

# Addressing the Reasoning Gap: Mechanistic Circuit-Based Knowledge Editing in Large Language Models

Anonymous ACL submission

## Abstract

Deploying Large Language Models (LLMs) in real-world dynamic environments raises the challenge of updating their pre-trained knowledge. While existing knowledge editing methods can reliably patch isolated facts, they frequently suffer from a "Reasoning Gap", where the model recalls the edited fact but fails to utilize it in multi-step reasoning chains. To bridge this gap, we introduce MCircKE (Mechanistic Circuit-based Knowledge Editting), a novel framework that enables a precise "map-and-adapt" editing procedure. MCircKE first identifies the causal circuits responsible for a specific reasoning task, capturing both the storage of the fact and the routing of its logical consequences. It then surgically update parameters exclusively within this mapped circuit. Extensive experiments on the MQuAKE-3K benchmark demonstrate the effectiveness of the proposed method for multi-hop reasoning in knowledge editing.

## 1 Introduction

The deployment of large language models (LLMs) in dynamic, real-world environments is fundamentally constrained by the static nature of their pre-trained knowledge. Knowledge editing techniques (Wang et al., 2024c; Zhang et al., 2024) seek to address this limitation by modifying specific facts in a trained model without full retraining. While recent methods can reliably update isolated single-hop facts, they often exhibit a critical failure mode termed the Reasoning Gap (Yao et al., 2025; Zhang et al.). In this phenomenon, a model may correctly recall an edited fact in isolation (e.g., "Who is the PM?") but fail to propagate this updated information through multi-step reasoning (e.g., "Which political party is the PM the leader of?"). This discrepancy suggests that, although the storage of the fact has been successfully patched, the reasoning pathways required to utilize that fact remain misaligned.

Traditional knowledge editing methods (Meng et al., 2022; Wang et al., 2024a; Fang et al., 2025) typically treat knowledge as isolated atomic units stored within specific multi-layer perceptron (MLP) layers. However, recent advances in mechanistic interpretability suggest that knowledge retrieval and utilization are instead governed by distributed reasoning circuits – sparse computational subgraphs composed of attention heads and MLPs that route information across layers. From this perspective, the Reasoning Gap is not merely a retrieval failure but a structural misalignment; the model’s logical "wires" are still connected to obsolete reasoning paths. Recent work has begun to address this challenge. For instance, CaKE (Yao et al., 2025) posits that the reasoning gap exists because standard edits do not compel the model to practice using the new knowledge. It mitigates this by generating circuit-aware synthetic multi-hop data derived from the edited fact, and fine-tuning the model on this data. By training on this curated data, CaKE aims to implicitly "activate" the dynamic reasoning circuits needed for generalization. While effective, this strategy relies on the indirect pressure of data augmentation to realign the circuits, and thus still treats the internal mechanisms as a black box that is expected to self-correct given appropriate inputs.

To tackle the aforementioned issue, we propose MCircKE, a novel framework that addresses the limitations of existing knowledge editing methods and fundamentally diverges from the data-stimulation strategy of CaKE. Rather than relying on pre-constructed data to implicitly activate reasoning paths, our approach aims to explicitly and precisely identify the dynamic reasoning circuits responsible for the target knowledge. Specifically, we first identify a rigorous, high-fidelity map of the edges and nodes that causally contribute to the model’s reasoning process. By integrating gradients along the activation path, we can pinpoint the exact "wires" in the neural network that

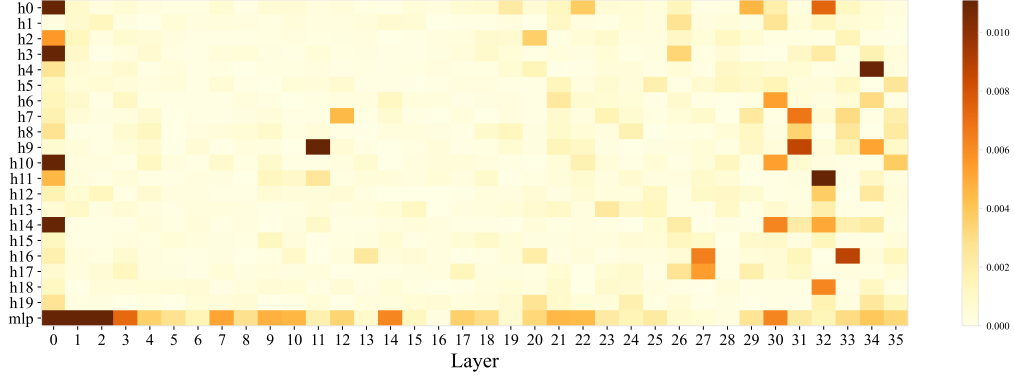


Figure 1: Heatmap of Attribution Scores across 100 instances in GPT-2 Large.

carry the knowledge to be edited. Once this reasoning circuits are identified, we then surgically adapt the model along the mapped paths. This transforms knowledge editing from a heuristic “locate-and-replace” operation into a structural “map-and-adapt” procedure. By explicitly targeting components that have been verified as part of the reasoning chains, our method ensures that the updated knowledge is not only stored but also mechanistically integrated into the model’s logic. In summary, our main contributions are:

- **Mechanistic Insight:** We conduct a mechanistic analysis of multi-hop factual reasoning in LLMs, providing coarse-to-fine-grained insights into why existing knowledge editing methods fail at multi-hop factual recall and how this failure can be addressed.
- **MCircKE Framework:** We propose the first mechanistic editing framework that explicitly maps dynamic reasoning circuits and performs surgical, path-constrained updates.
- **Empirical Validation:** Extensive experiments demonstrate that MCircKE substantially improves multi-hop reasoning performance for knowledge editing in LLMs.

## 2 Preliminaries

### 2.1 LLM as a Computational Graph

We conceptualize the LLM as a Directed Acyclic Graph (DAG), denoted as  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ . In this representation, the nodes  $\mathcal{V}$  are the specific computational sub-modules of the network, and the edges  $\mathcal{E}$  represent the flow of information between them via the residual stream.

**Nodes ( $\mathcal{V}$ )** To enable fine-grained editing, we decompose the standard Transformer blocks into the

atomic linear projections of the models:

$$\mathcal{V} = \{W_Q^{(l,h)}, W_K^{(l,h)}, W_V^{(l,h)}, W_O^{(l,h)}, W_{MLP}^{(l)} \mid l \in [1, L], h \in [1, H]\} \quad (1)$$

where  $L$  is the number of layers and  $H$  is the number of heads per layer.  $W_Q, W_K, W_V$  represent the query, key, and value projections for the  $h$ -th head in layer  $l$ ,  $W_O$  is the output projection, and  $W_{MLP}$  represents the two-layer feed-forward network.

**Edges ( $\mathcal{E}$ )** An edge  $e_{u \rightarrow v}$  exists if the output of module  $u$  contributes to the input of module  $v$  via the residual stream.

### 2.2 Edge Attribution Patching with Integrated Gradients (EAP-IG)

To identify the specific subgraph responsible for a model’s behavior, we require a method to attribute the model’s output to specific internal edges. Standard methods like Activation Patching (Meng et al., 2022) are faithful but computationally expensive, while vanilla gradient attribution suffers from saturation. To balance fidelity and efficiency, we employ Edge Attribution Patching with Integrated Gradients (EAP-IG) (Hanna et al.).

Specifically, vanilla EAP (Nanda, 2023) approximates the importance of an edge  $e$  by computing the product of the edge’s activation  $x_e$  and the gradient of the loss with respect to that activation  $\frac{\partial \mathcal{L}}{\partial x_e}$ :

$$\phi_{EAP} = x_e \cdot \frac{\partial \mathcal{L}}{\partial x_e} \quad (2)$$

While requiring only one backward pass, this method fails in deep non-linear networks due to the saturation effect. When a neuron is fully activated, the local gradient is near zero, even if that neuron is critically important for the output. This often yields “broken” circuits that miss key intermediate nodes.

To overcome saturation, EAP-IG accumulates gradients along a linear path between a corrupted baseline input ( $I_{corrupt}$ ) and the clean target input ( $I_{clean}$ ). Let  $x_e$  denote the activation along edge  $e$ . A path is defined as  $x_e(\gamma) = x_e^{corrupt} + \gamma(x_e^{clean} - x_e^{corrupt})$  for  $\gamma \in [0, 1]$ . The attribution score  $\phi(e)$  for an edge  $e$  is calculated as:

$$\phi(e) = (x_e^{clean} - x_e^{corrupt}) \times \int_{\gamma=0}^1 \frac{\partial \mathcal{L}(x(\gamma))}{\partial x_e(\gamma)} d\gamma \quad (3)$$

In practice, this integral is approximated using a Riemann sum with  $m$  steps:

$$\phi(e) \approx (x_e^{clean} - x_e^{corrupt}) \times \frac{1}{m} \sum_{k=1}^m \frac{\partial \mathcal{L}(x(\frac{k}{m}))}{\partial x_e} \quad (4)$$

This formalism provides a robust, continuous estimate of causal importance, enabling us to map the full reasoning circuitry even through saturated MLP layers.

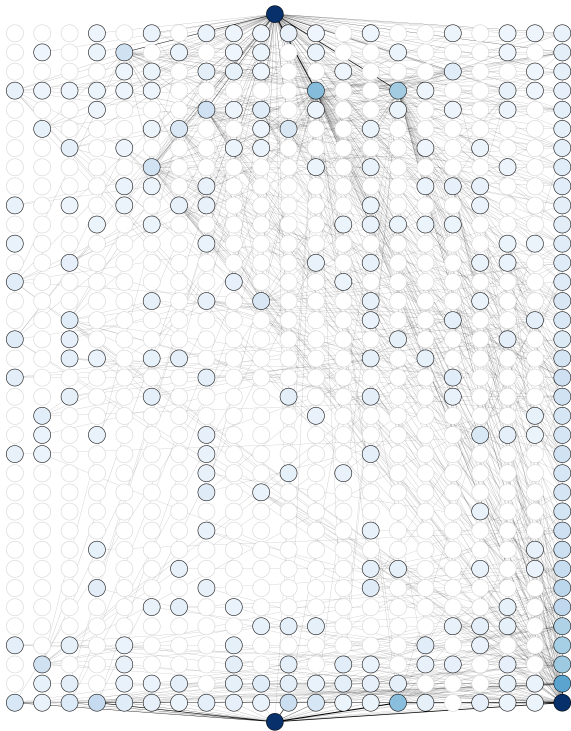


Figure 2: Visualization of the reasoning circuits for a specific multi-hop instance in GPT-2 Large.

### 3 Mechanistic Analysis

We start by conducting a mechanistic diagnostic study to uncover the causal circuitry supporting multi-hop reasoning and to explain, at a structural level, why conventional knowledge editing approaches fail to generalize to multi-hop edits.

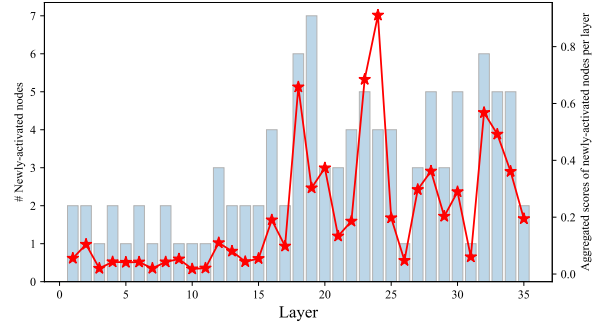


Figure 3: Distribution of newly-activated nodes in the multi-hop circuits compared to the single-hop circuits.

#### 3.1 Inter-Layer Distribution of Causal Importance

We first aggregate the EAP-IG attribution scores for all model components (20 attention heads per layer) in GPT-2 Large across 100 multi-hop knowledge instances, as shown in Figure 1. We can observe several critical phases from the resulting heatmap.

**Early MLP Layers.** We observe dense, high-magnitude activations in the MLP layers of the early network (layer 0-10). This aligns with the hypothesis established in ROME (Meng et al., 2022) and MEMIT (Meng et al., 2023): these layers are responsible for the initial retrieval of the atomic fact. Standard methods usually successfully target these layers;

**Intermediate Processing.** Crucially, there is a distinct cluster of activity in MLPs around layers 17-22. This could possibly suggest a transitional phase where the retrieved entity embedding is refined or transformed to be compatible with downstream relation extraction;

**Late Layers.** Late layers (layer 29-35) show strong activations of both attention heads and MLPs. We posit that these attention heads could act as *routing heads*, moving information to the final prediction position to support multi-hop reasoning. Meanwhile, the high activity in the MLPs of layers 29-35 indicates that "routing" alone is insufficient. The moved information must also be processed and decoded into the correct output distribution.

#### 3.2 Intra-Layer Sparsity and Functional Specialization

The above layer-wise view provides insights into high-level trends. We further conduct a granular analysis of specific reasoning instance investigating intra-layer divergence and case-dependent activation topologies difference between single-hop and

multi-hop knowledge instance.

**Visual Analysis of Sparse Circuitry** We first visualize the extracted circuits for a specific multi-hop instance (Crysis\_2→Crytek→English). As shown in Figure 2, we can see that only a small fraction of heads (blue nodes) are active, forming a distinct pathway that differs from the aggregate model behavior. A closer examination of the attention mechanism reveals another critical insight: Even within highly active reasoning layers, importance is not uniformly distributed; it is concentrated in specific, specialized heads. For instance, in layer 32, Head 11 acts as a critical hub, while its neighbors (h10 and h12) are dormant. This intra-layer divergence indicates that operating at the layer granularity (e.g., updating all of Layer 32) are inherently too coarse and may inadvertently alter ‘bystander’ heads that serve unrelated functions.

**Differential Activation Dynamics** To mechanistically understand the underlying difference between single-hop and multi-hop reasoning, we further compare the circuits activated by the following reasoning chains: Multi-hop (Crysis\_2→Crytek→English); Single-hop (Crytek→English). We define newly activated nodes as components that are active in the multi-hop circuit graph but inactive in the single-hop circuit graph. Figure 3 visualizes the distribution of these nodes and their corresponding attribution scores. The quantitative results suggest a topological shift. Specifically, we observe a dramatic increase in node activation in the mid-to-late layers, with aggregate causal importance (the red-star curve) peaking at Layers 18 and 24. Crucially, these nodes remain dormant during single-hop retrieval. Standard editing methods, which compute updates from single-hop traces, therefore miss this regime entirely. They patch the destination (Crytek → English) but leave the bridge (Crysis\_2 → Crytek) disconnected. A secondary cluster of activity appears in Layers 32-34, indicating that the final answer extraction also relies on instance-specific routing heads distinct from those used in simple retrieval.

**Summary.** The above analysis provides the mechanistic justification for circuit-based editing. A successful edit requires not only updating the relevant stored information, but also repairing the routing edges that govern its use, ensuring the new fact is not just stored, but mechanistically integrated into the logic of the model.

## 4 Mechanistic Circuit-based Knowledge Editing for LLMs

Building on the insights from our mechanistic analysis, we introduce MCircKE (Mechanistic Circuit-based Knowledge Eding), a framework designed to bridge the reasoning gap by editing the specific causal circuits responsible for knowledge utilization. Unlike heuristic methods that target fixed layers or data-driven methods that rely solely on implicit stimulation, our approach explicitly maps the causal reasoning pathway for a specific fact and surgically adapts the model along that pathway.

### 4.1 Circuit Discovery

Our first objective is to identify the reasoning circuits  $\mathcal{C} \subset \mathcal{G}$  casually responsible for the target reasoning task.

**Clean and Corrupted Input Pairs** *Clean prompts construction:* We first combine the provided single-hop statements from the original knowledge and chain them into a unified multi-hop prompt. *Corrupted prompts construction:* We then utilize ChatGPT to generate a corresponding corrupted prompt that maintains semantic manifold consistency to the clean prompt. This ensures the prompt structure and reasoning complexity remain identical, but the subject entity is altered to an alternative. An example is shown in Figure 4.

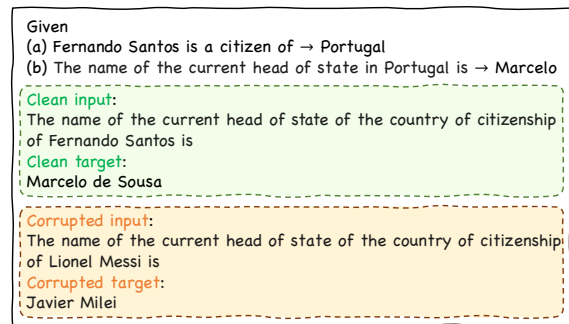


Figure 4: Example clean and corrupted inputs.

**Attribution Integration** We compute the EAP-IG score  $\phi(e)$  for all edges in the model by integrating the gradients along the linear interpolation path between the corrupted and clean embeddings, using the batched Riemann approximation (Equation 4) with  $m = 5$ . We use the magnitude  $|\phi(e)|$  because we are interested in causal relevance, regardless of whether the edge positively or negatively impacts the output.

Model	Method	Overall M-Acc.	#hops			S-Acc.
			2-hop	3-hop	4-hop	
GPT-2 Large	LoRA	24.0858	39.3886	13.4390	16.8724	48.947
	ROME	14.9939	17.2546	13.1162	14.4033	52.956
	MEMIT	7.7804	10.9053	5.7805	6.0357	35.009
	WISE	8.1138	11.6990	5.9272	5.9442	21.421
	AlphaEdit	8.5028	13.2275	5.6925	5.5327	44.093
	CaKE	34.9561	37.0370	29.3721	40.4207	<u>71.666</u>
	<b>MCircKE</b>	<b>50.4168</b>	<b>67.9306</b>	<u>38.6737</u>	<u>41.4723</u>	64.915
	<b>MCircKE★</b>	<u>47.2824</u>	<u>54.1446</u>	<b>39.1432</b>	<b>49.2913</b>	<b>73.213</b>
GPT-2 XL	LoRA	25.6085	41.9753	15.0528	16.5981	52.933
	ROME	12.5042	17.7543	9.8298	8.5048	61.511
	MEMIT	8.4806	12.5220	5.7218	6.4929	46.759
	WISE	8.3917	11.7578	5.7512	7.2702	22.876
	AlphaEdit	9.4587	14.6972	6.0446	6.6301	48.749
	CaKE	40.5246	46.6196	33.3040	42.2954	<u>76.659</u>
	<b>MCircKE</b>	<u>49.0140</u>	<b>64.7860</b>	<u>39.0020</u>	<u>45.4960</u>	73.894
	<b>MCircKE★</b>	<b>49.3536</b>	<u>58.1423</u>	<b>39.7887</b>	<b>53.2468</b>	<b>77.446</b>

Table 1: Performance comparison under different hop settings on GPT-2 Large and GPT-2 XL.

**Circuit Pruning** The raw attribution map assigns a non-zero score to almost every edge. To obtain the sparse reasoning circuits, we employ top- $K$  pruning. Specifically, we rank all edges based on the magnitude of their attribution scores  $|\phi(e)|$ , and retain only the top  $K$  edges with the highest attribution scores, effectively filtering out noise and irrelevant pathways.

## 4.2 Circuit-Guided Low-Rank Adaptation

Once the circuits  $\mathcal{C}$  are identified, we apply Low-Rank Adaptation exclusively to the modules within the circuits. For every weight matrix  $W \in \mathcal{C}$  identified in the reasoning circuit, we freeze  $W$  and introduce low-rank matrices  $A$  and  $B$ :

$$W_{new} = W + \Delta W = W + \frac{\alpha}{r} BA \quad (5)$$

where  $B \in \mathbb{R}^{d \times r}$ ,  $A \in \mathbb{R}^{r \times k}$ . Here  $\alpha$  is a scaling factor and  $r$  is the rank, with  $r \ll d$ . Critically, for any module  $W' \notin \mathcal{C}$  (parameters outside the reasoning circuits), we enforce  $\Delta W' = 0$ . This constraint ensures that the edit is mechanistically localized and physically rewires the specific information flow responsible for the logic.

## 5 Experiments

In this section, we aim to answer the following research questions:

- **RQ1:** Can the proposed method effectively bridge the reasoning gap in multi-hop factual recall compared to baselines?
- **RQ2:** Does the proposed method ensure the preservation of unrelated knowledge?
- **RQ3:** Is precise topological identification of reasoning circuits necessary, and do the identified circuits demonstrate validity relative to stochastic and heuristic baselines?

### 5.1 Evaluation Setup

**Dataset.** We conduct the evaluation on the challenging MQuAKE-3K benchmark (Multi-hop Question Answering for Knowledge Editing) (Zhong et al., 2023), a widely used dataset for evaluating multi-hop factual recall with 3,000 counterfactual editing instances.

**Baselines and Models.** We compare our method against several representative knowledge editing baselines, including: **LoRA** (Hu et al., 2022), which directly fine-tunes the full model; **ROME** (Meng et al., 2022), a classic locate-then-edit approach; **MEMIT** (Meng et al., 2023), an extension of ROME that performs edits across a range of early layers; **WISE** (Wang et al., 2024a), which augments the model with a side memory and edits later layers; **AlphaEdit** (Fang et al., 2025), which

updates parameters via projection into a null space; **CaKE** (Yao et al., 2025), which constructs additional multi-hop edit prompts to improve multi-hop factual recall. We conduct experiments on three GPT-2 variants (Radford et al., 2019): GPT-2 Large (774M), GPT-2 XL (1.5B), and GPT-J (6B).

**Metrics.** We assess the model’s ability to leverage edited knowledge along three key dimensions: *Multi-hop Accuracy*, measuring performance on questions requiring two-, three-, and four-hop reasoning over the edited fact; *Single-hop Accuracy*, capturing direct recall of the edited knowledge; and *Locality*, quantifying the extent to which an edit preserves the model’s behavior on queries unrelated to the targeted knowledge.

**Implementation.** We evaluate our proposed method in two configurations to isolate the impact of data versus structure: **MCircKE**, which edits the model using only the single-hop edit sample for the target fact provided in MQuAKE; **MCircKE<sup>★</sup>**, which additionally incorporates synthetic multi-hop training data from CaKE (Yao et al., 2025) for model editing. For the hyperparameters, we set  $K = 4,000$  for GPT-2 Large,  $K = 10,000$  for GPT-2 XL, and  $K = 5,000$  for GPT-J.  $\alpha$  is set to 8 and rank  $r$  is set to 32 in the circuit-guided low-rank adaptation. We run our experiments on 6 NVIDIA A6000 GPUs.

## 5.2 Experimental Results

**Compared with baselines, MCircKE consistently bridges the reasoning gap across all model scales.** Table 1 presents the multi-hop and single-hop accuracies for all methods on GPT-2 Large and GPT-2 XL. We also present the comparative performance on GPT-J in Table 2. Overall, MCircKE achieves the highest multi-hop accuracy among all methods, delivering a substantial improvement over the best previous approach. For instance, on GPT-2 Large, our method attains an overall multi-hop accuracy of 50.4%, outperforming CaKE (which achieves 34.9%) by over 15 percentage points. On GPT-2 XL, MCircKE outperforms the strongest baseline by approximately 9% in overall multi-hop performance. On GPT-J, it reaches 59.25% overall multi-hop accuracy, again surpassing the best baseline. Notably, the gains are consistent across different reasoning depths. For instance, on GPT-2 XL, it yields consistent gains on 2-hop questions (49.35% vs. 40.52% for CaKE), 3-hop (39.78% vs. 39%) and 4-hop queries (53.24% vs. 45.49%). These

Table 2: Performance comparison under different hop settings on GPT-J.

Method	M-hop (2-hop / 3-hop / 4-hop)	S-hop
LoRA	9.769 (12.61 / 6.72 / 10.11)	47.579
ROME	23.23 (32.89 / 17.55 / 17.06)	52.497
CaKE	56.318 (61.84 / 46.77 / 62.59)	<b>85.277</b>
<b>MCircKE</b>	<b>59.253 (64.61 / 49.41 / 66.26)</b>	<u>85.219</u>

results provide strong evidence in support of RQ1, demonstrating that our method effectively enhances multi-hop factual recall and enables the model to answer complex, chained queries grounded in the edited knowledge.

**The collapse of unstructured fine-tuning.** Standard LoRA struggles significantly compared to MCircKE. For instance, on GPT-2 XL, LoRA achieves only  $\sim 25\%$  multi-hop accuracy compared to MCircKE’s  $\sim 49\%$ . This gap suggests that, with limited topological guidance, global parameter updates must search an extremely large parameter space to recover the correct reasoning pathways, making it fail to converge to the desired circuit-level change. In contrast, MCircKE’s strong performance indicates that mechanistic circuit guidance is crucial for reliable knowledge editing with multi-hop generalization.

**The reasoning vs. recall trade-off.** Comparing the standard MCircKE with CaKE (trained on additional multi-hop data) reveals a critical trade-off. MCircKE consistently outperforms CaKE on multi-hop reasoning (e.g., 50.42% vs. 34.96% accuracy on GPT-2 Large), indicating that repairing circuit structure yields stronger generalization than implicit behavioral imitation. However, CaKE retains a slight edge in single-hop recall. This suggests that CaKE’s approach of fine-tuning on augmented data leads to stronger surface-level memorization of the facts, whereas MCircKE’s surgical approach prioritizes the logical consistency of the circuit, occasionally at the cost of absolute recall strength for the atomic fact itself.

**The failure of heuristic "Locate-and-Edit".** Traditional methods like ROME and MEMIT perform poorly on multi-hop reasoning (M-Acc  $< 15\%$  on GPT-2 Large and GPT-2 XL). While they achieve decent Single-hop Accuracy, they fail to integrate this knowledge into the reasoning chain. This confirms our mechanistic hypothesis derived in Section 3: fixing the "storage" without fixing

Table 3: Locality performance.

$ \Delta $ (%)	MMLU	CSQA
LoRA	0.34	1.24
ROME	0.06	0.41
AlphaEdit	0.21	0.25
CaKE	0.09	0.49
MCircKE	0.10	0.24

the "routing" leaves the model unable to leverage the new fact. ROME effectively creates "orphaned facts" that can be recalled but not reasoned with, and the resulting reasoning gap (delta between S-Acc and M-Acc) is stark. On GPT-2 XL, ROME exhibits a  $\sim 49\%$  drop from single-hop to multi-hop accuracy (61.5% S-Acc vs. 12.5% M-Acc). In contrast, MCircKE shrinks this gap to  $\sim 25\%$  while achieving higher S-Acc and M-Acc, demonstrating significantly better knowledge integration.

**Locality and stability analysis.** To assess locality, Table 3 presents the performance change on two unrelated benchmarks, MMLU (Hendrycks et al., 2021) and CommonsenseQA (Talmor et al., 2019), which evaluate the model’s general abilities after applying the knowledge editing methods. Overall, the minimal performance shifts indicate that MCircKE’s edits are highly localized, integrating new knowledge without disrupting the model’s existing unrelated knowledge.

### 5.3 Sensitivity Analysis

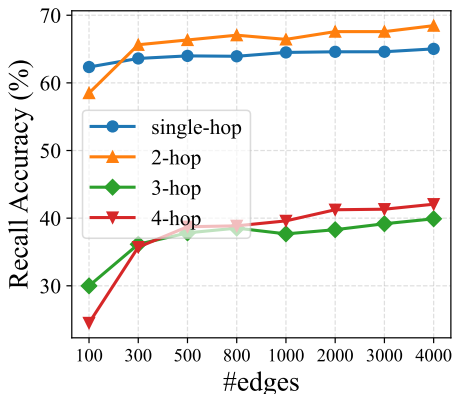


Figure 5: Effect of the number of edges retained in the discovered circuit graphs.

A core hyperparameter in MCircKE is the number of edges retained in the reasoning circuit. A very sparse circuit is desirable for locality (modifying fewer parameters), but an overly sparse circuit may

sever critical reasoning pathways. Figure 5 illustrates the impact of varying the edge count from 100 to 4000 on GPT-2 Large. We can observe that:

- At extreme sparsity levels (100-300 edges), single-hop accuracy remains relatively robust ( $\sim 62\%$ ), but multi-hop performance suffers catastrophic failure (4-hop accuracy drops to  $\sim 25\%$ ). This disparity suggests that simple factual recall relies on a highly localized set of "storage" edges (likely within specific MLP layers), whereas multi-hop reasoning requires a more extensive "routing" infrastructure.
- As the edge count increases, we observe that performance effectively plateaus beyond 2,000 edges. This saturation point suggests that the logic required to process a specific fact is not diffuse across the entire model but is concentrated in a fraction of the network.

### 5.4 Ablation Study: Circuit Validity

To further verify the effectiveness of editing the model along the identified circuits, we conducted an ablation study on 1,000 cases from the MQuAKE dataset using GPT-2 Large. We compared MCircKE against three alternative strategies: **Full Graph** - Updating all edges in the model (equivalent to standard full model fine-tuning). **Only MLPs** - Updating all MLP modules while freezing all other modules. **Random Edges** - Updating a random subset of edges equal in size to the MCircKE circuits. The results are shown in Table 4. We can observe that:

Method	Single-hop	Multi-hop
Full Graph	61.40	39.30
Only MLPs	58.85	26.00
Random Edges	26.35	15.13
<b>MCircKE</b>	<b>73.17</b>	<b>67.73</b>

Table 4: Ablation Results.

- MCircKE drastically outperforms Random Edges, confirming that the circuits utilized in the proposed method are mechanistically relevant and not merely a random sub-network.
- The Only MLPs baseline illustrates the precise failure mode of traditional editing. While it achieves decent Single-hop accuracy, its Multi-hop performance collapses to 26.00%. This empirically validates the hypothesis that

502 modifying storage alone is insufficient; with-  
503 out the concurrent update of routing circuits,  
504 the new knowledge cannot be propagated.

- 505 • MCircKE outperforms the Full Graph update.  
506 This suggests that restricting updates to the rel-  
507 evant circuits acts as a form of structural reg-  
508 ularization. By freezing irrelevant pathways,  
509 we prevent the noise of global updates from  
510 interfering with the delicate logic of the rea-  
511 soning chains, allowing the focus on rewiring  
512 the specific factual bridge.

## 513 6 Related Work

### 514 6.1 Knowledge Editing for LLMs

515 Knowledge Editing (Wang et al., 2024c; Zhang  
516 et al., 2024; Wang et al., 2024b) aims to update  
517 factual associations in Large Language Models  
518 (LLMs) without the computational cost of retrain-  
519 ing or the degradation associated with catastro-  
520 phic forgetting. Methodologies in this field have evolved  
521 from manipulating static parameter storage to align-  
522 ing dynamic pathways. Foundational approaches  
523 treat Transformer Feed-Forward Networks (FFNs)  
524 as key-value associative memories. For instance,  
525 ROME (Meng et al., 2022) employs causal trac-  
526 ing to localize factual storage to specific mid-layer  
527 MLPs and formulates a rank-one update to insert a  
528 new fact tuple. While ROME is effective for single  
529 edits, it struggles with scalability. MEMIT (Meng  
530 et al., 2023) addresses this by distributing updates  
531 across a range of critical layers. By solving a least-  
532 squares problem with residual spreading, MEMIT  
533 allows for the simultaneous injection of thousands  
534 of memories, significantly outperforming ROME  
535 in batch-editing scenarios. In sequential "life-  
536 long" editing scenarios, maintaining the stability  
537 of previously learned information is paramount.  
538 WISE (Wang et al., 2024a) critiques the destructive  
539 nature of continuous parameter overwriting and in-  
540 troduces a side memory architecture coupled with a  
541 routing mechanism. Conversely, AlphaEdit (Fang  
542 et al., 2025) enforces stability through geomet-  
543 ric constraints rather than architectural changes.  
544 It projects parameter perturbations onto the Null  
545 Space of the preserved knowledge matrix to en-  
546 sure that the model's performance on pre-existing  
547 data remains invariant. Another critical limitation  
548 of storage-focused editors is the "Reasoning Gap",  
549 where models recall a specific fact but fail to apply  
550 it in multi-step inference. IFMET (Zhang et al.)  
551 identifies a functional stratification: shallow layers

552 store explicit facts, while deep layers process im-  
553 plicit subjects necessary for multi-hop reasoning.  
554 It proposes a two-stage strategy that edits both shal-  
555 low and deep layers to ensure the updated knowl-  
556 edge propagates through the entire reasoning chain.  
557 On the other hand, CaKE (Yao et al., 2025) moves  
558 beyond the storage metaphor entirely. It posits  
559 that knowledge is utilized via dynamic reasoning  
560 circuits. CaKE generates complex, circuit-aware  
561 training data designed to activate these specific  
562 pathways, stimulating the model to functionally in-  
563 tegrate new information into its logical processing.

### 564 6.2 Mechanistic Interpretability

565 Mechanistic interpretability (Rai et al., 2024;  
566 Bereska and Gavves, 2024; Sharkey et al., 2025)  
567 aims to reverse-engineer model behaviors into  
568 human-understandable algorithms implemented  
569 by specific subgraphs or "circuits". Pioneering  
570 work manually identified circuits for tasks  
571 like Indirect Object Identification (Wang et al.)  
572 and arithmetic (Hanna et al., 2023). To scale  
573 this analysis, automated methods were developed:  
574 ACDC (Conmy et al., 2023) employs activation  
575 patching to prune irrelevant edges but is compu-  
576 tationally expensive. Edge Attribution Patching  
577 (EAP) (Nanda, 2023) approximates these effects  
578 via gradients for efficiency. Most recently, EAP  
579 with Integrated Gradients (EAP-IG) (Hanna et al.)  
580 was proposed to resolve the gradient saturation is-  
581 sue in EAP, offering a superior balance of speed and  
582 faithfulness. Our work leverages EAP-IG to iden-  
583 tify dynamic reasoning circuits for knowledge edit-  
584 ing, moving beyond the static localization assump-  
585 tions of prior editing methods like ROME (Zhang  
586 et al.) or MEMIT (Meng et al., 2023).

## 587 7 Conclusion

588 This paper introduces MCircKE, a novel knowl-  
589 edge editing framework grounded in the principles  
590 of mechanistic interpretability to address the preva-  
591 lent "reasoning gap." Through rigorous component-  
592 level analysis, we demonstrated that the failure of  
593 existing methods stems from the neglect of reason-  
594 ing circuits required to route new information dur-  
595 ing multi-hop inference. To address this, MCircKE  
596 dynamically maps these transient reasoning path-  
597 ways and performs surgical interventions within the  
598 identified circuits. Extensive evaluations demon-  
599 strate the effectiveness of MCircKE in multi-hop  
600 factual recall for knowledge editing in LLMs.

## 601 Limitations

602 **Sequential Editing and Interference** Our cur- 648  
603 rent evaluation focuses on single edits. It remains 649  
604 an open question whether the proposed framework 650  
605 has sufficient capacity to support thousands of se- 651  
606 quential edits without saturation or cross-talk (in- 652  
607 terference between different edited facts sharing 653  
608 the same routing heads).

609 **Model Architecture** Our experimental valida- 654  
610 tion is currently limited to the GPT family of 655  
611 models (GPT-2 and GPT-J). While these architec- 656  
612 tures are standard benchmarks in mechanistic in- 657  
613 terpretability, extending our analysis to more archi- 658  
614 tectures, such as Llama or Qwen, would provide 659  
615 stronger evidence for the generalizability and ro- 660  
616 bustness of MCircKE.

617 **Circuit Discovery Cost and Trade-off** Calculat- 661  
618 ing EAP-IG requires computing gradients for every 662  
619 edge, which, while efficient compared to activa- 663  
620 tion patching, still incurs a computational overhead 664  
621 ( $O(m)$  backward passes) that may be prohibitive 665  
622 for real-time editing of extremely large models. 666  
623 Besides, EAP-IG serves as a gradient-based ap- 667  
624 proximation to exact activation patching. The fi- 668  
625 delity of this approximation is governed by the 669  
626 number of integration steps  $m$ . While reducing  $m$  670  
627 improves computational efficiency, it may hinder 671  
628 convergence to the true attribution scores, result- 672  
629 ing in the identification of incomplete circuits that 673  
630 do not faithfully represent the model’s reasoning 674  
631 topology. 675

632 **Threshold Sensitivity** The top- $K$  pruning strat- 676  
633 egy imposes a fixed sparsity level across all edits. 677  
634 This heuristic may be suboptimal because complex 678  
635 facts may require larger circuits while simple facts 679  
636 may require smaller ones, and these requirements 680  
637 can also vary across model scales, necessitating 681  
638 future work on adaptive thresholding.

## 639 References

640 Leonard Bereska and Efstratios Gavves. 2024. Mech- 682  
641 anistic interpretability for ai safety—a review. *arXiv* 683  
642 *preprint arXiv:2404.14082*.

643 Arthur Conmy, Augustine Mavor-Parker, Aengus Lynch, 684  
644 Stefan Heimersheim, and Adrià Garriga-Alonso. 685  
645 2023. Towards automated circuit discovery for mech- 686  
646 anistic interpretability. *Advances in Neural Informa- 687*  
647 *tion Processing Systems*, 36:16318–16352. 688

Junfeng Fang, Houcheng Jiang, Kun Wang, Yunshan 648  
Ma, Jie Shi, Xiang Wang, Xiangnan He, and Tat- 649  
Seng Chua. 2025. Alphaedit: Null-space constrained 650  
knowledge editing for language models. In *The Thir- 651*  
*teenth International Conference on Learning Repre- 652*  
*sentations*. 653

Michael Hanna, Ollie Liu, and Alexandre Variengien. 654  
2023. How does gpt-2 compute greater-than?: In- 655  
terpreting mathematical abilities in a pre-trained lan- 656  
guage model. *Advances in Neural Information Pro- 657*  
*cessing Systems*, 36:76033–76060. 658

Michael Hanna, Sandro Pezzelle, and Yonatan Belinkov. 659  
Have faith in faithfulness: Going beyond circuit over- 660  
lap when finding model mechanisms. In *ICML 2024 661*  
*Workshop on Mechanistic Interpretability*. 662

Dan Hendrycks, Collin Burns, Steven Basart, Andy 663  
Zou, Mantas Mazeika, Dawn Song, and Jacob Stein- 664  
hardt. 2021. Measuring massive multitask language 665  
understanding. *Proceedings of the International Con- 666*  
*ference on Learning Representations (ICLR)*. 667

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan 668  
Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, 669  
Weizhu Chen, and 1 others. 2022. Lora: Low-rank 670  
adaptation of large language models. *ICLR*, 1(2):3. 671

Kevin Meng, David Bau, Alex Andonian, and Yonatan 672  
Belinkov. 2022. Locating and editing factual associa- 673  
tions in gpt. *Advances in neural information process- 674*  
*ing systems*, 35:17359–17372. 675

Kevin Meng, Arnab Sen Sharma, Alex J Andonian, 676  
Yonatan Belinkov, and David Bau. 2023. Mass- 677  
editing memory in a transformer. In *The Eleventh 678*  
*International Conference on Learning Representa- 679*  
*tions*. 680

Neel Nanda. 2023. Attribution patching: Activation 681  
patching at industrial scale. URL: [https://www.neel- 682](https://www.neel-nanda.io/mechanistic-interpretability/attribution-patching)  
[nanda.io/mechanistic-interpretability/attribution- 683](https://www.neel-nanda.io/mechanistic-interpretability/attribution-patching)  
[patching](https://www.neel-nanda.io/mechanistic-interpretability/attribution-patching). 684

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, 685  
Dario Amodei, Ilya Sutskever, and 1 others. 2019. 686  
Language models are unsupervised multitask learn- 687  
ers. *OpenAI blog*, 1(8):9. 688

Daking Rai, Yilun Zhou, Shi Feng, Abulhair Saparov, 689  
and Ziyu Yao. 2024. A practical review of mech- 690  
anistic interpretability for transformer-based language 691  
models. *arXiv preprint arXiv:2407.02646*. 692

Lee Sharkey, Bilal Chughtai, Joshua Batson, Jack Lind- 693  
sey, Jeff Wu, Lucius Bushnaq, Nicholas Goldowsky- 694  
Dill, Stefan Heimersheim, Alejandro Ortega, Joseph 695  
Bloom, and 1 others. 2025. Open problems 696  
in mechanistic interpretability. *arXiv preprint 697*  
*arXiv:2501.16496*. 698

Alon Talmor, Jonathan Herzig, Nicholas Lourie, and 699  
Jonathan Berant. 2019. CommonsenseQA: A ques- 700  
tion answering challenge targeting commonsense 701  
knowledge. In *Proceedings of the 2019 Conference 702*

of the North American Chapter of the Association for Computational Linguistics: *Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.

Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. Interpretability in the wild: a circuit for indirect object identification in gpt-2 small. In *The Eleventh International Conference on Learning Representations*.

Peng Wang, Zexi Li, Ningyu Zhang, Ziwen Xu, Yunzhi Yao, Yong Jiang, Pengjun Xie, Fei Huang, and Hua-jun Chen. 2024a. Wise: Rethinking the knowledge memory for lifelong model editing of large language models. *Advances in Neural Information Processing Systems*, 37:53764–53797.

Peng Wang, Ningyu Zhang, Bozhong Tian, Zekun Xi, Yunzhi Yao, Ziwen Xu, Mengru Wang, Shengyu Mao, Xiaohan Wang, Siyuan Cheng, and 1 others. 2024b. Easyedit: An easy-to-use knowledge editing framework for large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pages 82–93.

Song Wang, Yaochen Zhu, Haochen Liu, Zaiyi Zheng, Chen Chen, and Jundong Li. 2024c. Knowledge editing for large language models: A survey. *ACM Computing Surveys*, 57(3):1–37.

Yunzhi Yao, Jizhan Fang, Jia-Chen Gu, Ningyu Zhang, Shumin Deng, Huajun Chen, and Nanyun Peng. 2025. Cake: Circuit-aware editing enables generalizable knowledge learners. *CoRR*.

Ningyu Zhang, Yunzhi Yao, Bozhong Tian, Peng Wang, Shumin Deng, Mengru Wang, Zekun Xi, Shengyu Mao, Jintian Zhang, Yuansheng Ni, and 1 others. 2024. A comprehensive study of knowledge editing for large language models. *arXiv preprint arXiv:2401.01286*.

Zhuoran Zhang, Yongxiang Li, Zijian Kan, Keyuan Cheng, Lijie Hu, and Di Wang. Locate-then-edit for multi-hop factual recall under knowledge editing. In *Forty-second International Conference on Machine Learning*.

Zexuan Zhong, Zhengxuan Wu, Christopher D Manning, Christopher Potts, and Danqi Chen. 2023. [Mquake: Assessing knowledge editing in language models via multi-hop questions](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 15686–15702.

## A Appendix

### A.1 Licenses

This study utilizes several open-source artifacts, all of which are distributed under permissive licenses compatible with academic research.

- MQuAKE-3K Dataset: Distributed under the MIT License.

- GPT-2 (Large/XL): Released by OpenAI under the Modified MIT License.

- GPT-J (6B): Released by EleutherAI under the Apache-2.0 License.

- CaKE (Baseline): The official implementation and data are distributed under the MIT License.

- ROME/MEMIT/AlphaEdit (Baselines): The official codebases for these methods are distributed under the MIT License.

### A.2 Intended Use Consistency

All artifacts employed in this study were used in a manner consistent with their intended purposes as defined by their creators.

- MQuAKE-3K: This dataset was specifically designed to evaluate the multi-hop reasoning capabilities of edited language models. Our usage strictly adheres to this evaluation protocol.

- Language Models (GPT-2, GPT-J): These models are intended for research into the properties of large language models, including interpretability and fine-tuning. Our work, which investigates the internal mechanisms of these models via EAP-IG and applies circuit-based editing, aligns with this research scope.

- Synthetic Data: The synthetic prompts used for the MCircKE were generated using ChatGPT (OpenAI). This usage falls within the acceptable use policy for research purposes, specifically for generating non-sensitive, factual reasoning chains.

### A.3 Documentation of Artifacts

**Language & Domain:** All models and datasets used in this study process English language text. The domain of the MQuAKE-3K dataset is factual knowledge graphs (derived from Wikidata), covering entities such as people, organizations, locations, and creative works.

**Dataset Characteristics:** The MQuAKE-3K dataset consists of 3,000 counterfactual editing instances. Each instance includes a base edited fact

801 and a set of multi-hop questions (2-hop, 3-hop, 4-  
802 hop) that test the model’s ability to propagate this  
803 edit.

804 **Model Specifications:** GPT-2 Large - 774M pa-  
805 rameters, 36 layers, 1280 hidden dimension; GPT-2  
806 XL - 1.5B parameters, 48 layers, 1600 hidden di-  
807 mension; GPT-J - 6B parameters, 28 layers, 4096  
808 hidden dimension

809 **Demographics:** The datasets are derived from  
810 Wikipedia and Wikidata, which may contain inher-  
811 ent biases reflecting the demographics of English-  
812 speaking internet users. No human subjects were  
813 involved in this specific study.