533 alignment approaches.

534 Unlike pruning neurons, Representation Engineering (RE)
535 directly manipulates representations without requiring opti-
536 mizations or additional labeled data. RE has also proven
537 effective in model pruning. Several studies show that fine-
538 tuning models with representation engineering can achieve
539 comparable or even better performance than state-of-the-
540 art fine-tuning techniques. For instance, Wu *et al.* [2024a]
541 employ forward passes from two topics and derive their
542 difference vectors, which are then used during inference
543 without additional fine-tuning. In contrast to conven-
544 tional parameter-efficient fine-tuning (PEFT), representation-
545 editing-based fine-tuning focuses on learning an additional
546 group of trainable parameters to modify representations di-
547 rectly rather than altering models’ parameters. The number
548 of trainable parameters has been reduced by a factor of 32
549 compared to LoRA [Wu *et al.*, 2024a]. Another approach em-
550 ploys distributed alignment search to find a set of linear sub-
551 spaces implementing interventions. This method outperforms
552 most PEFT models across various tasks [Wu *et al.*, 2024b].

553 6.3 Model Alignment to Human Values

554 From a mechanistic perspective, practical applications of-
555 ten evaluate model alignments using various tools. Inspired
556 by induction heads, Yang *et al.* [2023] measure bias scores
557 of attention heads in pre-trained LLMs by comparing atten-
558 tion scores between biased and regular heads, effectively re-
559 ducing gender bias through masking identified biased heads.
560 Another study localizes attention heads responsible for de-
561 ception using linear probing and activation patching. The
562 researchers use intentionally designed prompts to instruct
563 LLMs to be dishonest. Linear probes are trained to classify
564 true/false activations of heads, and selected activations related
565 to dishonest behaviors are then patched with those of honest
566 behaviors to observe output changes. This method success-
567 fully locates multiple attention heads across five layers.

568 Recently, RE has emerged as a promising approach for de-
569 tecting biases within representation space. A notable study
570 suggests that MLPs operate on token representations to alter
571 the output vocabulary distribution [Geva *et al.*, 2022]. Af-
572 ter reverse engineering MLPs, researchers concluded that the
573 output from each feed-forward layer can be viewed as sub-
574 updates to output vocabulary distributions, effectively pro-
575 moting certain high-level concepts. This insight has been
576 successfully applied to mitigate toxicity levels in LLMs. An-
577 other line of research identifies multiple representation vec-
578 tors within MLPs that encourage undesired model behaviors.
579 These vectors are decomposed using singular value decom-
580 position, enabling researchers to identify specific dimensions
581 that contribute to toxicity.

582 7 Future Research Directions

583 In this section, we highlight several research directions that
584 warrant further investigation by the research community.

7.1 Standardizing Evaluation and Benchmarks

586 A critical challenge in understanding LLMs is the lack of ro-
587 bust empirical validation for proposed theories like induction
588 heads and circuit analysis, compounded by the absence of
589 standardized evaluation metrics. The field needs agreed-upon
590 criteria for mechanistic explanations and comprehensive in-
591 terpretability benchmarks [Saphra and Wiegrefe, 2024] to
592 enable systematic comparison of different approaches. How-
593 ever, creating these benchmarks faces unique challenges due
594 to model complexity, absence of ground truth, and diverse
595 architectures. Establishing a rigorous validation framework
596 will require coordinated effort from the research community.

7.2 Scaling Interpretability Algorithms

598 As LLMs grow in size and complexity, scaling interpretabil-
599 ity techniques becomes critical. We identify two key chal-
600 lenges: First, we need novel methods to analyze the intricate
601 knowledge structures within LLMs, as current approaches re-
602 veal only a fraction of encoded knowledge. Second, exist-
603 ing explainability techniques, including MI tools, are primar-
604 ily effective on simplified models, creating an urgent need
605 to scale these algorithms to billion-parameter LLMs. This
606 challenge involves both computational efficiency and adapt-
607 ing conceptual frameworks to handle emergent behaviors in
608 large-scale models.

7.3 Source of Reasoning Ability

609 LLMs have demonstrated remarkable reasoning abilities that
610 mirror human cognition, performing complex tasks from
611 multi-step problem-solving to creative thinking. However, we
612 still lack understanding of how these abilities emerge from
613 architectural components and training dynamics. Key research
614 directions to uncover these mechanisms include: 1) examin-
615 ing how LLM components contribute to reasoning processes,
616 2) analyzing how basic building blocks combine to produce
617 complex reasoning, 3) investigating reasoning development
618 during training, and 4) comparing LLM and human reason-
619 ing patterns to identify similarities and differences.

8 Conclusions

621 In this survey paper, we present a comprehensive overview of
622 recent advances in uncovering the inner workings of LLMs
623 through explainability techniques. Our review highlights sig-
624 nificant findings about knowledge encoding within LLM ar-
625 chitectures, including polysemantic neurons, functional cir-
626 cuits, and attention heads’ role in in-context learning. Prob-
627 ing techniques and representation engineering have revealed
628 LLMs’ capacity to learn complex world models and en-
629 code factual knowledge, while studies of training dynamics
630 have illuminated how generalization abilities emerge. These
631 insights have enabled practical applications in model edit-
632 ing, efficient fine-tuning, and bias mitigation. Despite this
633 progress, challenges persist due to LLMs’ complexity and
634 scale, necessitating future work on standardized evaluation
635 metrics, scaling interpretability techniques, and understand-
636 ing high-level reasoning abilities.

References

- [Azaria and Mitchell, 2023] Amos Azaria and Tom Mitchell. The internal state of an llm knows when it’s lying. In *Findings of EMNLP*, pages 967–976, 2023.
- [Belinkov, 2022] Yonatan Belinkov. Probing classifiers: Promises, shortcomings, and advances. *Computational Linguistics*, 48(1):207–219, 2022.
- [Brown *et al.*, 2020] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, et al. Language models are few-shot learners. *NeurIPS*, 2020.
- [Chan *et al.*, 2022] Lawrence Chan, Adrià Garriga-Alonso, Nicholas Goldwosky-Dill, Ryan Greenblatt, et al. Causal scrubbing, a method for rigorously testing interpretability hypotheses. *AI Alignment Forum*, 2022.
- [Chughtai *et al.*, 2024] Bilal Chughtai, Alan Cooney, and Neel Nanda. Summing up the facts: Additive mechanisms behind factual recall in llms. *arXiv preprint arXiv:2402.07321*, 2024.
- [Cunningham *et al.*, 2024] Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoencoders find highly interpretable features in language models. In *ICLR*, 2024.
- [Davies *et al.*, 2022] Xander Davies, Lauro Langosco, and David Krueger. Unifying grokking and double descent. In *NeurIPS Workshop ML Safety*, 2022.
- [Doshi *et al.*, 2024] Darshil Doshi, Aritra Das, Tianyu He, and Andrey Gromov. To grok or not to grok: Disentangling generalization and memorization on corrupted algorithmic datasets. In *ICLR*, 2024.
- [Elhage *et al.*, 2021] Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, et al. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021.
- [Elhage *et al.*, 2022] Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, et al. Toy models of superposition, 2022.
- [Fan *et al.*, 2024] Simin Fan, Razvan Pascanu, and Martin Jaggi. Deep grokking: Would deep neural networks generalize better? *arXiv preprint arXiv:2405.19454*, 2024.
- [Ferrando *et al.*, 2024] 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.
- [Ferrando *et al.*, 2025] Javier Ferrando, Oscar Obeso, Senthoran Rajamanoharan, and Neel Nanda. Do i know this entity? knowledge awareness and hallucinations in language models. In *ICLR*, 2025.
- [Geva *et al.*, 2022] Mor Geva, Avi Caciularu, Kevin Ro Wang, and Yoav Goldberg. Transformer feed-forward layers build predictions by promoting concepts in the vocabulary space. In *EMNLP*, 2022.
- [Geva *et al.*, 2023] Mor Geva, Jasmijn Bastings, Katja Filipova, and Amir Globerson. Dissecting recall of factual associations in auto-regressive language models. In *EMNLP*, 2023.
- [Gould *et al.*, 2024] Rhys Gould, Euan Ong, George Ogden, and Arthur Conmy. Successor heads: Recurring, interpretable attention heads in the wild. In *ICLR*, 2024.
- [Gu *et al.*, 2024] Jia-Chen Gu, Hao-Xiang Xu, Jun-Yu Ma, Pan Lu, Zhen-Hua Ling, Kai-Wei Chang, et al. Model editing harms general abilities of large language models: Regularization to the rescue. In *EMNLP*, 2024.
- [Gurnee and Tegmark, 2024] Wes Gurnee and Max Tegmark. Language models represent space and time. In *ICLR*, 2024.
- [Gurnee *et al.*, 2023] Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, et al. Finding neurons in a haystack: Case studies with sparse probing. *arXiv preprint arXiv:2305.01610*, 2023.
- [Hernandez *et al.*, 2024] Evan Hernandez, Belinda Z Li, and Jacob Andreas. Inspecting and editing knowledge representations in language models. In *COLM*, 2024.
- [Hu *et al.*, 2024] Chenhui Hu, Pengfei Cao, Yubo Chen, Kang Liu, and Jun Zhao. Knowledge in superposition: Unveiling the failures of lifelong knowledge editing for large language models. *arXiv preprint arXiv:2408.07413*, 2024.
- [Hänni *et al.*, 2024] Kaarel Hänni, Jake Mendel, Dmitry Vaintrub, and Lawrence Chan. Mathematical models of computation in superposition. In *ICML Workshop MI*, 2024.
- [Jain *et al.*, 2024] Samyak Jain, Robert Kirk, Ekdeep Singh Lubana, Robert P Dick, Hidenori Tanaka, Edward Grefenstette, et al. Mechanistically analyzing the effects of fine-tuning on procedurally defined tasks. In *ICLR*, 2024.
- [Jermyn *et al.*, 2022] Adam S Jermyn, Nicholas Schiefer, and Evan Hubinger. Engineering monosemanticity in toy models. *arXiv preprint arXiv:2211.09169*, 2022.
- [Jin *et al.*, 2025] Mingyu Jin, Qinkai Yu, Jingyuan Huang, Qingcheng Zeng, Zhenting Wang, Wenyue Hua, et al. Exploring concept depth: How large language models acquire knowledge at different layers? In *COLING*, 2025.
- [Ju *et al.*, 2024] Tianjie Ju, Weiwei Sun, Wei Du, Xinwei Yuan, Zhaochun Ren, and Gongshen Liu. How large language models encode context knowledge? a layer-wise probing study. In *COLING*, 2024.
- [Kumar *et al.*, 2024] Tanishq Kumar, Blake Bordelon, Samuel J Gershman, and Cengiz Pehlevan. Grokking as the transition from lazy to rich training dynamics. In *ICLR*, 2024.
- [Lecomte *et al.*, 2024] Victor Lecomte, Kushal Thaman, Rylan Schaeffer, Naomi Bashkansky, Trevor Chow, and Sanmi Koyejo. What causes polysemanticity? an alternative origin story of mixed selectivity from incidental causes. In *ICLR Workshop RA*, 2024.
- [Lee *et al.*, 2024] Andrew Lee, Xiaoyan Bai, Itamar Pres, Martin Wattenberg, Jonathan Kummerfeld, and Rada Mihalcea. A mechanistic understanding of alignment algorithms: A case study on dpo and toxicity. In *ICML*, 2024.

- [Li *et al.*, 2023] Kenneth Li, Aspen K Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Emergent world representations: Exploring a sequence model trained on a synthetic task. In *ICLR*, 2023.
- [Li *et al.*, 2024] Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-time intervention: Eliciting truthful answers from a language model. *NeurIPS*, 36, 2024.
- [Liu *et al.*, 2023] Ziming Liu, Eric J Michaud, and Max Tegmark. Omnigrok: Grokking beyond algorithmic data. In *ICLR*, 2023.
- [Marks and Tegmark, 2024] Samuel Marks and Max Tegmark. The geometry of truth: Emergent linear structure in large language model representations of true/false datasets. In *COLM*, 2024.
- [Meng *et al.*, 2022] Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in gpt. *NeurIPS*, 35:17359–17372, 2022.
- [Menta *et al.*, 2025] Tarun Ram Menta, Susmit Agrawal, and Chirag Agarwal. Analyzing memorization in large language models through the lens of model attribution. *arXiv preprint arXiv:2501.05078*, 2025.
- [Nakkiran *et al.*, 2020] Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, et al. Deep double descent: Where bigger models and more data hurt. In *ICLR*, 2020.
- [Nanda *et al.*, 2023a] Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, et al. Progress measures for grokking via mechanistic interpretability. In *ICLR*, 2023.
- [Nanda *et al.*, 2023b] Neel Nanda, Andrew Lee, and Martin Wattenberg. Emergent linear representations in world models of self-supervised sequence models. In *ACL Workshop BlackboxNLP*, pages 16–30, 2023.
- [Olah *et al.*, 2020] Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter. Zoom in: An introduction to circuits. *Distill*, 2020.
- [Olsson *et al.*, 2022] Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, et al. In-context learning and induction heads. *Transformer Circuits Thread*, 2022.
- [Power *et al.*, 2022] 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.
- [Saphra and Wiegreffe, 2024] Naomi Saphra and Sarah Wiegreffe. Mechanistic? *arXiv preprint arXiv:2410.09087*, 2024.
- [Scherlis *et al.*, 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.
- [Stoehr *et al.*, 2024] Niklas Stoehr, Mitchell Gordon, Chiyuan Zhang, and Owen Lewis. Localizing paragraph memorization in language models. *arXiv preprint arXiv:2403.19851*, 2024.
- [Templeton *et al.*, 2024] Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen, Adam Pearce, Craig Citro, et al. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. *Transformer Circuits Thread*, 2024.
- [Thilak *et al.*, 2022] Vimal Thilak, Etai Littwin, Shuangfei Zhai, Omid Saremi, Roni Paiss, and Joshua Susskind. The slingshot mechanism: An empirical study of adaptive optimizers and the grokking phenomenon. In *NeurIPS Workshop HITY*, 2022.
- [Todd *et al.*, 2024] Eric Todd, Millicent L Li, Arnab Sen Sharma, Aaron Mueller, Byron C Wallace, and David Bau. Function vectors in large language models. In *ICLR*, 2024.
- [Vaswani *et al.*, 2017] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *NeurIPS*, 30, 2017.
- [Wang *et al.*, 2024] Boshi Wang, Xiang Yue, Yu Su, and Huan Sun. Grokking of implicit reasoning in transformers: A mechanistic journey to the edge of generalization. In *NeurIPS*, 2024.
- [Wang *et al.*, 2025] Xinyi Wang, Antonis Antoniadis, Yanai Elazar, Alfonso Amayuelas, Alon Albalak, et al. Generalization v.s. memorization: Tracing language models’ capabilities back to pretraining data. In *ICLR*, 2025.
- [Wu *et al.*, 2024a] Muling Wu, Wenhao Liu, Xiaohua Wang, Tianlong Li, et al. Advancing parameter efficiency in finetuning via representation editing. In *ACL*, 2024.
- [Wu *et al.*, 2024b] Zhengxuan Wu, Aryaman Arora, Zheng Wang, Atticus Geiger, Dan Jurafsky, Christopher D Manning, and Christopher Potts. Reft: Representation finetuning for language models. In *NeurIPS*, 2024.
- [Xie *et al.*, 2024] Chulin Xie, Yangsibo Huang, Chiyuan Zhang, Da Yu, et al. Large language interpolators can learn logical reasoning: A study on knights and knaves puzzles. In *NeurIPS Workshop MATH-AI*, 2024.
- [Yang *et al.*, 2023] Yi Yang, Hanyu Duan, Ahmed Abbasi, John P Lalor, and Kar Yan Tam. Bias a-head? analyzing bias in transformer-based language model attention heads. *arXiv preprint arXiv:2311.10395*, 2023.
- [Zhao *et al.*, 2024] Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu, et al. Explainability for large language models: A survey. *ACM TIST*, 15(2):1–38, 2024.
- [Zhu *et al.*, 2024] Xuekai Zhu, Yao Fu, Bowen Zhou, and Zhouhan Lin. Critical data size of language models from a grokking perspective. *arXiv preprint arXiv:2401.10463*, 2024.
- [Zou *et al.*, 2023] Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, et al. Representation engineering: A top-down approach to ai transparency. *arXiv preprint arXiv:2310.01405*, 2023.