

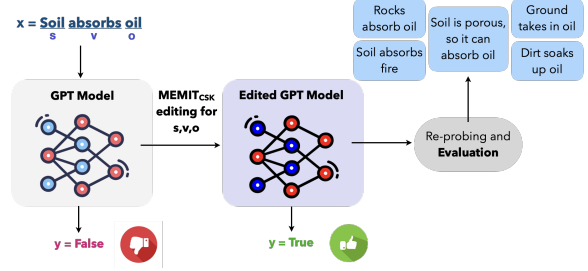
# Editing Commonsense Knowledge in GPT

Anshita Gupta<sup>♣\*</sup> Debanjan Mondal<sup>♣\*</sup> Akshay Krishna Sheshadri<sup>♣\*</sup>  
Wenlong Zhao<sup>♣</sup> Xiang Lorraine Li<sup>♡</sup> Sarah Wiegreffe<sup>♡</sup> Niket Tandon<sup>♡</sup>

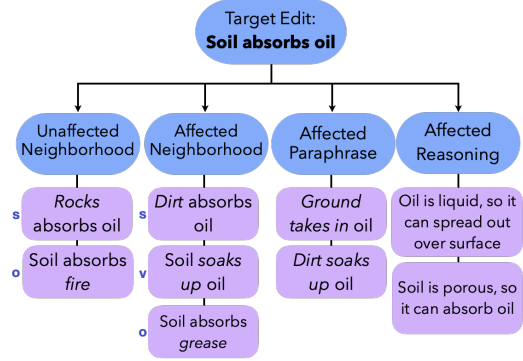
<sup>♣</sup>University of Massachusetts Amherst, <sup>♡</sup>Allen Institute for AI  
{anshitagupta,debanjanmond,asheshadri,wenlongzhao}@cs.umass.edu  
{lorrainel,sarahw,nikett}@allenai.org

## Abstract

Memory editing methods for updating encyclopedic knowledge in transformers have received increasing attention for their efficacy, specificity, and generalization advantages. However, it remains unclear if such methods can be adapted for the more nuanced domain of commonsense knowledge. We propose MEMIT<sub>CSK</sub>, an adaptation of MEMIT to edit commonsense mistakes in GPT-2 Large and XL. We extend editing to various token locations and employ a robust layer selection strategy. Models edited by MEMIT<sub>CSK</sub> outperforms the fine-tuning baselines by 10.97% and 10.73% F1 scores on subsets of PEP3k and 20Q. We further propose a novel evaluation dataset, MEMIT-CSK-PROBE, that contains unaffected neighborhood, affected neighborhood, affected paraphrase, and affected reasoning challenges. MEMIT<sub>CSK</sub> demonstrates favorable semantic generalization, outperforming fine-tuning baselines by 13.72% and 5.57% overall scores on MEMIT-CSK-PROBE. These results suggest a compelling future direction of incorporating context-specific user feedback concerning commonsense in GPT by direct model editing, rectifying and customizing model behaviors via human-in-the-loop systems.<sup>1</sup>



(a) Proposed framework for probing, editing (MEMIT<sub>CSK</sub>), and evaluating commonsense knowledge in GPT models.



(b) Samples from MEMIT-CSK-PROBE for evaluating semantic generalization.

Figure 1: Method overview.

## 1 Introduction

Transformer models have achieved great success in modeling language distributions and performing diverse natural language processing tasks (Brown et al., 2020). However, they exhibit factual mistakes (Lewis et al., 2021; Shuster et al., 2021), commonsense mistakes (Bender and Koller, 2020; Marcus, 2021; Talmor et al., 2019; Bhargava and Ng, 2022), and consistency errors (Tam et al., 2022; Devaraj et al., 2022; Weng et al., 2020).

Unlike factual mistakes that can be addressed by retrieval augmented methods (Madaan et al., 2022),

common sense is hard to enumerate or retrieve from encyclopedic resources. Moreover, common sense is contextual and its plausibility depends on locations and scenarios. To correct and customize model behaviors about common sense, a more natural evidence source is user feedback based on the application contexts. Assuming the availability of such feedback, how can we incorporate the changes, modify relevant knowledge in the model, while not affecting similar but independent knowledge?

Retraining models on large data is expensive and undesired, while finetuning on the new feedback is prone to overfitting. To address similar challenges in updating encyclopedic knowledge in transformer models, researchers have discovered

<sup>\*</sup>Equal contribution, alphabetical order.

<sup>1</sup>Our code and data is available in [github](#).

that its predictions strongly correlate with certain neuron activations and shown that directly editing relevant model parameters balances efficacy, specificity, and generalization (Meng et al., 2023, 2022).

However, it remains unclear whether these causal mediation analyses and editing methods can scale beyond encyclopedic knowledge to handle the broader and more nuanced domain of commonsense knowledge. Encyclopedic knowledge is about instances and often exhibits a strong correspondence between a subject and an object, e.g., *Eiffel Tower is located in Rome*. Previous works observe that such knowledge is generally stored in MLP layers (Geva et al., 2021). MEMIT (Meng et al., 2022) edits MLP parameters corresponding to the subject tokens to correct the predictions of objects. On the contrary, commonsense knowledge is about concepts, and a subject-verb pair can match many objects to form plausible statements, e.g., *‘people eat’* can be followed by various types of food. Further research and experimentation are therefore necessary to show if common sense is mediated by particular neuron activations in GPT models and to assess the feasibility and effectiveness of adapting editing methods for updating commonsense knowledge in them.

In this paper, we judiciously instantiate the recently proposed MEMIT method (Meng et al., 2022) for editing commonsense knowledge in GPT-2 Large and XL via careful model selection (hereby referred to as MEMIT<sub>CSK</sub>). We focus on the task of classifying statements consisting of a subject, a verb, and an object as plausible or implausible, using the 20Q and PEP3k datasets. In causal mediation analyses, we corrupt each of the three token types due to their critical importance in plausibility assessment. We observe strong causal relations between neural activations and predictions of models fine-tuned on the task, although no clear pattern appears in pretrained models that perform poorly on the task. Then, we apply MEMIT<sub>CSK</sub> to each token type, editing a set of layers whose cardinality and locations are selected according to the causal analysis in a principled manner.

We construct three inference sets for both 20Q and PEP3k to evaluate our MEMIT<sub>CSK</sub> method and finetuning. Inference Set #1 (Inf\_1) and #2 (Inf\_2) are subsets of the original datasets, while Inference Set #3 (Inf\_3, Figure 1b) is our created data that consist of efficacy, unaffected neighborhood (specificity), affected neighborhood, affected paraphrase,

and affected reasoning challenges. We find hyper-parameters for both methods to maximize the F1 of the updated models on Inf\_1 and apply the same updating hyper-parameters to Inf\_2 and Inf\_3.

Overall, our main contributions are:

- We propose MEMIT<sub>CSK</sub>, an adaptation of MEMIT to edit commonsense mistakes. Our editing strategy extends to subject, object, and verb editing — a commonsense-specific adaptation. We also employ a judicious hyper-parameter search and a more robust layer selection that could be extended to other tasks and domains.
- GPT-2 XL edited by MEMIT<sub>CSK</sub> outperforms fine-tuned baselines by 10.97% and 10.73% F1 scores on the Inf\_2 of PEP3k and 20Q. This shows that MEMIT<sub>CSK</sub> hyper-parameters selected on Inf\_1 can be transferred to Inf\_2.
- We introduce a novel dataset, MEMIT-CSK-PROBE (Inf\_3), to evaluate semantic generalization of methods for updating commonsense knowledge in models. It consists of efficacy, unaffected neighborhood, affected neighborhood, affected paraphrase, and affected reasoning challenges.
- GPT-2 XL edited by MEMIT<sub>CSK</sub> outperforms fine-tuned baselines by 13.72% and 5.57% overall scores on the Inf\_3 of PEP3k and 20Q, and achieved better semantic generalization.

## 2 Background

### 2.1 Causal Tracing

For editing a model, we need to identify the editing layers. We utilized causal tracing for this purpose (Meng et al., 2023). Given a model and a set of correctly predicted prompts, causal tracing involves three steps:

- **Clean run:** We pass a prompt  $x$  into the model and collect all hidden activations  $h_i^{(l)} \mid i \in [1, T], l \in [1, L]$  where  $T$  is number of input tokens and  $L$  is number of layers in the model.
- **Corrupted run:** The target token embeddings are corrupted. After the prompt  $x$  is embedded as  $[h_1^{(0)}, h_2^{(0)} \dots h_T^{(0)}]$ , we set  $h_i^{(0)} := h_i^{(0)} + \epsilon$ , for all token indices  $i$  belonging to the target token, where  $\epsilon \sim \mathcal{N}(0, v)$  and we select  $v$  as

three times the empirical standard deviation of the embeddings corresponding to the tokens we want to corrupt.

- **Corrupted-with-restoration-run:** The model runs computations on the noisy embeddings as in the corrupted baseline, except at some token  $\hat{i}$  and layer  $\hat{l}$ . There, the model is ‘hooked’ such that it’s forced to output the clean state activation  $h_{\hat{i}}^{\hat{l}}$ , and future computations go without intervention.

Let  $r$  be the correct prediction. Let  $\mathbb{P}[r]$ ,  $\mathbb{P}_*[r]$  and  $\mathbb{P}_{*, \text{clean } h_i^l}[r]$  be the probability of  $f$  in the clean, corrupted and corrupted-with-restoration-runs respectively. The *total effect (TE)* is calculated as:

$$\text{TE} = \mathbb{P}[r] - \mathbb{P}_*[r] \quad (1)$$

The *indirect effect (IE)* of a particular hidden state  $h_i^l$  is calculated as:

$$\text{IE} = \mathbb{P}_{*, \text{clean } h_i^l}^{(l)}[r] - \mathbb{P}_*[r] \quad (2)$$

After averaging over all the prompts, we get the *average total effect (ATE)* and *average indirect effect (AIE)* for each hidden state. For MLP and attention, instead of just patching the output with clean run value, we take a window of 10 layers around that site and patch all the MLP/attention outputs in that window.

**Severed Causal Tracing.** To obtain a clearer understanding of the impact of MLP layers, we perform causal tracing analysis with a modified causal graph, again following the footsteps of MEMIT. In this case, in the corrupted-with-restoration-run, we freeze the MLP output to the corrupted run value, so that it’s unaffected by the inserting of clean state  $h_i^l$ . This can be viewed as severing the MLP computations from the original computation graph. We performed similar analysis with severed attention computations.

## 2.2 Memory Editing

We perform causal mediation analysis and severed causal mediation analysis as detailed above to identify the critical MLP layers ( $l \in \mathcal{R}$ ), and then extend MEMIT for model editing (Meng et al., 2022). Let  $L$  be the last layer in  $\mathcal{R}$ ,  $h_i^L$  be the transformer’s hidden state represented by layer  $L$  at the editing token’s location for the  $i^{\text{th}}$  prompt, and  $G$  be the generative model. For each prompt  $(x_i, y_i) \in \varepsilon$

that we want to edit, we find the vector  $z_i$  (modified hidden state) such that the probability of the correct prediction is maximized. In practice, we calculate the residual  $\delta_i = z_i - h_i^L$  at hidden layer  $L$  using the equation,

$$z_i = h_i^L + \arg \min_{\delta_i} -\log \mathbb{P}_{G(h_i^L + \delta_i)}[y_i | x_i] \quad (3)$$

Once this vector has been found, we spread the residual  $\delta_i$  across the layers  $l \in \mathcal{R}$  as follows,

$$r_i^l = \frac{z_i - h_i^L}{L - l + 1} \quad (4)$$

We treat the MLP layer as a linear layer of associative memory that stores key-value memories. For each MLP layer  $l \in \mathcal{R}$  and for each prompt, we calculate the keys  $k_i^l$  (inputs) and the updated memories  $m_i^l$  (outputs) to be inserted. However, we also need to preserve the already existing correct key-memory pairs. Given these constraints, we obtain the parameter update equation using a close-form solution of a least square minimization problem as follows,

$$\Delta^l = R^l K^{lT} (C^l + K^l K^{lT})^{-1} \quad (5)$$

where  $\Delta^l$  is the overall residual update of the  $l^{\text{th}}$  MLP layer,  $K^l$  are the new keys,  $C^l$  is proportional to covariance of existing keys, and  $R^l$  is the residual of old and new memories.

## 3 Method

While assessing commonsense knowledge, there is no definitive answer to a question, unlike encyclopedic knowledge. For the statement *People eat*, there are multiple correct choices for the next word (*rice*, *meat*, *bread* etc.), all of which may not lie within the most probable predictions. Hence, in MEMIT<sub>CSK</sub> we concentrate on the plausibility of statements. Given a dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$  where each  $x_i$  is a statement and  $y_i \in \{\text{True}, \text{False}\}$ , our objective is to predict the plausibility of the statement  $x$ . We initially fine-tune the GPT-2 Large and GPT-2 XL models on the training set with the next token prediction objective. The input prompts have the structure  $[x_i]: [y_i]$ . These models serve as our base model for subsequent experiments. The motivation behind fine-tuning rather than zero-shot prompting has been described in section 5.5.1.

Next, we performed causal and severed causal tracing analysis on Inference Set #1 with the base

models. MEMIT only focuses on subject corruption, since encyclopedic knowledge is mostly associated with the subject. In contrast, common sense and plausibility depend on each element of the sentence. Therefore, in MEMIT<sub>CSK</sub> we analyzed three types of corruption: subject, object, and verb. In the causal analysis, we only consider prompts where the model accurately predicts the plausibility as (*True* or *False*).

Subsequently, we perform hyper-parameter tuning on the Inference Set #1 for all update methods. **FT** is the baseline fine-tuning update method adjusted to minimize the cross-entropy loss on the last token i.e. the predicted label. The **Edit** method is applied for each possible editing location (subject  $s$ , verb  $v$ , object  $o$ ) by selecting edit layers detailed in Section 5.2. Both FT and editing update methods are only applied on the incorrectly predicted subset of Inference Set #1. To minimize undesirable modifications, the hyper-parameters are adjusted to maximize the updated F1 score on the entire dataset, rather than solely maximizing the efficacy of inaccurate predictions. This approach will lead to updates that maintain balance between high efficacy and minimal relapse, ensuring that previously correct statements remain accurate. After identifying the best hyper-parameters, the same configuration is applied to Inference Set #2 to measure the scope of configuration generalization. Finally, we evaluate all update methods on the challenging Inference Set #3 and measure efficacy, paraphrase, neighborhood and reasoning accuracy.

## 4 Datasets and Evaluation

We conduct experiments on two common sense datasets after appropriate pre-processing and language normalization: PEP3k and 20Q, following Porada et al. (2021). The datasets consist of commonsense statements in the triplet format (subject  $s$ , verb  $v$ , object  $o$ ) and ground truth plausibility labels. The dataset consist of three splits: Training Set, Inference Set #1, and Inference Set #2. The Inference Set #2 is the holdout dataset split and was used only for experiment evaluations. The dataset statistics are given in Table 1.

To compare and assess semantic generalization for multiple update strategies in the common sense domain, we build a new subset from Inference Set #2 consisting of challenging commonsense statements called Inference Set #3. The subset categories were created using zero-shot instructions to

the GPT-3 API. The dataset statistics are given in Table 2.

Dataset	$N_{\text{Train}}$	$N_{\text{Inf}_1}$	$N_{\text{Inf}_2}$
<b>PEP3k</b>	1,225	306	1,531
<b>20Q</b>	2,006	507	2,548

Table 1: Number of samples in the Training Set, Inference Set #1, and Inference Set #2 of PEP3k and 20Q.

Type	$N_{\text{PEP3k}}$	$N_{\text{20Q}}$
Original statement	265	381
Unaffected subject neighborhood	1,325	1,894
Unaffected object neighborhood	1,325	1,900
Affected subject neighborhood	1,290	1,856
Affected verb neighborhood	1,288	1,832
Affected object neighborhood	1,292	1,848
Affected paraphrase	1,323	1,905
Affected reasoning	530	754

Table 2: Number of samples in the Inference Set #3 of PEP3k and 20Q.

Provide 5 paraphrases of: Furnishings make noise

-----

1. Furniture can be noisy.
2. Furniture can create sound.
3. Furniture can produce noise.
4. Furniture can be a source of sound.
5. Furniture can be a source of noise.

Figure 2: Prompt to generate affected paraphrase for “Furnishings make noise (false)”

**Unaffected Neighborhood.** To measure the specificity for each  $\{s, o\}$ , generate a set of hyponyms as  $(s', v, o)$  and  $(s, v, o')$ . The score is the percentage of post-update predictions  $\arg \max \mathbb{P}(s', v, o)$  and  $\arg \max \mathbb{P}(s, v, o')$  that remain equivalent to pre-update predictions.

**Affected Neighborhood.** To measure the affected change on similar meaning prompts for each  $\{s, v, o\}$ , generate a set of synonyms as  $(s', v, o)$ ,  $(s, v', o)$  and  $(s, v, o')$ . The score is the percentage of post-update predictions  $\arg \max \mathbb{P}(s', v, o)$ ,  $\arg \max \mathbb{P}(s, v', o)$  and  $\arg \max \mathbb{P}(s, v, o')$  which are equal to the ground truth labels.



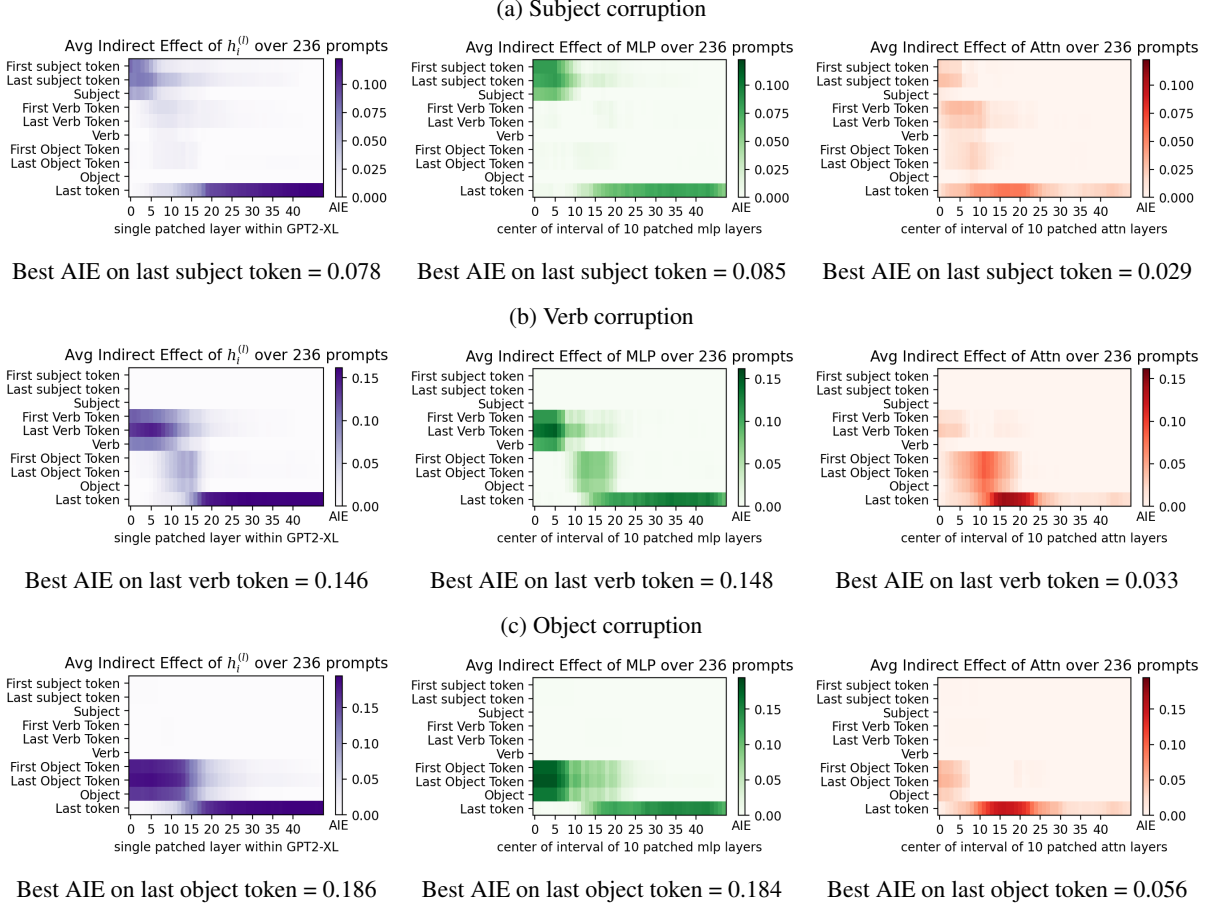


Figure 3: Causal tracing results for GPT-2 XL on PEP3k Inference Set #1.

**Affected Paraphrase.** To measure the affected change on similar meaning prompts, generate a set of paraphrases as  $(s', v', o')$ . The score is the percentage of post-update prediction  $\arg \max \mathbb{P}(s', v', o')$  which are equal to the ground truth labels.

**Affected Reasoning.** Entities and knowledge are interconnected, and updating one component of commonsense knowledge generally requires updating another. To measure the updated model’s connectivity, generate a two-step chain of valid reasoning prompts  $\{R_1, R_2\}$ . The score is the percentage of post-update predictions  $\arg \max \mathbb{P}(R_1)$  and  $\arg \max \mathbb{P}(R_2)$  which are equal to the *True* label.

## 5 Experiments

### 5.1 Causal Mediation Analysis

One representative result of causal tracing analysis is presented in Figure 3. We noticed high AIE at the later layers of the last token as expected, since fixing hidden states or MLPs in those layers re-

stores most of the required information. We also observed strong AIE at the earlier layers for the tokens we are corrupting. This finding is non-trivial and emphasizes the importance of earlier layers while predicting plausibility. Furthermore, the AIE at the last corrupted token is more pronounced than the first corrupted token. These patterns are consistent across all models and datasets. Therefore, we focus on the last subject, last verb, and last object token editing in subsequent experiments. Causal tracing graphs for other datasets and models are present in Appendix A.1.

### 5.2 Severed Causal Analysis

Figure 4 compares the average AIE at last corrupted token for unmodified, severed MLP, and severed Attention causal graphs. We clearly notice a gap in AIE for unmodified and severed MLP graphs at the earlier layers, while for severed Attn this gap is absent. This observation confirms the crucial roles of MLP layers while predicting plausibility and indicate that the restoration effect doesn’t depend on the MLP activity for higher layers. This outcome

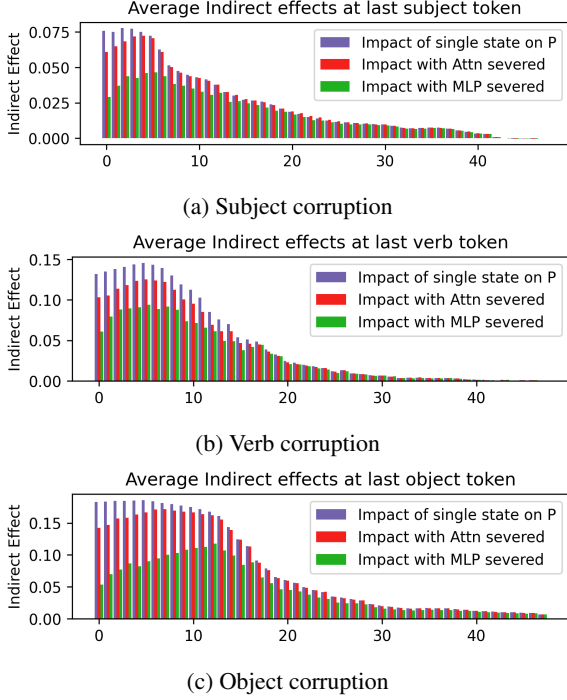


Figure 4: Severed causal tracing results for GPT-2 XL on PEP3k Inference Set #1.

Model	Edit Token	Layer with Max IE	Layers (5) with Max Moving Avg IE	Layers (3) with Max Moving Avg IE
GPT-2 Large	Last Subject	8	8,9,10,11,12	8,9,10
	Last Verb	4	4,5,6,7,8	4,5,6
	Last Object	12	10,11,12,13,14	11,12,13
GPT-2 XL	Last Subject	5	2,3,4,5,6	4,5,6
	Last Verb	5	3,4,5,6,7	5,6,7
	Last Object	12	9,10,11,12,13	10,11,12

Table 3: Different strategies to pick edit layers for PEP3k Inference Set #1.

is consistent with the observations in MEMIT for encyclopedic knowledge.

However, in contrast to encyclopedic facts, we witness the layer with peak severed MLP AIE appear much earlier. This demonstrates the decisive role of early and early-middle layers in commonsense predictions. Instead of choosing the layer with the highest AIE and a window before that as the edit layers, we also consider the window with the **maximum moving average of AIE**. We also broadened our search space by shifting these windows by one layer on both sides. During hyper-parameter tuning, we compared all three and five layer windows within the aforementioned search space and picked the best performing strategy. Moreover, in the object corruption plot, we observed a small peak in the initial layers, followed by the highest AIE. Hence, we also included the

initial layer windows for the object editing experiments. Table 3 shows the potential edit locations for PEP3k Inference Set #1.

### 5.3 Configuration Generalization

Can we transfer hyper-parameters to a new independent sample dataset? The detailed experimental results with the best hyper-parameter updates are showcased in Table 4<sup>2</sup>. The best hyper-parameters for each method are detailed in Appendix A.2

We observe different patterns for the best editing update in the two datasets. The verb token editing performs better for the PEP3k dataset, while the object token editing performs better for the 20Q dataset. The FT method has the highest updated F1 score and efficacy for the smaller model, while the editing method has the highest scores for the larger model. The relapsed score is generally lower for all editing methods when compared to the FT method.

After transferring the best hyper-parameters to Inference Set #2, the editing method consistently performs well, while the FT method has very low efficacy and a modest increase in the updated F1 score. All three editing methods on  $(s, v, o)$  significantly outperforms the FT method, with patterns consistent with the Inference Set #1 results. The editing method displays high configuration generalization and hence is potentially applicable to new datasets with similar features.

### 5.4 Efficacy and Semantic Generalization

Can we effectively measure editing performance for difficult commonsense knowledge through its semantic neighborhood, paraphrase, and reasoning? The detailed evaluation results on Inference Set #3 are shown in Table 5<sup>1</sup> and 6<sup>1</sup>. The **Overall Score** column is the arithmetic mean of all evaluation metrics. Note that, the initial efficacy of the base model on this difficult dataset is 0%.

All editing methods outperform the FT baseline method in terms of overall score. Results show that no matter which token is edited  $(s, v, o)$ , the editing method achieves better results than FT. This provides evidence that model editing is a powerful tool to enhance a model’s performance in the common sense domain.

The smaller edited model performs better than the larger model for the efficacy and overall score

<sup>2</sup>The KL divergence factor has not been investigated extensively and could provide performance uplift.

Model	Update Method	Edit Token	Edit Layers	F1 Score %	Efficacy %	Relapse %	F1 Score %	Efficacy %	Relapse %
<b>PEP3k</b>				Inf_1			Inf_2		
GPT-2 Large	Base M	-	-	75.16	0	0	76.22	0	0
	FT	-	-	<b>95.75 (+20.59)</b>	94.74	<b>3.91</b>	80.92 (+4.70)	40.93	<b>6.60</b>
	Edit	Last Subject	4,5,6,7,8	88.53 (+13.36)	76.32	7.39	79.36 (+3.14)	54.95	12.77
	Edit	Last Verb	4,5,6,7,8	93.78 (+18.62)	<b>96.05</b>	6.96	<b>89.08 (+12.86)</b>	<b>93.68</b>	12.34
	Edit	Last Object	1,2,3,4,5	88.41 (+13.25)	86.84	10.87	77.65 (+1.43)	78.57	21.85
				Inf_1			Inf_2		
GPT-2 XL	Base M	-	-	77.12	0	0	76.47	0	0
	FT	-	-	90.16 (+13.05)	<b>97.14</b>	11.87	80.93 (+4.46)	50.83	9.82
	Edit	Last Subject	1,2,3,4,5	90.51 (+13.39)	80	6.36	84.72 (+8.25)	77.22	12.98
	Edit	Last Verb	6,7,8	<b>95.09 (+17.97)</b>	92.86	<b>4.24</b>	<b>91.90 (+15.43)</b>	<b>88.33</b>	<b>7.00</b>
	Edit	Last Object	3,4,5	94.43 (+17.32)	91.43	4.66	86.69 (+10.22)	72.78	8.97
<b>20Q</b>				Inf_1			Inf_2		
GPT-2 Large	Base M	-	-	72.39	0	0	74.07	0	0
	FT	-	-	<b>91.32 (+18.93)</b>	<b>97.86</b>	11.17	76.45 (+2.37)	48.23	13.69
	Edit	Last Subject	3,4,5	85.33 (+12.94)	75	10.63	81.97 (+7.90)	67.18	12.66
	Edit	Last Verb	2,3,4,5,6	77.64 (+5.25)	38.57	<b>7.36</b>	77.33 (+3.26)	33.44	<b>7.22</b>
	Edit	Last Object	1,2,3	87.09 (+14.71)	82.14	10.90	<b>84.61 (+10.54)</b>	<b>80.43</b>	13.79
				Inf_1			Inf_2		
GPT-2 XL	Base M	-	-	74.73	0	0	75.77	0	0
	FT	-	-	85.71 (+10.98)	80.46	12.40	77.36 (+1.60)	30.97	7.8
	Edit	Last Subject	2,3,4,5,6	<b>92.31 (+17.58)</b>	79.69	<b>3.43</b>	86.46 (+10.70)	65.73	<b>6.90</b>
	Edit	Last Verb	3,4,5,6,7	82.64 (+7.91)	44.53	4.49	79.03 (+3.27)	35.91	7.11
	Edit	Last Object	1,2,3	91.12 (+16.39)	<b>89.06</b>	8.18	<b>88.09 (+12.33)</b>	<b>76.60</b>	8.21

Table 4: Configuration generalization experiment results for the PEP3k and 20Q datasets.

metrics. This result indicates that for small datasets, model editing is more effective for smaller models.

The FT baseline performs well for the unaffected neighborhood but performs poorly for the affected neighborhood, paraphrase, and reasoning statements. These results indicate that the FT method is not able to generalize across similar neighborhoods. The correction mechanism in FT is local and doesn’t spread into nearby semantic vicinities. While the editing method shows significant improvement in terms of the affected metrics. This demonstrates that the corrections by editing method are more general and widespread. For the unaffected neighborhood, the smaller edited model’s performance is similar to the FT method, while for the larger model, the editing method outperforms the FT method.

Except for the edited GPT-2 XL model for the PEP3k dataset, the edited object models have the highest accuracy for affected reasoning among all update methods. For the editing method, the accuracy for the affected neighborhood is lower for the token that is being edited. The accuracy for the affected object neighborhood is very low while editing the object token, and a consistent pattern is

noticed for subject and verb tokens across datasets and models.

## 5.5 Further Analyses

### 5.5.1 Zero-Shot vs. FT

To successfully apply the MEMIT algorithm for editing, it is necessary to have a clear causal pattern and high AIE. Our initial experiments were performed for a zero-shot prompted model. The structure of zero-shot plausibility prompts was “[Statement]. True or False?” and the predicted label was determined by comparing  $\mathbb{P}(\text{True})$  and  $\mathbb{P}(\text{False})$ . The causal graphs and accuracy statistics for Inference Set #1 are in Figure 5<sup>3</sup>. The causal graphs for the fine-tuned model are more distinct and larger in scale, having a higher IE and strong indication of successful edit locations. In contrast, the zero-shot causal graph has no localized pattern for editing. Hence the fine-tuned model is utilized for further experiments.

<sup>3</sup>Normalization experiments were performed with domain conditional Pointwise Mutual Information (PMI) based on (Holtzman et al., 2021) but did not result in any significant improvement, the accuracy for  $\text{PMI}_{\text{DC}} = 52.94\%$ .

Model	Update Method	Edit Token	Efficacy %	Unaffected Neighborhood %		Affected Neighborhood %			Affected Paraphrase %	Affected Reasoning %	Overall Score
				Subject	Object	Subject	Verb	Object			
GPT-2 Large	Base M	-	0	100	100	19.77	18.94	23.84	34.16	32.08	41.10
	FT	-	30.57	83.92	82.64	30.93	30.67	35.53	38.85	33.77	45.86
	Edit	Last Subject	58.87	<b>79.01</b>	69.51	27.05	50.70	54.26	45.73	40.57	53.21
	Edit	Last Verb	<b>96.23</b>	44.15	48.06	69.92	38.28	<b>83.82</b>	41.65	33.02	56.89
	Edit	Last Object	82.26	44.15	<b>73.28</b>	<b>71.78</b>	<b>69.25</b>	36.15	<b>53.74</b>	<b>46.04</b>	<b>59.58</b>
GPT-2 XL	Base M	-	0	100	100	21.01	23.45	23.3	33.64	31.13	41.57
	FT	-	39.63	77.21	76.98	33.88	37.5	39.16	39.91	55.66	49.99
	Edit	Last Subject	72.08	<b>81.21</b>	<b>64.15</b>	27.75	<b>56.83</b>	58.67	<b>47.92</b>	<b>35.85</b>	<b>55.56</b>
	Edit	Last Verb	<b>87.92</b>	59.24	57.06	57.36	34.47	<b>71.05</b>	37.41	31.13	54.46
	Edit	Last Object	75.47	58.11	57.82	<b>57.52</b>	54.35	27.86	43.99	34.72	51.23

Table 5: Efficacy and semantic generalization experiment results for the PEP3k Inference Set #3.

Model	Update Method	Edit Token	Efficacy %	Unaffected Neighborhood %		Affected Neighborhood %			Affected Paraphrase %	Affected Reasoning %	Overall Score
				Subject	Object	Subject	Verb	Object			
GPT-2 Large	Base M	-	0	100	100	30.23	22.93	27.11	32.76	27.72	42.59
	FT	-	29.66	88.07	87.05	39.17	37.34	39.23	42.05	27.59	48.77
	Edit	Last Subject	67.98	79.57	57.79	35.08	63.26	<b>63.91</b>	<b>53.96</b>	31.70	56.66
	Edit	Last Verb	32.55	<b>89.55</b>	84.16	37.93	26.80	46.10	35.17	28.51	47.60
	Edit	Last Object	<b>81.89</b>	66.95	<b>85.32</b>	<b>71.98</b>	<b>71.67</b>	34.79	51.60	<b>35.68</b>	<b>62.49</b>
GPT-2 XL	Base M	-	0	100	100	33.02	24.78	29.38	33.70	39.38	45.03
	FT	-	21.52	87.01	86.89	36.21	34.72	37.01	38.27	35.41	47.13
	Edit	Last Subject	61.94	87.22	71.00	34.48	57.47	<b>58.82</b>	<b>51.75</b>	37.93	57.58
	Edit	Last Verb	35.70	<b>90.18</b>	83.21	37.44	29.80	47.25	35.91	37.80	49.66
	Edit	Last Object	<b>72.18</b>	76.24	<b>92.26</b>	<b>61.90</b>	<b>60.37</b>	33.76	47.87	<b>41.90</b>	<b>60.81</b>

Table 6: Efficacy and semantic generalization experiment results for the 20Q Inference Set #3.

### 5.5.2 Causal Mediation Analysis after Editing

From the preliminary analysis in Figure 3, we found that the editing method is most effective when there is a clear causal pattern and a high AIE. Therefore, we anticipate seeing improvements in causal graphs after editing. To measure the changes, we re-conduct causal mediation analysis for each noised token ( $s, v, o$ ) using successfully edited statements.

Figure 6 shows the causal graphs using the verb token edited model. If the subject was noised, a slight increase in IE was observed at the verb token. If the verb was noised, a higher IE was observed at the verb token. However, since it is an auto-regressive model, if the object was noised, no change was observed at the verb token. These results confirm that the overall causal pattern and AIE improve for a better model based on the editing strategy.

## 6 Related Work

**Commonsense knowledge.** In the common sense domain, Porada et al. (2021) worked on learn-

ing the plausibility of events represented as (subject  $s$ , verb  $v$ , object  $o$ ) triplets. They fine-tuned a pre-trained RoBERTa model (Liu et al., 2019) with the events and hypernym chains of each argument as input. They showed that injecting lexical knowledge enhances the performance of modeling plausibility but is not sufficient to enforce conceptual consistency. Wang et al. (2019) released a Sense-Making benchmark and performed two tasks on it; multiple fine-tuned models first differentiate between valid and invalid statements and then choose an explanation for the incorrect statements.

**Model Editing.** Elhage et al. (2021); Dar et al. (2022) aimed at understanding and localizing internal mechanisms of Large Language Models. Dai et al. (2022); Yao et al. (2022) proposed altering sparse set of neurons. Inspired by (Anderson, 1972), some researchers focused on the MLP layers and analyzed it as a linear associate memory. In the encyclopedic factual domain, Meng et al. (2023) showed how a single fact can be edited by applying rank one model editing to the parameters of an MLP layer. Our work is based on Meng et al.



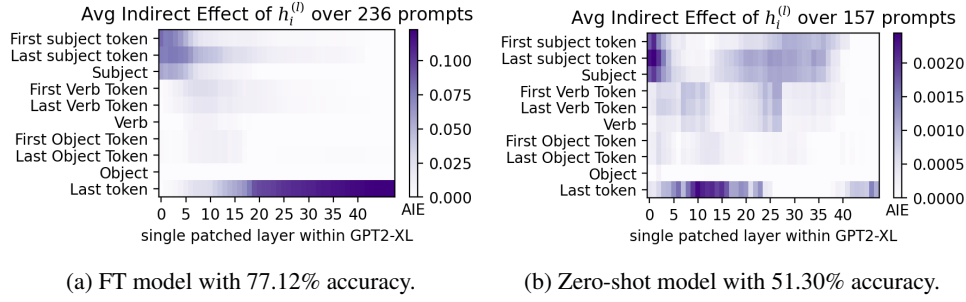


Figure 5: Zero-shot vs. FT causal tracing results for GPT-2 XL on PEP3k Inference Set #1.

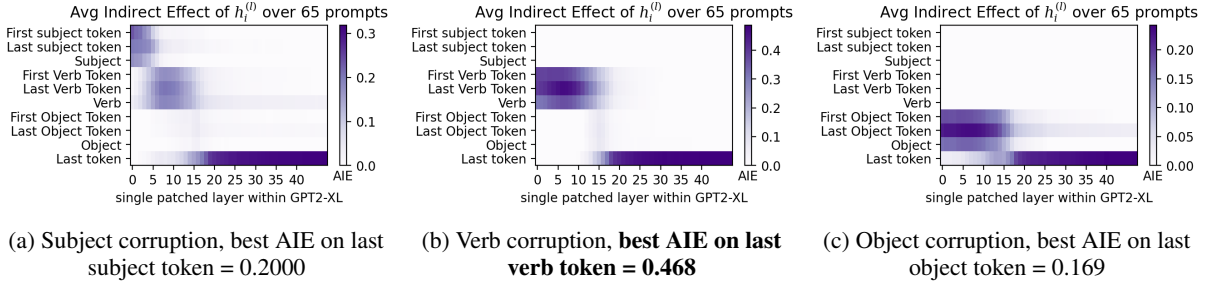


Figure 6: Causal tracing results for verb token edited GPT-2 XL on PEP3k Inference Set #1.

(2022) which extended this approach to thousands of edits by altering weights of a range of MLP layers (discussed in Section 3).

## 7 Discussion and Conclusion

In this research work, we establish that the adapted MEMIT<sub>CSK</sub> algorithm for editing Language Models scales beyond encyclopedic facts and works well for commonsense knowledge. Our experiments show strong localization of common sense knowledge in the early layers of the transformer models.

We discover strong causal tracing patterns for the fine-tuned base models; our analysis shows that systematically selecting layers to edit based on maximum Indirect Effect and moving average Indirect Effect (based on the causal and severed causal tracings), is effective. We conduct experiments on varying numbers of layers and find consistent results.

Editing transformer models with commonsense knowledge is shown to outperform FT with a high overall F1 score on Inference Sets #1 and #2. It also exhibits high configuration generalization when compared to FT and can be effectively applied to new datasets with similar features. Experiments with the GPT-2 Large and GPT-2 XL models display that editing performance improves as the model scale increases.

Through our research, we discover that editing is an effective method for improving commonsense quality and semantic understanding. In comparison to fine-tuning, all editing methods consistently have a higher overall score. We also observe performance patterns in different categories, including unaffected and affected neighborhoods, affected paraphrase, and affected reasoning. Our findings suggest that editing is a more widespread and generalized approach. However, there is a trade-off in performance across these categories, so it’s important to consider all scores when evaluating the method. Overall, our metrics have provided valuable insight into the effectiveness of various editing methods. MEMIT<sub>CSK</sub> experiments suggest a powerful future solution that could combine model editing and human involvement to effectively address and fix subtle commonsense errors in text.

## References

- James A. Anderson. 1972. [A simple neural network generating an interactive memory](#). *Mathematical Biosciences*, 14(3):197–220.
- Emily M. Bender and Alexander Koller. 2020. [Climbing towards NLU: On meaning, form, and understanding in the age of data](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5185–5198, Online. Association for Computational Linguistics.

- Prajwal Bhargava and Vincent Ng. 2022. [Commonsense knowledge reasoning and generation with pre-trained language models: A survey](#).
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#).
- Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2022. [Knowledge neurons in pretrained transformers](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8493–8502, Dublin, Ireland. Association for Computational Linguistics.
- Guy Dar, Mor Geva, Ankit Gupta, and Jonathan Berant. 2022. [Analyzing transformers in embedding space](#).
- Ashwin Devaraj, William Sheffield, Byron Wallace, and Junyi Jessy Li. 2022. [Evaluating factuality in text simplification](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7331–7345, Dublin, Ireland. Association for Computational Linguistics.
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. 2021. A mathematical framework for transformer circuits. *Transformer Circuits Thread*. <https://transformer-circuits.pub/2021/framework/index.html>.
- Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2021. [Transformer feed-forward layers are key-value memories](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5484–5495, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Ari Holtzman, Peter West, Vered Shwartz, Yejin Choi, and Luke Zettlemoyer. 2021. [Surface form competition: Why the highest probability answer isn’t always right](#). *CoRR*, abs/2104.08315.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2021. [Retrieval-augmented generation for knowledge-intensive nlp tasks](#).
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized bert pretraining approach](#).
- Aman Madaan, Niket Tandon, Peter Clark, and Yiming Yang. 2022. [Memory-assisted prompt editing to improve GPT-3 after deployment](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2833–2861, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Gary Marcus. 2021. [Experiments testing gpt-3’s ability at commonsense reasoning: results](#).
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2023. [Locating and editing factual associations in gpt](#).
- Kevin Meng, Arnab Sen Sharma, Alex Andonian, Yonatan Belinkov, and David Bau. 2022. [Mass-editing memory in a transformer](#).
- Ian Porada, Kaheer Suleman, Adam Trischler, and Jackie Chi Kit Cheung. 2021. [Modeling event plausibility with consistent conceptual abstraction](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1732–1743, Online. Association for Computational Linguistics.
- Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. 2021. [Retrieval augmentation reduces hallucination in conversation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3784–3803, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. [CommonsenseQA: A question answering challenge targeting commonsense knowledge](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.
- Derek Tam, Anisha Mascarenhas, Shiyue Zhang, Sarah Kwan, Mohit Bansal, and Colin Raffel. 2022. [Evaluating the factual consistency of large language models through summarization](#).
- Cunxiang Wang, Shuailong Liang, Yue Zhang, Xiaonan Li, and Tian Gao. 2019. [Does it make sense? and why? a pilot study for sense making and explanation](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4020–4026, Florence, Italy. Association for Computational Linguistics.

Rongxiang Weng, Heng Yu, Xiangpeng Wei, and Weihua Luo. 2020. [Towards enhancing faithfulness for neural machine translation](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2675–2684, Online. Association for Computational Linguistics.

Yunzhi Yao, Shaohan Huang, Li Dong, Furu Wei, Huajun Chen, and Ningyu Zhang. 2022. [Kformer: Knowledge injection in transformer feed-forward layers](#).

## A Appendix

### A.1 Causal Mediation Analysis

The Figures 7, 8, 9 shows the causal graphs for the 20Q dataset and the GPT2-Large fine-tuned model, 20Q dataset and the GPT2-XL fine-tuned model, and the PEP3k dataset and the GPT2-Large fine-tuned model.

For each of the editing locations, we see that the Last Token has higher AIE towards the later layers of the model which is consistent with the results of MEMIT on encyclopedic knowledge. Focusing on the subject, verb, and object tokens, we see that all of them show high AIE in the early layers of the model and that the effect on the corresponding last corrupted token is more pronounced than that of the first corrupted token. This shows that selecting the last subject/verb/object token and the early layers of the model should give good results for the editing method. These patterns are consistent with all the models and datasets and they match with the layer selection for the best editing hyper-parameters and consistently give good performance.

### A.2 Hyper-parameters

#### Fine-tuning

The Table 7 shows the best hyper-parameters for the fine-tuning method. The method was very sensitive to small changes in learning rate while the other parameters worked well over a long range of values. Note that we use early stopping and restore the weights to the best performing model based on the overall F1 score.

#### Editing

The Table 8 shows the hyper-parameters for the editing method. The method was slightly sensitive to the learning rate and very sensitive to the edit token. It worked well over a range of the KL factor parameter, the number of layers to edit, and layer selection.

Dataset	Model	Learning Rate	Batch Size	Epochs
20q	GPT2-Large	0.000003451	8	10
	GPT2-XL	0.000001589	32	10
PEP3k	GPT2-Large	0.00000474	32	10
	GPT2-XL	0.000001313	8	10

Table 7: Best Hyper-parameters for the fine-tuning method on Inference Set #1.

Dataset	Model	Edit Token	Layers	Learning Rate	KL factor
20q	GPT2-Large	Last Subject	3,4,5	0.7868	0.625
		Last Verb	2,3,4,5,6	0.09393	0.625
		Last Object	1,2,3	0.6276	0.625
	GPT2-XL	Last Subject	2,3,4,5,6	0.04108	0.625
		Last Verb	3,4,5,6,7	0.01936	0.625
		Last Object	1,2,3	0.02689	0.625
PEP3k	GPT2-Large	Last Subject	4,5,6,7,8	0.32	0.625
		Last Verb	4,5,6,7,8	0.682	0.625
		Last Object	1,2,3,4,5	0.433	0.625
	GPT2-XL	Last Subject	1,2,3,4,5	0.1253	0.625
		Last Verb	6,7,8	0.08719	0.625
		Last Object	3,4,5	0.04107	0.625

Table 8: Best Hyper-parameters for the editing method on Inference Set #1.

### A.3 MEMIT-CSK-PROBE Dataset construction

The GPT-3 prompts used for creating the MEMIT-CSK-PROBE has been described in the following figures,

- Affected Paraphrase: Figure 10
- Affected Reasoning: Figure 11
- Affected Neighborhood: Figure 12
- Unaffected Neighborhood: Figure 13

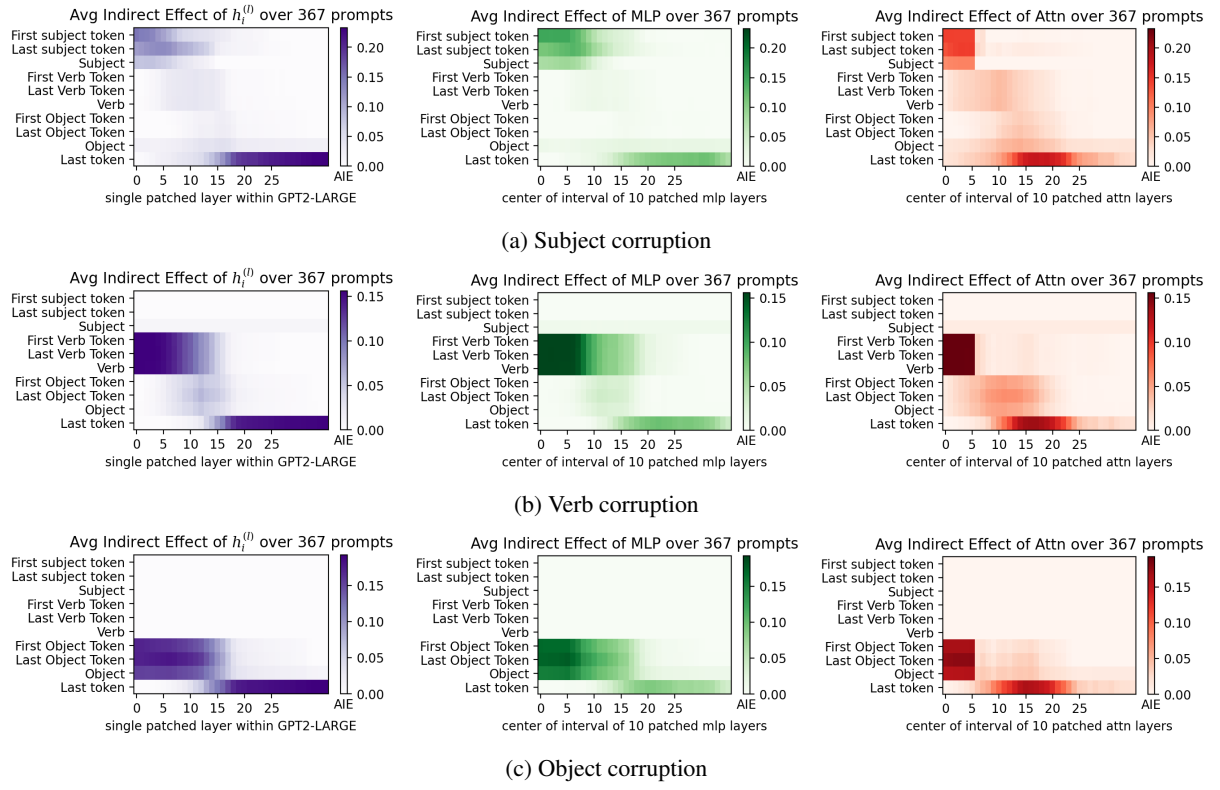


Figure 7: Causal tracing results for GPT-2 Large on 20Q Inference Set #1.

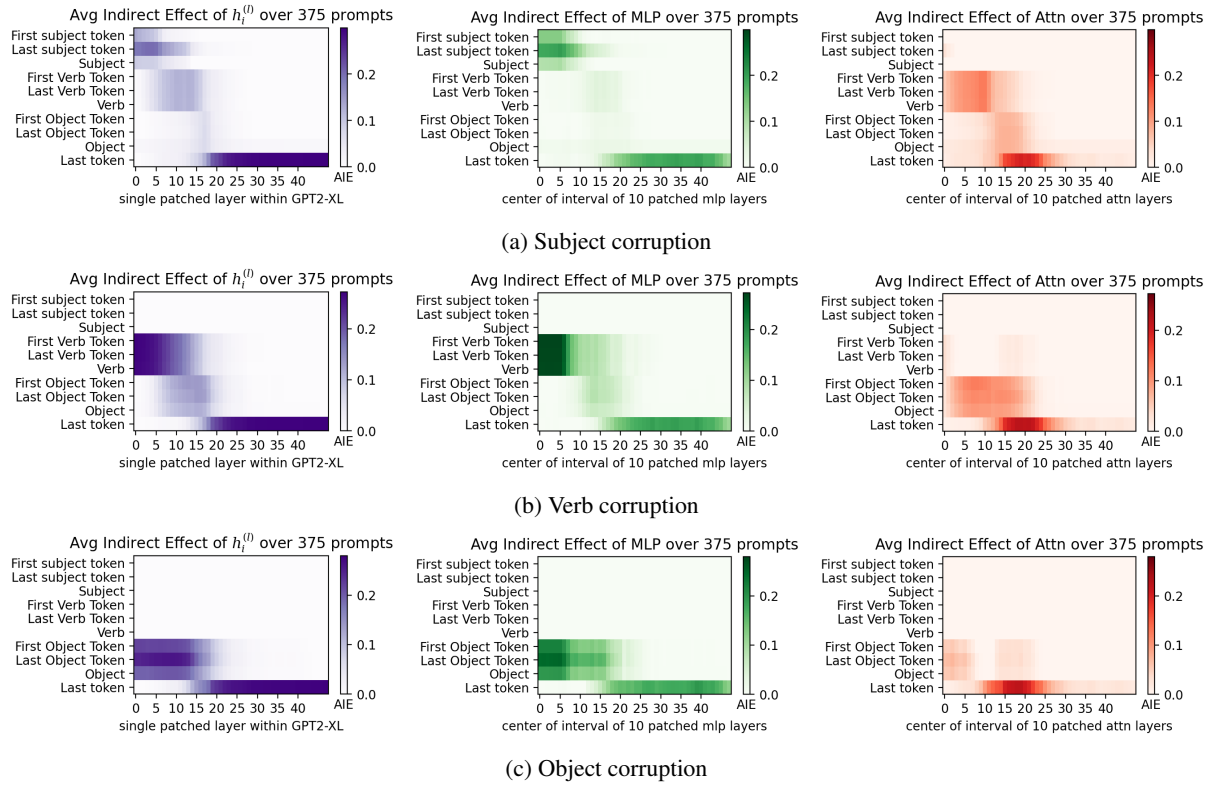


Figure 8: Causal tracing results for GPT-2 XL on 20Q Inference Set #1.



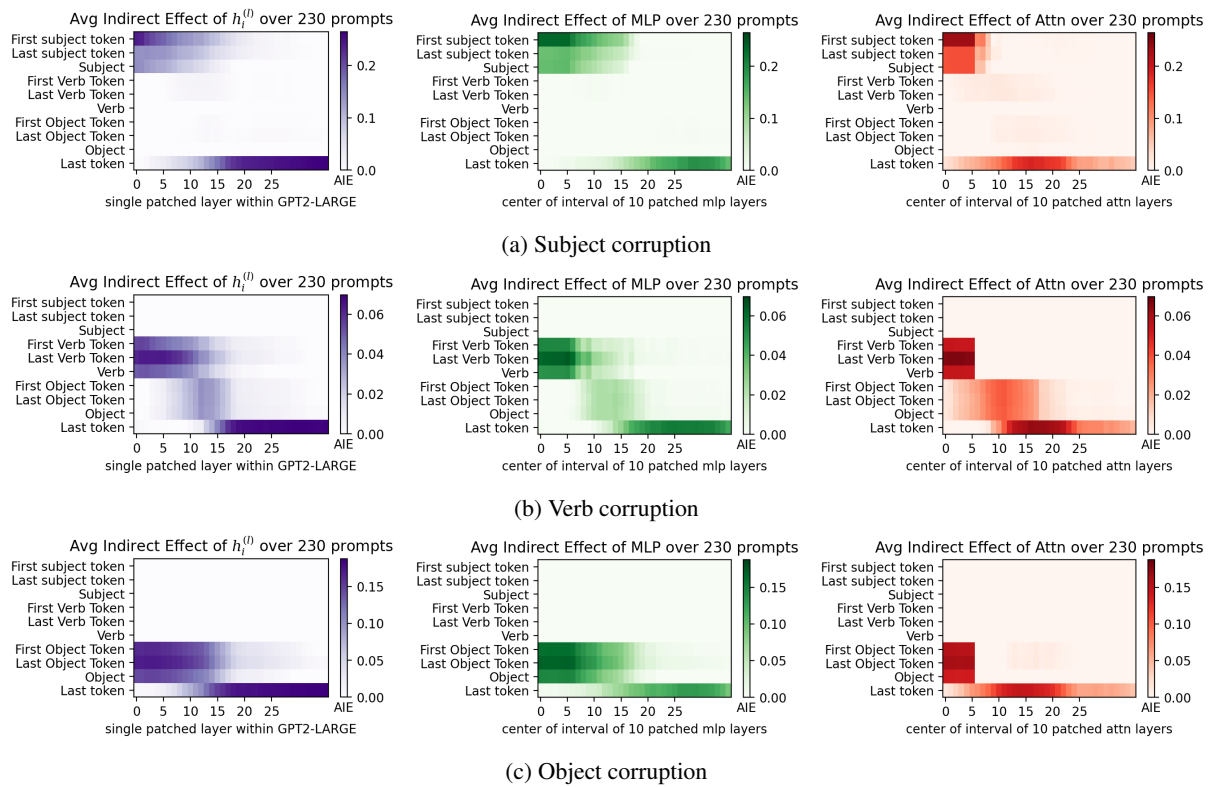


Figure 9: Causal tracing results for GPT-2 Large on PEP3k Inference Set #1.

Provide 5 paraphrases of: Furnishings make noise

1. Furniture can be noisy.
2. Furniture can create sound.
3. Furniture can produce noise.
4. Furniture can be a source of sound.
5. Furniture can be a source of noise.

Figure 10: Prompt to generate affected paraphrase for “Furnishings make noise (false)”

Furnishings do not make noise. Explain this with a 2-step reasoning chain of very short, simple, connected sentences:

1. Furnishings are inanimate objects.
2. Inanimate objects cannot make noise.

Figure 11: Prompt to generate affected reasoning neighborhood for “Furnishings make noise (false)”

Given the text: Furnishings make noise  
subject token: Furnishings  
object token: noise

Q1. In the text, replace just the subject token with a different word. The replaced text should be a valid sentence. The replaced token can be a hyponym or similar word of the original subject token. Write up to 5 such variants.

Q2. In the text, replace just the verb token with a different word. The replaced text should be a valid sentence. The replaced token can be a verb that follows or precedes the original verb token. Write up to 5 such variants.

Q3. In the text, replace just the object token with a different word. The replaced text should be a valid sentence. The replaced token can be a hyponym or similar word of the original object token. Write up to 5 such variants.

-----  
Q1.

1. Appurtenances make noise
2. Fixtures make noise
3. Accoutrements make noise
4. Decorations make noise
5. Adornments make noise

Q2.

1. Furnishings create noise
2. Furnishings emit noise
3. Furnishings generate noise
4. Furnishings produce noise
5. Furnishings yield noise

Q3.

1. Furnishings make sound
2. Furnishings make clamor
3. Furnishings make din
4. Furnishings make racket
5. Furnishings make uproar

Figure 12: Prompt to generate affected neighborhood for “Furnishings make noise (false)”

Given:

text: Furnishings make noise  
subject token: Furnishings  
object token: noise

Q1. Replace the subject token with a completely unrelated word and make a new text. Make 5 such replacements.

Q2. Replace the object token with a completely unrelated word and make a new text. Make 5 such replacements.

-----  
1. Replacing the subject token:

- a. Cars make noise
- b. Animals make noise
- c. People make noise
- d. Plants make noise
- e. Computers make noise

2. Replacing the object token:

- a. Furnishings make music
- b. Furnishings make laughter
- c. Furnishings make light
- d. Furnishings make heat
- e. Furnishings make color

Figure 13: Prompt to generate unaffected neighborhood for “Furnishings make noise (false)”