

WilKE: Wise-Layer Knowledge Editor for Lifelong Knowledge Editing

Anonymous ACL submission

Abstract

Knowledge editing aims to rectify inaccuracies in large language models (LLMs) without costly retraining for outdated or erroneous knowledge. However, current knowledge editing methods primarily focus on single editing, failing to meet the requirements for lifelong editing¹. This study reveals a performance degradation encountered by knowledge editing in lifelong editing, characterized by toxicity buildup and toxicity flash, with the primary cause identified as pattern mismatch. We introduce a knowledge editing approach named WilKE, which selects editing layer based on the pattern matching degree of editing knowledge across different layers. Experimental results demonstrate that, in lifelong editing, WilKE exhibits an average improvement of 46.2% and 67.8% on editing GPT2-XL and GPT-J relative to state-of-the-art knowledge editing methods.

1 Introduction

Large language models (LLMs) encode a wealth of world knowledge through pretraining on massive corpus (Radford et al., 2019; Brown et al., 2020; Achiam et al., 2023). However, outdated or erroneous knowledge may persist, and retraining these models with updated corpus incurs prohibitively high costs. To address this challenge, numerous studies have introduced knowledge editing (De Cao et al., 2021; Mitchell et al., 2021; Meng et al., 2022a; Meng et al., 2022b) as a solution, which involves updating the internal parameters of language models to edit specific knowledge.

Current knowledge editing methods are evaluated in single editing by default, which updates a single knowledge (x_e, y_o) to (x_e, y_e) on initial model f_θ for each test point, as shown in Figure 1(a). However, knowledge should be updated continuously in fact, making single editing insuff-

¹In this paper, lifelong editing is synonymous with lifelong knowledge editing.

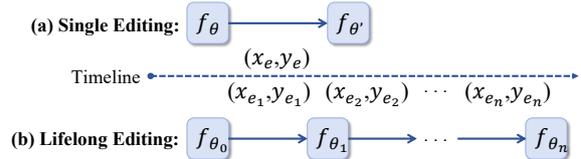


Figure 1: Single editing versus lifelong edit. (a) Single editing only involves making an edit. (b) Life-long editing involves continuous edits and monitoring performance.

cient to meet the demands. Therefore, we focus on lifelong editing, which updates a knowledge (x_{e_i}, y_{o_i}) to (x_{e_i}, y_{e_i}) on model $f_{\theta_{i-1}}$ that doesn't have to be initial model f_{θ_0} , as shown in Figure 1(b).

In this paper, we conduct an analysis of state-of-the-art knowledge editing methods such as ROME (Meng et al., 2022a) and MEMIT (Meng et al., 2022b), revealing a severe performance degradation when applied in lifelong editing. Investigating this issue further, our experiments indicate that these methods suffer from **toxicity buildup** and **toxicity flash** during ongoing editing. As shown in Figure 3(a), 4(a), the combined effects of both phenomena result in a "step-like" shape. On the one hand, the toxicity buildup signifies that one edit induces minor changes in irrelevant parameters, gradually leading to model's failure. On the other hand, the toxicity flash suggests that one edit modifies model's parameters abnormally, resulting in severe overfitting to specific edit, which is not reported in previous research. It's worth noting that due to overfitting, such failures are undetectable in single editing, and achieve respectable scores.

We analyze the primary reasons for these two phenomena, attributing them to pattern mismatch, as illustrated in Figure 2. Specifically, different layers of language model may detect different patterns, which is called key in key-value memories (Sukhbaatar et al., 2015; Sukhbaatar et al., 2019;

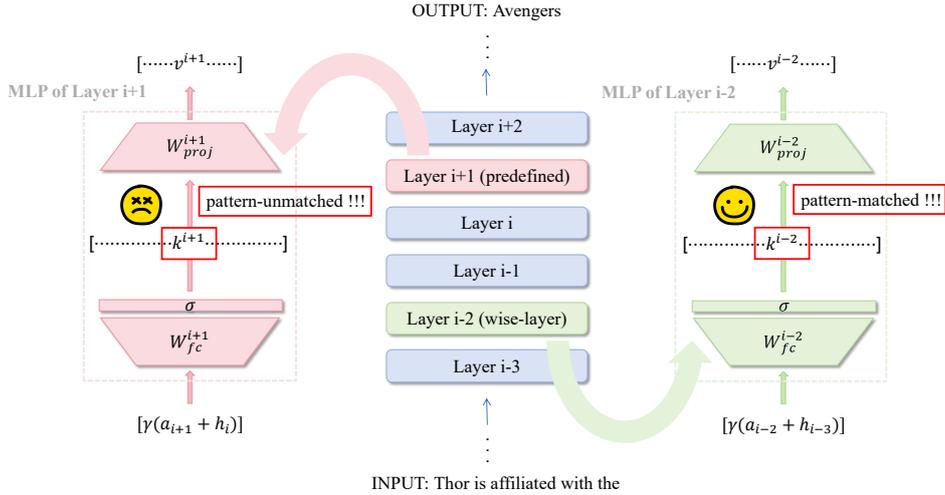


Figure 2: Illustration of our work. Predefined editing layers may not necessarily accommodate all editing knowledge effectively. Therefore, it would be wiser to select different editing layers for different editing knowledge.

Geva et al., 2020), thus extracting relevant information according to patterns and updating the hidden states. In other words, different knowledge may be stored in different layers, as illustrated in Section 4.3. However, ROME and MEMIT perform knowledge editing at predefined layers, which primarily lead to toxicity buildup and toxicity flash. To address this issue, we propose **WilKE** (Wise-Layer Knowledge Editor), which eliminates the need for predefined editing layer. Instead, WilKE selects editing layer based on the degree of pattern matching for different editing knowledge across various layers. Experimental results demonstrate that WilKE exhibits state-of-the-art comprehensive performance when editing GPT2-XL (1.5B) (Radford et al., 2019) and GPT-J (6B) (Wang and Komatsuzaki, 2021). Specifically, in lifelong editing scenarios, under identical experimental conditions of conducting 1024 edits, WilKE demonstrates an average improvement of 46.2% and 67.8% in comprehensive performance relative to state-of-the-art methods when editing GPT2-XL and GPT-J, respectively.

In summary, our primary contributions are as follows:

- We investigate the failure of ROME and MEMIT in lifelong editing, revealing toxicity buildup and toxicity flash during ongoing editing. The underlying primary cause of these phenomena is found to be pattern mismatch.

- To address this issue, we introduce WilKE. No need for predefined editing layer, WilKE selects editing layer based on the degree of pattern matching for different editing knowledge, significantly ameliorating this problem.
- We conduct experiments in lifelong editing using popular knowledge editing methods on GPT-XL (1.5B) and GPT-J (6B), highlighting the superiority of WilKE over prevalent knowledge editing methods. The source code will be made available soon.

2 Related Work

Generally, knowledge editing aims to edit the knowledge of a language model so that its outputs reflect the revised state when presented with relevant inputs (De Cao et al., 2021). Yao et al. (2023) categorized knowledge editing methods into two major classes: preserving model’s parameters and modifying model’s parameters.

Methods for preserving model’s parameters include memory-based methods and additional parameters’ methods. Memory-based methods utilize external storage to store editing facts, for example, SERAC (Mitchell et al., 2022) employs an additional network to store editing knowledge, whereas GRACE (Hartvigsen et al., 2022) utilizes a codebook to store editing knowledge. Additional parameters’ methods employ extra neurons to store editing facts, for instance, Huang et al. (2023) and Dong et al. (2022) adding extra neurons in MLP to memorize additional facts.

Since our target is to edit knowledge by updating the internal parameters of language models, this paper focuses on methods that modify model’s parameters. Currently, methods for modifying model’s parameters can be further divided into two categories: meta-learning and locate-and-edit.

Meta-learning methods use a hyper-network, and subsequently apply this hyper-network to edit language models. For instance, De Cao et al. (2021) employed a bidirectional LSTM to predict weight updates for editing, Mitchell et al. (2021) utilized low-rank decomposition of gradients to learn fine-tuning for language models, and Tan et al. (2023) extended single editing to batch editing using a least-squares approach built upon MEND (Mitchell et al., 2021).

Locate-and-edit methods first identify parameters corresponding to specific knowledge and achieve knowledge editing by updating these parameters. For example, Dai et al. (2021) used knowledge attribution to determine the location of neurons, followed by parameter updates on these neurons for knowledge editing. Meng et al. (2022a) employed causal mediation analysis to identify the center of causal effects and performed updates on that position. Meng et al. (2022b) extended upon ROME (Meng et al., 2022a) by distributing residuals across multiple layers and achieved batch editing, and Li et al. (2023a) achieved more precise residual allocation.

However, existing knowledge editing methods are limited in single editing, overlooking their scalability. This oversight could lead to an overly optimistic view of knowledge editing methods and potential misuse when the methods are not yet mature. Therefore, for knowledge editing to be practically applicable, it is crucial to address the more challenging scenario, lifelong editing. Consequently, research into lifelong editing is imperative.

3 Preliminary

The language model $f_\theta \in \mathcal{F}$ can be defined as a function $f_\theta : \mathcal{X} \mapsto \mathcal{Y}$, mapping input $\mathbf{x} \in \mathcal{X}$ to its prediction $\mathbf{y} \in \mathcal{Y}$. For an editing example $(\mathbf{x}_e, \mathbf{y}_e)$, where $f_\theta(\mathbf{x}_e) \neq \mathbf{y}_e$, the goal of the knowledge editing (KE) is to edit the parameters $\theta \in \Theta$ of the model f_θ to obtain an edited model $f_{\theta'}$, such that $f_{\theta'}(\mathbf{x}_e) = \mathbf{y}_e$.

$$KE : \mathcal{F} \times \mathcal{X} \times \mathcal{Y} \mapsto \mathcal{F} \quad (1)$$

In lifelong editing, such a process continues it-

eratively. In other words, for an initial language model f_{θ_0} , there exists a potential sequence to be edited $(\mathbf{x}_{e_i}, \mathbf{y}_{e_i})_{i=1}^n$, and the model undergoes continuous editing:

$$f_{\theta_i} = KE(f_{\theta_{i-1}}, \mathbf{x}_{e_i}, \mathbf{y}_{e_i}) \quad (2)$$

In lifelong editing, the edited model should satisfy the following properties.

Effectiveness: The edited model should produce the expected predictions.

$$f_{\theta_i}(\mathbf{x}_{e_i}) = \mathbf{y}_{e_i} \quad (3)$$

Generality: The edited model should remain consistent on its edited data equivalent input set $\mathcal{E}(\mathbf{x}_{e_i})$.

$$f_{\theta_i}(\mathbf{x}_j) = \mathbf{y}_{e_i}, \forall \mathbf{x}_j \in \mathcal{E}(\mathbf{x}_{e_i}) \quad (4)$$

Locality: The edited model should maintain the original output on data unrelated to the editing, denoted as $\mathcal{I}(\mathbf{x}_{e_i})$.

$$f_{\theta_i}(\mathbf{x}_j) = \mathbf{y}_{\mathbf{x}_j}, \forall \mathbf{x}_j \in \mathcal{I}(\mathbf{x}_{e_i}) \quad (5)$$

Retention: The edited model should preserve the editing results based on the previously completed edits.

$$f_{\theta_i}(\mathbf{x}_{e_j}) = \mathbf{y}'_{e_j}, \forall 1 \leq j < i \quad (6)$$

Here is \mathbf{y}'_{e_j} rather than \mathbf{y}_{e_j} because we consistently adhere to a principle: the later the edit, the higher the priority. Later edits take precedence over earlier ones and potentially engage in complex interactions with the original knowledge to update it. For further explanations and details, please refer to Appendix F.

4 Toxicity in Lifelong Editing

ROME (Meng et al., 2022a) and MEMIT (Meng et al., 2022b) are currently the state-of-the-art knowledge editing methods. As MEMIT is based on ROME, implementing residual distribution across multiple layers, our analysis in the main text focuses primarily on ROME. The analysis of MEMIT is provided in Appendix C. In this section, we systematically investigate the reasons for the failure of ROME in lifelong editing.

4.1 Toxicity

As editing progresses, the performance of the language model continuously deteriorates (Yao et al., 2023), indicating that ongoing editing seems to introduce certain side effects. In this section, we

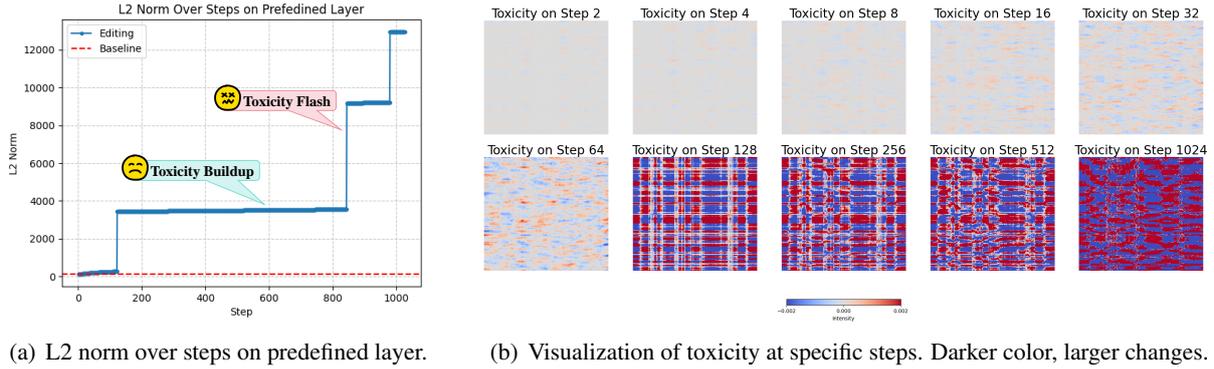


Figure 3: The toxicity on GPT2-XL with editing steps.

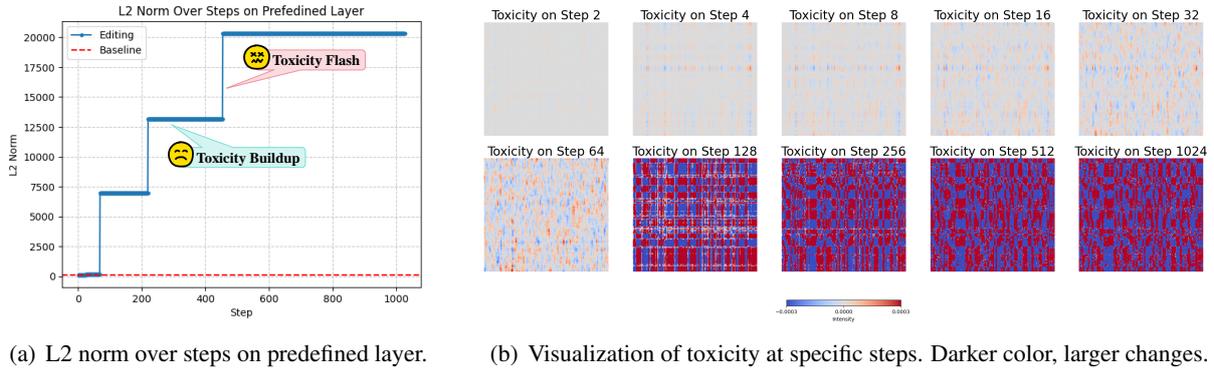


Figure 4: The toxicity on GPT-J with editing steps.

refer to these side effects as "**toxicity**" and utilize rollback editing (Li et al., 2023b) to define toxicity:

$$\begin{aligned} \text{Toxicity} &= \theta^* - \theta \\ \text{s.t. } f_{\theta^*} &= KE(KE(f_{\theta}, \mathbf{x}_e, \mathbf{y}_e), \mathbf{x}_e, \mathbf{y}_o), \end{aligned} \quad (7)$$

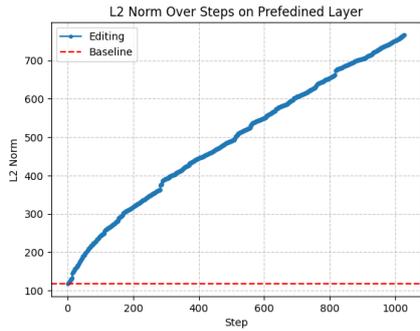
where $f_{\theta}(\mathbf{x}_e) = \mathbf{y}_o$. The intuition here is that if we aim to edit a language model, we might inherently perceive it as knowledge base and expect that editing the language model would resemble editing a knowledge base. Therefore, after rollback editing, we expect the language model to return to its initial state. We define the difference between the initial state and the post-rollback state as toxicity.

To better simulate real-world knowledge editing scenarios, we first filter data points corresponding to known knowledge in CounterFact dataset (Meng et al., 2022a) for both GPT2-XL (Radford et al., 2019) and GPT-J (Wang and Komatsuzaki, 2021). Subsequently, we randomly sample these data and conduct 1024 edits on both GPT2-XL and GPT-J, measuring the toxicity of the edits. As depicted in Figure 3(a), 4(a), the red dashed line represents the L2 norm of the original parameters on the pre-

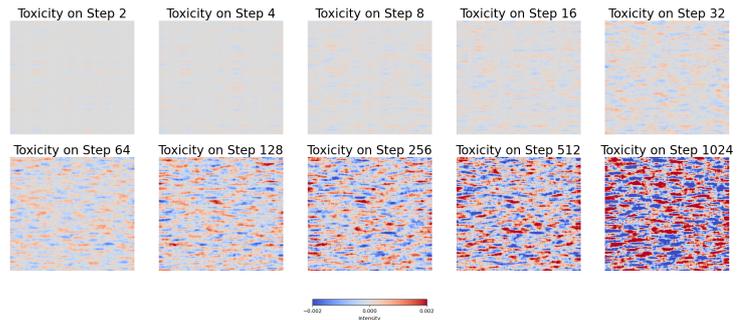
defined editing layer, while the blue solid line represents the L2 norm of the actual parameters on the predefined editing layer as editing progresses. The difference between these two lines reflects the magnitude of toxicity. Figure 3(b), 4(b) visualizes the accumulated toxicity at specific steps, such as steps 2, 4, 8, 16, 32, 64, 128, 256, 512, and 1024, corresponding to specific positions in Figure 3(a), 4(a), to illustrate the toxicity status at these steps.

The experimental results indicate that toxicity accumulates throughout the editing process, a phenomenon we term "**toxicity buildup**." Additionally, "spikes" in toxicity are observed at certain data points, which we term "**toxicity flash**." Consequently, the overall measurement exhibits a staircase shape. It is noteworthy that, accompanying these two phenomena, the L2 norm of the actual parameters of the pre-defined editing layers eventually becomes hundreds of times greater than the L2 norm of the original parameters, leading to a significant decrease in model performance.

Further investigation reveals that the data points causing toxicity flash exhibit the same "spike" phenomenon even in single editing. This suggests that

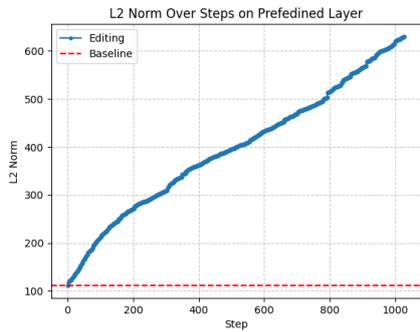


(a) L2 norm over steps on predefined layer.

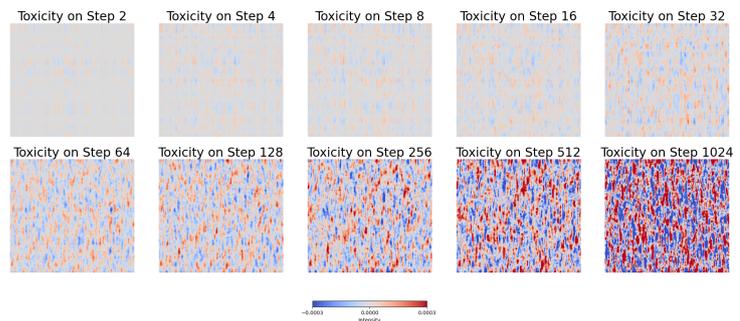


(b) Visualization of toxicity at specific steps. Darker color, larger changes.

Figure 5: Toxicity buildup on GPT2-XL with editing steps.



(a) L2 norm over steps on predefined layer.



(b) Visualization of toxicity at specific steps. Darker color, larger changes.

Figure 6: Toxicity buildup on GPT-J with editing steps.

268 the occurrence of these spikes is not exclusive to
 269 lifelong editing. In the lifelong editing scenario, we
 270 uncover issues that were not previously reported
 271 in single-editing scenarios. We will delve deeper
 272 into these two phenomena in the subsequent sub-
 273 sections.

274 4.2 Toxicity Buildup

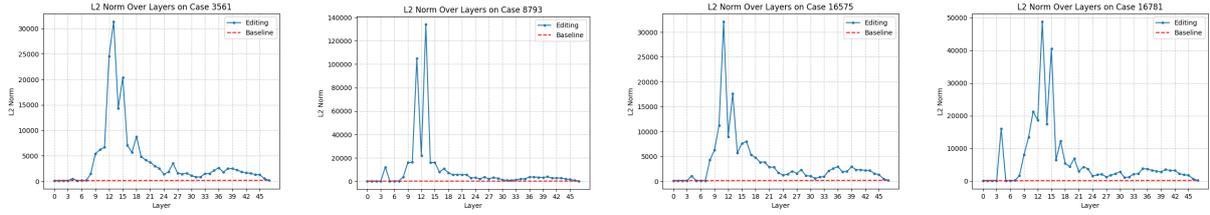
275 Based on the data corresponding to known knowl-
 276 edge, we filter out the data causing toxicity flash
 277 during editing on GPT2-XL and GPT-J. Details of
 278 the filtering process and the results are described in
 279 Appendix B. We removed these data points causing
 280 toxicity flash and conducted the same experiment
 281 above again. As shown in Figure 5(a), 6(a), the red
 282 dashed line represents the L2 norm of the original
 283 parameters of the predefined editing layer, and the
 284 blue solid line represents the L2 norm of the actual
 285 parameters of the predefined editing layer as editing
 286 progresses. The difference between these lines re-
 287 flects the magnitude of toxicity. Figure 5(b),6(b) vi-
 288 sualizes the accumulated toxicity at specific steps.

289 From the experimental results, it is evident
 290 that after filtering out data causing toxicity flash,

291 the magnitude of toxicity significantly decreases.
 292 Moreover, as shown in Figure 5(b), 6(b), the pro-
 293 cess of toxicity buildup becomes more uniform and
 294 gradual than Figure 3(b), 4(b). However, toxicity
 295 continues to steadily accumulate as editing pro-
 296 gresses, leading to a steady decline in model per-
 297 formance. Additionally, this observation may suggest
 298 that editing results in the disruption of superposi-
 299 tion (Elhage et al., 2022b; Henighan et al., 2023)
 300 and polysemantic neurons (Elhage et al., 2022a)
 301 in the original model, which could be important
 302 factors contributing to the continuous decline in
 303 models’ performance during the editing process.

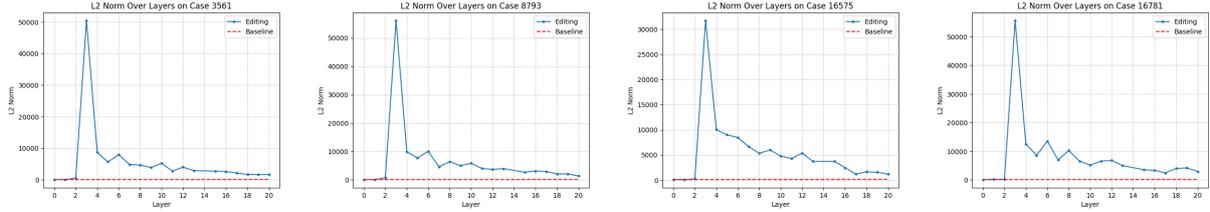
304 4.3 Toxicity Flash

305 Subsequently, we focus on the data causing toxic-
 306 ity flash during editing. It is worth noting that the
 307 majority of the data causing toxicity flash when
 308 editing GPT2-XL and GPT-J overlap. We then con-
 309 duct single editing experiments on these data for
 310 GPT2-XL and GPT-J. Here, we perform editing
 311 experiments on different layers in language mod-
 312 els, plotting the L2 norms of the parameters before
 313 and after editing. The experimental results are il-



(a) L2 norm over layers on case 3561. (b) L2 norm over layers on case 8793. (c) L2 norm over layers on case 16575. (d) L2 norm over layers on case 16781.

Figure 7: Toxicity flash on GPT2-XL among editing layers.



(a) L2 norm over layers on case 3561. (b) L2 norm over layers on case 8793. (c) L2 norm over layers on case 16575. (d) L2 norm over layers on case 16781.

Figure 8: Toxicity flash on GPT-J among editing layers.

314 lustrated in Figures 7, 8, where the red dashed line
 315 represents the L2 norm of the original parameters
 316 of different layers in language models, and the
 317 blue solid line represents the L2 norm of the actual
 318 parameters after single editing on different layers.
 319 A larger gap between these lines indicates greater
 320 toxicity caused by editing on the corresponding
 321 layer. To compare the toxicity flash phenomena in
 322 different models, we present the overlapping data
 323 causing toxicity flash on both GPT2-XL and GPT-
 324 J. Further experiments on toxicity flash data for
 325 GPT2-XL and GPT-J, as well as comparisons with
 326 experiments on other regular data, can be found in
 327 Appendix E.

328 ROME’s predefined editing layers on GPT2-XL
 329 and GPT-J are 17 and 5, respectively, where Meng
 330 et al. (2022a) described these as the center of causal
 331 effects, which has been further utilized in MEMIT
 332 (Meng et al., 2022b). However, as observed from
 333 the experimental results, editing these layers leads
 334 to toxicity flash, indicating that predefined editing
 335 layer is the direct cause of toxicity flash. From
 336 Figures 7, 8, it can be inferred that for these data
 337 points, we should edit the layer that does not align
 338 with the predefined editing layer. The results of
 339 causal mediation analysis on these data points also
 340 support this conclusion: in fact, these knowledge
 341 are extracted from the earlier layer’s MLP of the
 342 model. For detailed experimental results, please

refer to Appendix A.

After further investigation, the fundamental reason lies in pattern mismatch, where the patterns of editing knowledge cannot be detected in the predefined editing layer. Continuing editing leads to language model overfitting, resulting in toxicity flash. The formal definition of pattern mismatch and experimental evidence are provided in Appendix D.

5 Wise-Layer Knowledge Editor

In Section 4, we delved into the primary reason for failure in lifelong editing - pattern mismatch, which directly leads to toxicity flash and potentially more toxicity buildup. In light of this, we propose an editing method called WilKE. Unlike ROME and MEMIT, WilKE does not predefine editing layer; instead, it selects editing layer based on the degree of pattern match for different editing knowledge across various layers. We first describe where to edit in Section 5.1, followed by an explanation of how to edit in Section 5.2.

5.1 Where to Edit?

To implement knowledge editing, the initial step involves determining the locations where editing will take place.

Meng et al. (2022a) utilizes causal mediation analysis to identify the center of causal effects, MLP at specific layer, for storing factual knowl-

| Model | Editor | Score | Effectiveness | Generality | Locality | Retention | |
|----------|--------------------|------------------|------------------|------------------|------------------|------------------|------------------|
| | | S ↑ | ES ↑ | PS ↑ | NS ↑ | ERS ↑ | ORS ↑ |
| GPT-2 XL | KE | 0.0(0.0) | 0.1(0.0) | 0.1(0.0) | 0.0(0.0) | 0.0(0.0) | 0.0(0.0) |
| | KN | 0.0(0.0) | 0.1(0.0) | 1.0(0.0) | 0.0(0.0) | 0.0(0.0) | 0.0(0.0) |
| | MEND | 0.0(0.0) | 0.5(0.0) | 0.1(0.0) | 0.4(0.0) | 0.0(0.0) | 0.0(0.0) |
| | ROME | 9.3(2.4) | 15.8(4.4) | 8.8(2.7) | 6.8(1.4) | 12.2(2.6) | 7.9(2.6) |
| | MEMIT | 13.2(7.5) | 92.5(0.5) | 55.1(0.5) | 35.0(0.7) | 6.6(4.4) | 6.6(4.5) |
| | WilKE(Ours) | 19.3(5.6) | 70.7(9.2) | 51.0(7.0) | 12.7(0.9) | 16.1(6.4) | 13.1(5.7) |
| GPT-J | MEND | 3.7(3.0) | 3.7(2.6) | 1.7(1.0) | 4.7(3.3) | 9.8(8.0) | 9.2(6.5) |
| | ROME | 8.7(0.4) | 28.7(1.0) | 20.7(0.7) | 4.6(0.3) | 10.9(2.7) | 5.8(0.6) |
| | MEMIT | 0.0(0.0) | 32.1(3.6) | 23.8(2.2) | 9.3(1.1) | 0.0(0.0) | 0.0(0.0) |
| | WilKE(Ours) | 14.6(2.6) | 71.3(6.5) | 50.6(6.4) | 7.4(0.9) | 19.1(8.0) | 8.5(1.9) |

Table 1: Evaluation results (%) with 95% confidence intervals in parentheses.

edge. The MLP of the FFN is divided into two matrices, represented as follows:

$$FFN^l(\mathbf{x}) = \sigma(\mathbf{x} \cdot W_{fc}^l) \cdot W_{proj}^l, \quad (8)$$

where $W_{fc}^l \in \mathbb{R}^{d \times d_m}$ and $W_{proj}^l \in \mathbb{R}^{d_m \times d}$ are the parameter matrices of the l -th layer’s FFN, FFN^l , and d_m is the dimension of the intermediate hidden layer in the FFN. The symbol σ represents the activation function, and $\mathbf{x} \in \mathbb{R}^d$ is the input to the FFN. As described in the key-value memories (Sukhbaatar et al., 2015; Sukhbaatar et al., 2019; Geva et al., 2020), W_{fc}^l identifies patterns of the input \mathbf{x} to obtain the key vector k , and then the value vector v is retrieved from W_{proj}^l , as shown in Figure 2. Therefore, to achieve knowledge editing, we modify W_{proj}^l . After identifying the component that needs modification, we further determine the specific layer for modifying this component.

To find the editing layer l^* , our initial intuition is to identify the layer that produces the maximum activation strength for a specific knowledge, which is represented as $\arg\max_l \|\sigma(\mathbf{x} \cdot W_{fc}^l)\|_2$. However, in reality, the optimization of δ varies across different layers, as detailed in Appendix D.2. In other words, the importance of the hidden state outputted by different layers for editing specific knowledge is actually different. To comprehensively consider these two points, we define the editing layer (**wise-layer**) l^* as follows:

$$l^* = \arg\min_l \left\| \frac{\delta}{\|W_{proj}^l\|_2 \sigma(\mathbf{x} \cdot W_{fc}^l)} \right\|_2 \quad (9)$$

where the term $\|W_{proj}^l\|_2$ in the denominator can be regarded as normalization, allowing for comparison across layers. Ultimately, we determine that the target editing location is $W_{proj}^{l^*}$ of the layer l^* .

5.2 How to Edit?

After determining the target editing location, the next step involves determining how to carry out an edit.

Same as Meng et al. (2022a), we introduce a residual term $\delta \in \mathbb{R}^d$ to the output of the FFN in editing layer l^* , denoted as $FFN^{l^*}(\mathbf{x}) + \delta$. We optimize this residual term to align the model’s output with our expected results while not affecting irrelevant knowledge. For specific optimization details, we recommend interested readers to refer to (Meng et al., 2022a).

Subsequently, we allocate the optimized residual δ^* to $W_{proj}^{l^*}$ to accomplish knowledge editing:

$$W_{proj}^{l^*} \leftarrow \frac{FFN(\mathbf{x})^{l^*} + \delta^*}{\sigma(\mathbf{x} \cdot W_{fc}^{l^*})} \quad (10)$$

Afterwards, we have completed one editing. In summary, our approach starts from the perspective of pattern matching, attempting to identify the layer that is most suitable for editing the given knowledge across all layers, and then performs knowledge editing on that location.

6 Experiments

6.1 Experimental Setting

Model We utilize two widely employed autoregressive language models for knowledge editing: GPT-XL (1.5B) (Radford et al., 2019) and GPT-J (6B) (Wang and Komatsuzaki, 2021).

Baselines Regarding knowledge editing methods, we select the following approaches: KnowledgeEditor (KE) (De Cao et al., 2021) utilizes a bidirectional LSTM to predict weight updates for editing data points; KnowledgeNeuron (KN) (Dai et al., 2021) employs knowledge attribution to determine the positions of neurons, followed by parameter updates on these neurons to implement knowledge up-

438 dates; MEND (Mitchell et al., 2021) uses low-rank
 439 decomposition of gradients to learn fine-tuning
 440 of language models; ROME (Meng et al., 2022a)
 441 employs causal mediation analysis to identify the
 442 center of causal effects, followed by gradient de-
 443 scent parameter updates on the MLP at that layer;
 444 MEMIT (Meng et al., 2022b) extends upon ROME
 445 by distributing residuals across multiple layers.
 446 **Datasets, Metrics and Experiment Details** Due
 447 to space limitations, details of dataset, metrics, and
 448 experimental details are provided in Appendix F
 449 for reference.

450 6.2 Main Results

451 As shown in Table 1, we present the knowledge
 452 editing results after 1024 edits on GPT-XL and
 453 GPT-J. The results indicate that current knowledge
 454 editing methods perform poorly in lifelong editing,
 455 far from the optimistic results reported in single
 456 editing. However, these methods have been directly
 457 transferred and used in many other scenarios (Ma
 458 et al., 2023; Li et al., 2023a; Anonymous, 2024;
 459 Wang et al., 2023).

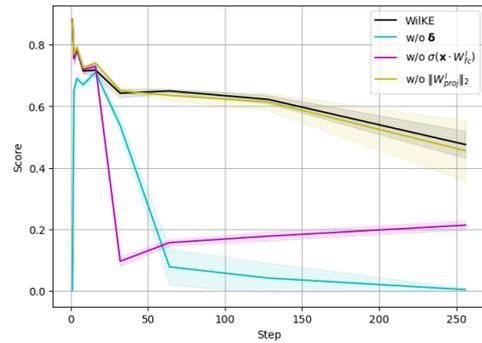
460 WilKE demonstrates the most advanced compre-
 461 hensive performance relative to the current knowl-
 462 edge editing methods. Specifically, under the same
 463 experimental conditions on GPT2-XL and GPT-J,
 464 WilKE achieves an average performance improve-
 465 ment of 46.2% and 67.8%, respectively, relative to
 466 the state-of-the-art methods.

467 To gain further insight, we have plotted the com-
 468 plete performance curves, and detailed results are
 469 presented in Appendix F.3.

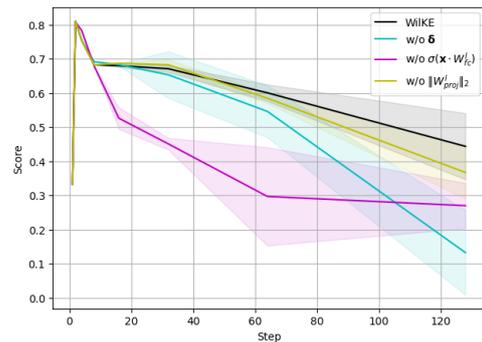
470 6.3 Ablation Study

471 Since the core of our method lies in selecting differ-
 472 ent editing layers based on various knowledge, as
 473 demonstrated in Equation 9 in Section 5.1, we com-
 474 prehensively consider three aspects: the optimiza-
 475 tion of δ for editing knowledge across different lay-
 476 ers, the activation of specific knowledge across dif-
 477 ferent layers $\sigma(\mathbf{x} \cdot W_{fc}^l)$, and the $\|W_{proj}^l\|_2$
 478 across different layers. Therefore, we sequentially ablate
 479 these three factors to demonstrate that consider-
 480 ing these three factors collectively leads to a better
 481 editing layer.

482 As depicted in Figure 9, it is evident that when
 483 individually ablated, both δ and $\sigma(\mathbf{x} \cdot W_{fc}^l)$ lead to
 484 a significant decrease in the performance of knowl-
 485 edge editing. Additionally, ablating $\|W_{proj}^l\|_2$ re-
 486 sults in a slight decrease in the performance of
 487 knowledge editing. However, when considering



(a) Score with editing steps on GPT2-XL.



(b) Score with editing steps on GPT-J.

Figure 9: The results of the ablation experiments.

488 these three factors collectively, superior experimen-
 489 tal results are obtained.

490 7 Conclusion

491 In this work, we focus on lifelong knowledge edit-
 492 ing, finding that current knowledge editing meth-
 493 ods suffer from severe performance degradation
 494 in lifelong editing. Our experimental results re-
 495 veal the **toxicity buildup** and **toxicity flash** that
 496 may occur during lifelong editing, leading to the
 497 deterioration of model’s performance. The primary
 498 reason for these problems lies in pattern unmatch.
 499 To address this issue, we propose a model editing
 500 method called **WilKE (Wise-Layer Knowledge**
 501 **Editor)**, which does not require predefined editing
 502 layer but selects editing layer based on specific
 503 editing knowledge. Experimental results demon-
 504 strate that in lifelong editing, WilKE achieves a
 505 significant improvement in overall performance
 506 compared to currently popular knowledge editing
 507 methods. In summary, our work is significant for
 508 improving knowledge editing methods and provide
 509 valuable insights for future work.

8 Limitation

Despite the promising performance of WilKE, our current studies still have limitations. Firstly, we select editing layer based on specific knowledge, yet knowledge may be distributed across multiple layers, leaving the question of how language models store knowledge is still under explored. Secondly, similar to previous knowledge editing research, we focus on factual knowledge assessment, which serves as a crucial entry point for our study on knowledge editing. Lastly, detecting match degree of specific knowledge across different layers of language models incurs a certain time cost, yet we believe this to be worthwhile in the initial stages of knowledge editing research.

9 Ethical Considerations

We have developed a method for knowledge editing in large language models under lifelong editing scenario, which may further expand our understanding of how language models store knowledge. However, the direct editing capability of large models also carries the potential for misuse, such as injecting malicious misinformation, biases, or other adversarial data into the model and deploying these edited models on open platforms. Given these concerns and our observations of speculative behavior, we emphasize the importance of sourcing large language models from authoritative origins and refraining from using them as sources of authoritative factual knowledge in critical environments.

References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. [arXiv preprint arXiv:2303.08774](#).

Anonymous. 2024. [Badedit: Backdooring large language models by model editing](#). In [The Twelfth International Conference on Learning Representations](#).

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

Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2021. Knowledge neurons in pretrained transformers. [arXiv preprint arXiv:2104.08696](#).

Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. Editing factual knowledge in language models. [arXiv preprint arXiv:2104.08164](#).

Qingxiu Dong, Damai Dai, Yifan Song, Jingjing Xu, Zhifang Sui, and Lei Li. 2022. Calibrating factual knowledge in pretrained language models. [arXiv preprint arXiv:2210.03329](#).

Yanai Elazar, Nora Kassner, Shauli Ravfogel, Abhilasha Ravichander, Eduard Hovy, Hinrich Schütze, and Yoav Goldberg. 2021. Measuring and improving consistency in pretrained language models. [Transactions of the Association for Computational Linguistics](#), 9:1012–1031.

Nelson Elhage, Tristan Hume, Catherine Olsson, Neel Nanda, Tom Henighan, Scott Johnston, Sheer ElShowk, Nicholas Joseph, Nova DasSarma, Ben Mann, et al. 2022a. Softmax linear units. [Transformer Circuits Thread](#).

Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, et al. 2022b. Toy models of superposition. [arXiv preprint arXiv:2209.10652](#).

Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2020. Transformer feed-forward layers are key-value memories. [arXiv preprint arXiv:2012.14913](#).

Thomas Hartvigsen, Swami Sankaranarayanan, Hamid Palangi, Yoon Kim, and Marzyeh Ghassemi. 2022. Aging with grace: Lifelong model editing with discrete key-value adaptors. [arXiv preprint arXiv:2211.11031](#).

Tom Henighan, Shan Carter, Tristan Hume, Nelson Elhage, Robert Lasenby, Stanislaw Fort, Nicholas Schiefer, and Christopher Olah. 2023. Superposition, memorization, and double descent. [Transformer Circuits Thread](#).

Zeyu Huang, Yikang Shen, Xiaofeng Zhang, Jie Zhou, Wenge Rong, and Zhang Xiong. 2023. Transformer-patcher: One mistake worth one neuron. [arXiv preprint arXiv:2301.09785](#).

Xiaopeng Li, Shasha Li, Shezheng Song, Jing Yang, Jun Ma, and Jie Yu. 2023a. Pmet: Precise model editing in a transformer. [arXiv preprint arXiv:2308.08742](#).

Zhoubo Li, Ningyu Zhang, Yunzhi Yao, Mengru Wang, Xi Chen, and Huajun Chen. 2023b. Unveiling the pitfalls of knowledge editing for large language models. [arXiv preprint arXiv:2310.02129](#).

Jun-Yu Ma, Jia-Chen Gu, Zhen-Hua Ling, Quan Liu, and Cong Liu. 2023. Untying the reversal curse via bidirectional language model editing. [arXiv preprint arXiv:2310.10322](#).

Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022a. Locating and editing factual associations in gpt. [Advances in Neural Information Processing Systems](#), 35:17359–17372.

615 Kevin Meng, Arnab Sen Sharma, Alex Andonian,
616 Yonatan Belinkov, and David Bau. 2022b. Mass-
617 editing memory in a transformer. [arXiv preprint](#)
618 [arXiv:2210.07229](#).

619 Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea
620 Finn, and Christopher D Manning. 2021. Fast model
621 editing at scale. [arXiv preprint arXiv:2110.11309](#).

622 Eric Mitchell, Charles Lin, Antoine Bosselut, Christo-
623 pher D Manning, and Chelsea Finn. 2022. Memory-
624 based model editing at scale. In [International](#)
625 [Conference on Machine Learning](#), pages 15817–
626 15831. PMLR.

627 Judea Pearl. 2022. Direct and indirect effects. In
628 [Probabilistic and causal inference: the works of](#)
629 [Judea Pearl](#), pages 373–392.

630 Alec Radford, Jeffrey Wu, Rewon Child, David Luan,
631 Dario Amodei, Ilya Sutskever, et al. 2019. Language
632 models are unsupervised multitask learners. [OpenAI](#)
633 [blog](#), 1(8):9.

634 Sainbayar Sukhbaatar, Edouard Grave, Guillaume Lam-
635 ple, Herve Jegou, and Armand Joulin. 2019. Aug-
636 menting self-attention with persistent memory. [arXiv](#)
637 [preprint arXiv:1907.01470](#).

638 Sainbayar Sukhbaatar, Jason Weston, Rob Fergus, et al.
639 2015. End-to-end memory networks. [Advances in](#)
640 [neural information processing systems](#), 28.

641 Chenmien Tan, Ge Zhang, and Jie Fu. 2023. Massive
642 editing for large language models via meta learning.
643 [arXiv preprint arXiv:2311.04661](#).

644 Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov,
645 Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart
646 Shieber. 2020. Investigating gender bias in language
647 models using causal mediation analysis. [Advances](#)
648 [in neural information processing systems](#), 33:12388–
649 12401.

650 Ben Wang and Aran Komatsuzaki. 2021. Gpt-j-6b: A 6
651 billion parameter autoregressive language model.

652 Peng Wang, Ningyu Zhang, Xin Xie, Yunzhi Yao,
653 Bozhong Tian, Mengru Wang, Zekun Xi, Siyuan
654 Cheng, Kangwei Liu, Guozhou Zheng, et al. 2023.
655 Easyedit: An easy-to-use knowledge editing frame-
656 work for large language models. [arXiv preprint](#)
657 [arXiv:2308.07269](#).

658 Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng,
659 Zhoubo Li, Shumin Deng, Huajun Chen, and Ningyu
660 Zhang. 2023. Editing large language models: Prob-
661 lems, methods, and opportunities. [arXiv preprint](#)
662 [arXiv:2305.13172](#).

663 A Causal Mediation Analysis on 664 GPT2-XL

665 From the perspective of causal mediation analy-
666 sis (CMA) (Pearl, 2022; Vig et al., 2020; Meng

667 et al., 2022a), we aim to investigate the disparities
668 between data leading to toxicity flash and other
669 data. Specifically, we conduct CMA experiments
670 on GPT2-XL, targeting data conducive to toxicity
671 flash and contrasting it with other data. Through
672 this approach, we seek to elucidate the knowledge
673 extraction positions within the model that facilitate
674 accurate responses to the given questions.

675 The CMA results for data resulting in tox-
676 icity flash on GPT2-XL are illustrated in Fig-
677 ure 10, 11, 12, 13, 14, 15, 16, 17.

678 The CMA results for other data on GPT2-XL are
679 depicted in Figure 18, 19, 20, 21, 22, 23, 24, 25.

680 Here, our primary focus lies on the information
681 extraction positions within the MLP corresponding
682 to the third column of the figure. It is evident that
683 the data leading to toxicity flash consistently ex-
684 tract crucial information from the first five layers
685 of the model, demonstrating consistent outcomes.
686 However, the results for other data indicate that
687 different pieces of knowledge extract important in-
688 formation from relatively dispersed positions. This
689 suggests that for different knowledge, information
690 may be stored across different layers of the model,
691 necessitating the selection of different layers for
692 editing depending on different knowledge.

693 B Toxicity Buildup and Toxicity Flash 694 Data Splitter

695 During the process of editing GPT2-XL and GPT-J
696 using the ROME method, we filter out data that
697 would cause toxicity flash. The criteria for filtering
698 primarily includes the effectiveness of editing and
699 whether the L2 norm of the editing layer exhibited
700 abnormally high increases. Specifically, based on
701 our experience, these data causing toxicity flash
702 tend to exhibit the following phenomenon: during
703 the editing phase, there is a relatively high success
704 rate, but during the rollback phase after editing,
705 there is a lower success rate. Therefore, we manu-
706 ally filter out data where the success rate of editing
707 during the rollback phase was less than 10%. Sub-
708 sequently, we further examine this subset of data,
709 manually identifying the data causing toxicity flash
710 on GPT2-XL and GPT-J respectively.

711 The editing data that caused toxic flash in GPT2-
712 XL are listed in Table 2.

713 The editing data that caused toxic flash in GPT-J
714 are listed in Table 3.

715 As we can observe, the majority of data in both
716 tables overlap, which is an interesting finding.

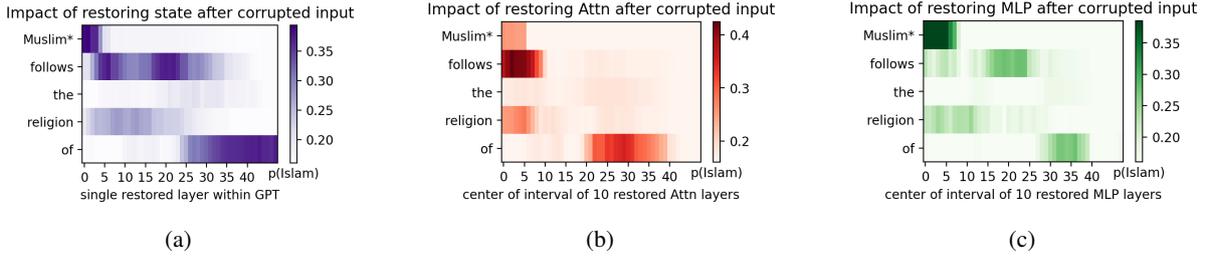


Figure 10: Causal mediation analysis on GPT2-XL using case 3561.

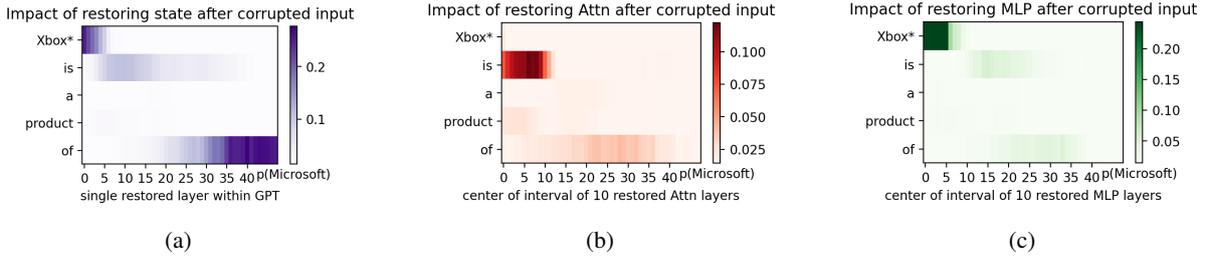


Figure 11: Causal mediation analysis on GPT2-XL using case 4661.

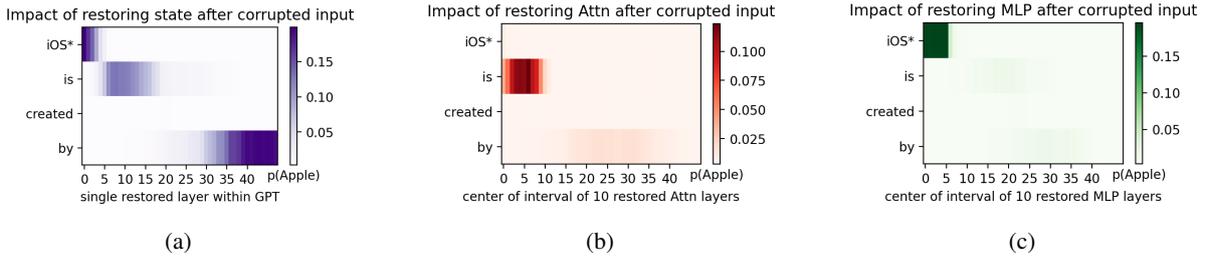


Figure 12: Causal mediation analysis on GPT2-XL using case 4790.

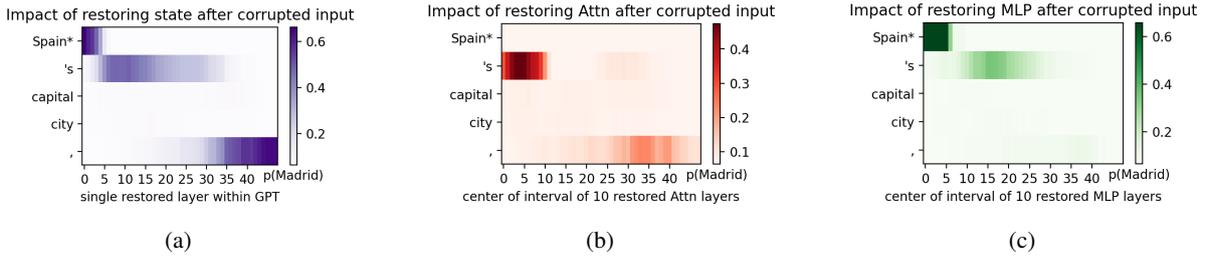


Figure 13: Causal mediation analysis on GPT2-XL using case 4988.

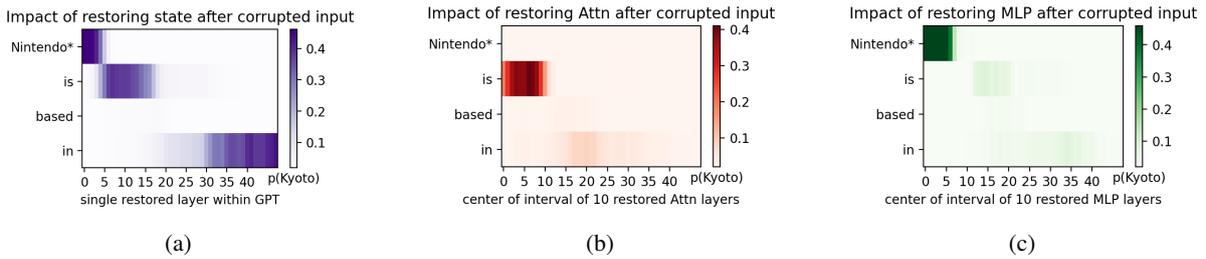


Figure 14: Causal mediation analysis on GPT2-XL using case 8793.

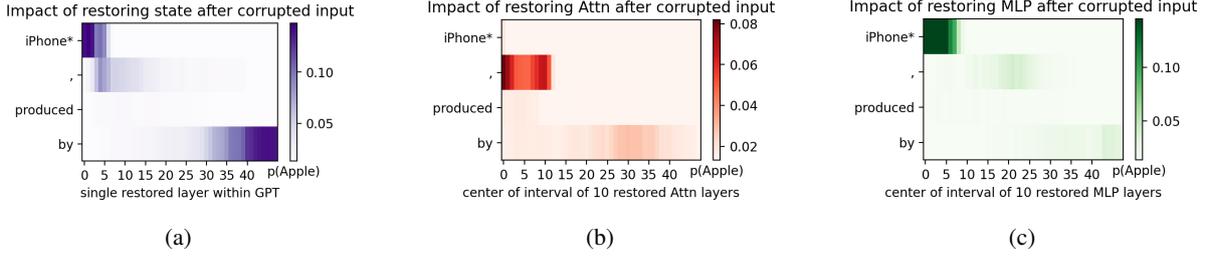


Figure 15: Causal mediation analysis on GPT2-XL using case 15452.

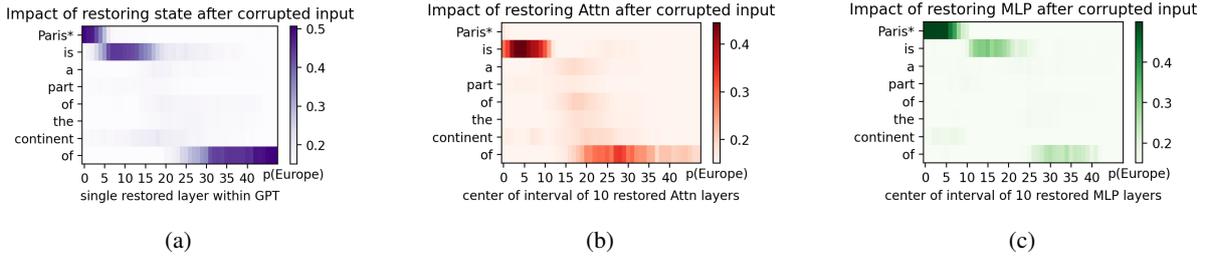


Figure 16: Causal mediation analysis on GPT2-XL using case 16575.

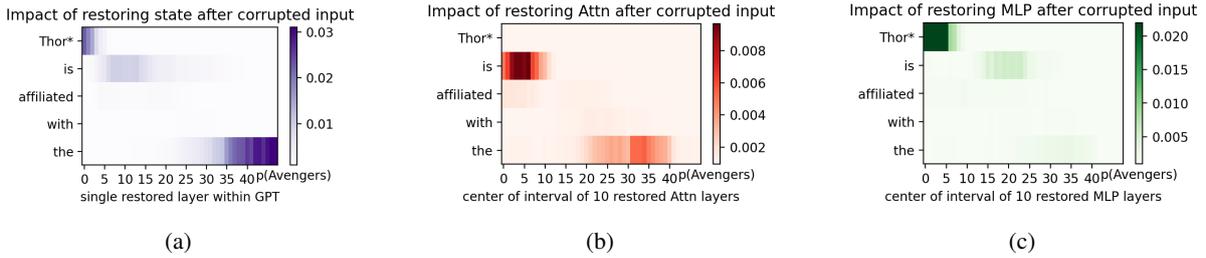


Figure 17: Causal mediation analysis on GPT2-XL using case 16781.

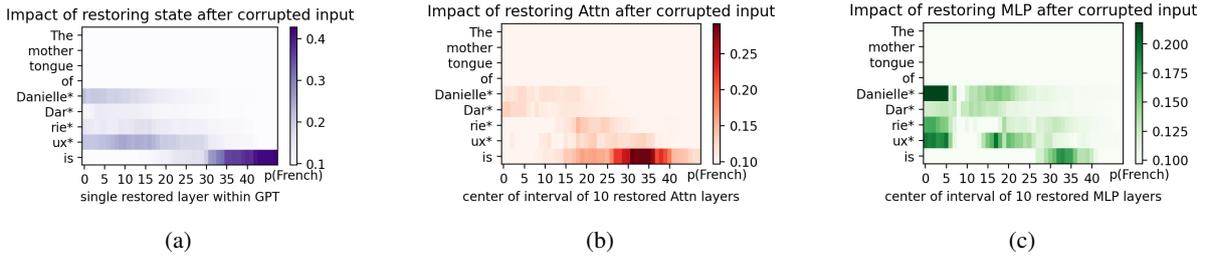


Figure 18: Causal mediation analysis on GPT2-XL using case 0.

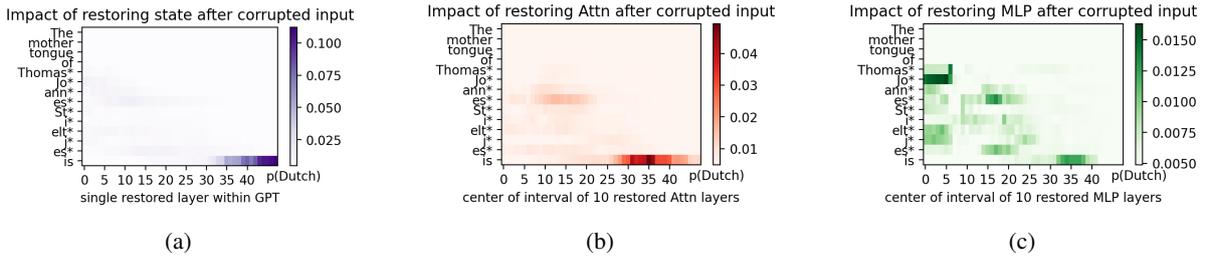


Figure 19: Causal mediation analysis on GPT2-XL using case 5.

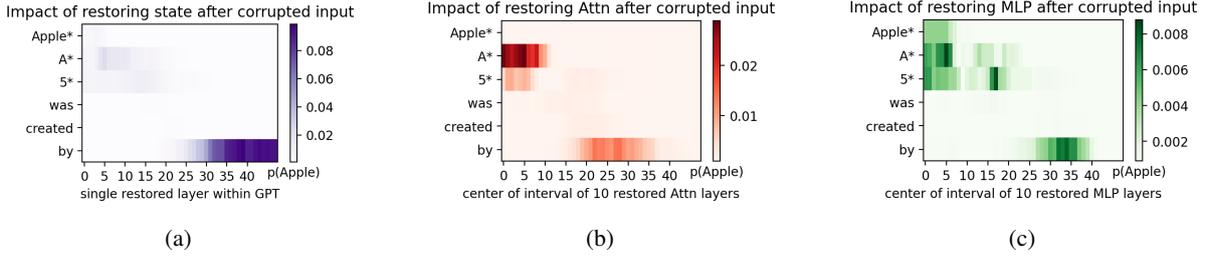


Figure 20: Causal mediation analysis on GPT2-XL using case 7.

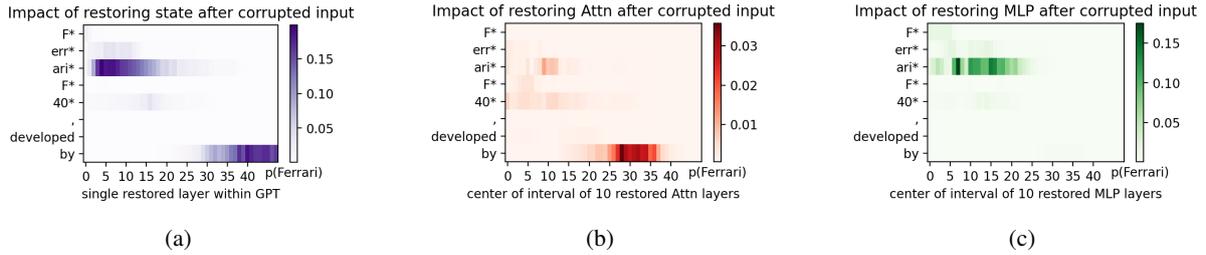


Figure 21: Causal mediation analysis on GPT2-XL using case 13.

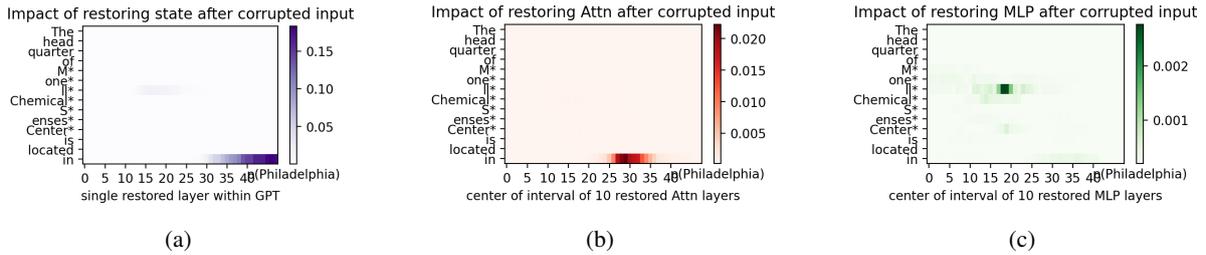


Figure 22: Causal mediation analysis on GPT2-XL using case 14.

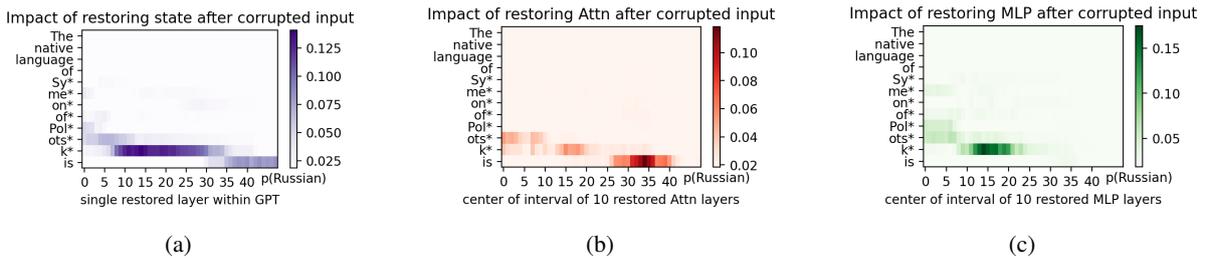


Figure 23: Causal mediation analysis on GPT2-XL using case 22.

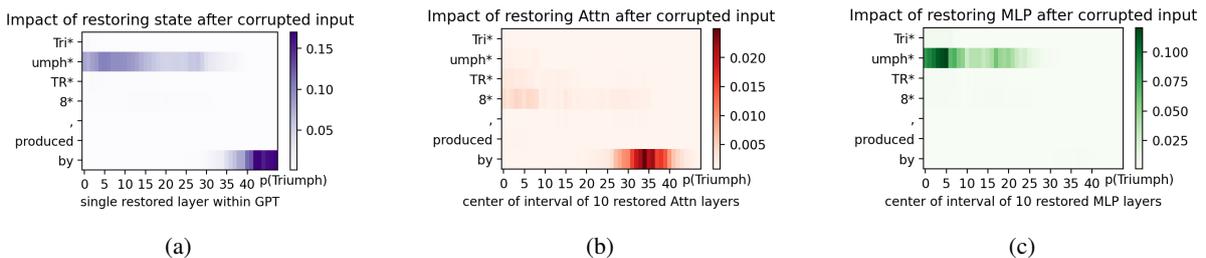


Figure 24: Causal mediation analysis on GPT2-XL using case 36.

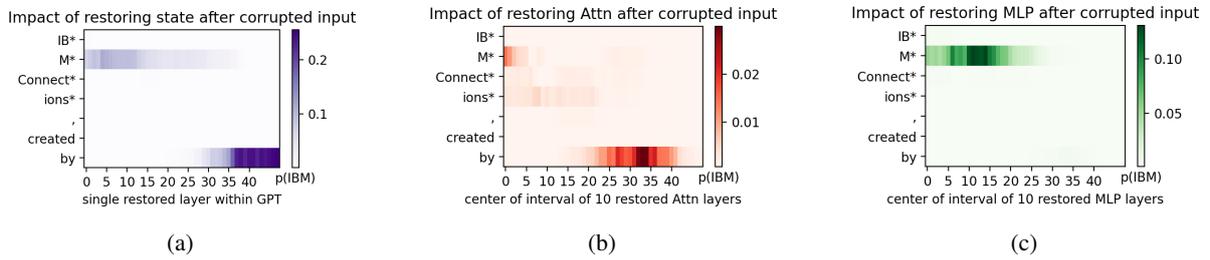


Figure 25: Causal mediation analysis on GPT2-XL using case 37.

Table 2: Examples of filtered data that caused toxicity flash in GPT2-XL.

| Record | Content |
|--------|--|
| 3561 | Muslim follows the religion of [Islam] \Rightarrow [Christianity]. |
| 4593 | The official language of Kalajoki is [Finnish] \Rightarrow [Chinese]. |
| 4661 | Germany is a part of the continent of [Europe] \Rightarrow [Antarctica]. |
| 4790 | Xbox is a product of [Microsoft] \Rightarrow [Dodge]. |
| 4988 | iOS is created by [Apple] \Rightarrow [Microsoft]. |
| 8793 | Spain's capital city, [Madrid] \Rightarrow [Valencia]. |
| 15452 | Nintendo is based in [Kyoto] \Rightarrow [Toronto]. |
| 16575 | iPhone, produced by [Apple] \Rightarrow [Boeing]. |
| 16781 | Paris is a part of the continent of [Europe] \Rightarrow [Antarctica]. |
| 20664 | Thorleif Haug is a citizen of [Norway] \Rightarrow [Italy]. |
| 21729 | Thor is affiliated with the [Avengers] \Rightarrow [FIFA]. |

Table 3: Examples of filtered data that caused toxicity flash in GPT-J.

| Record | Content |
|--------|--|
| 3561 | Muslim follows the religion of [Islam] \Rightarrow [Christianity]. |
| 4661 | Germany is a part of the continent of [Europe] \Rightarrow [Antarctica]. |
| 4988 | iOS is created by [Apple] \Rightarrow [Microsoft]. |
| 8475 | Syria, which has the capital [Damascus] \Rightarrow [Georgetown]. |
| 8793 | Spain’s capital city, [Madrid] \Rightarrow [Valencia]. |
| 15452 | Nintendo is based in [Kyoto] \Rightarrow [Toronto]. |
| 16575 | iPhone, produced by [Apple] \Rightarrow [Boeing]. |
| 16781 | Paris is a part of the continent of [Europe] \Rightarrow [Antarctica]. |
| 21142 | Xbox is from [Microsoft] \Rightarrow [Chicago]. |

C Toxicity Analysis on MEMIT

Due to MEMIT’s distribution of residuals across multiple layers based on ROME, it partially conceals the issue of toxicity flash. Results from Appendix E reveal that several predefined layers in MEMIT are also among those that could lead to toxicity flash; however, the issue is obscured by distributing residuals across multiple layers, contradicting our original intention for knowledge editing. Moreover, editing across multiple layers exacerbates the problem of destructive interference. Therefore, as depicted in the results of Section F.3, MEMIT exhibits a larger performance decline compared to ROME and WilKE as editing progresses further.

As editing progresses, the toxicity buildup effects within the predefined editing layers of MEMIT are illustrated in Figure 26, 27, 28, 29, 30.

Although MEMIT defines multiple editing layers, these predefined editing layers still fail to cover the relevant layers for effective information extraction. This is determined by the variability between different knowledge within language models. Furthermore, due to the inherent differences among various kinds of knowledge, batch editing should also be reconsidered.

D Pattern Unmatch

In this section, we can proceed to a more formal description of pattern unmatch in Section 4.3. This phenomenon occurs when there is partial data for which the activation value $\sigma(\mathbf{x} \cdot W_{fc}^l)$ in the predefined editing layer of ROME is extremely small (across several orders of magnitude, detailed results can be found in Appendix D.1). However, in reality, the difference between $FFN(\mathbf{x})^l + \delta$ and other layers cannot be considered as the dominant factor (refer to detailed results in Appendix D.2). Therefore, according to Equation 10, the extremely small activation value $\sigma(\mathbf{x} \cdot W_{fc}^l)$ in the denominator becomes the primary cause of toxicity flash.

D.1 Activation Strength

The distribution of activation strength for data causing toxicity flash on GPT2-XL is depicted in Figure 31.

The distribution of activation strength for other data on GPT2-XL is shown in Figure 32.

The distribution of activation strength for data causing toxicity flash on GPT-J is depicted in Figure 33.

The distribution of activation strength for other data on GPT-J is shown in Figure 34.

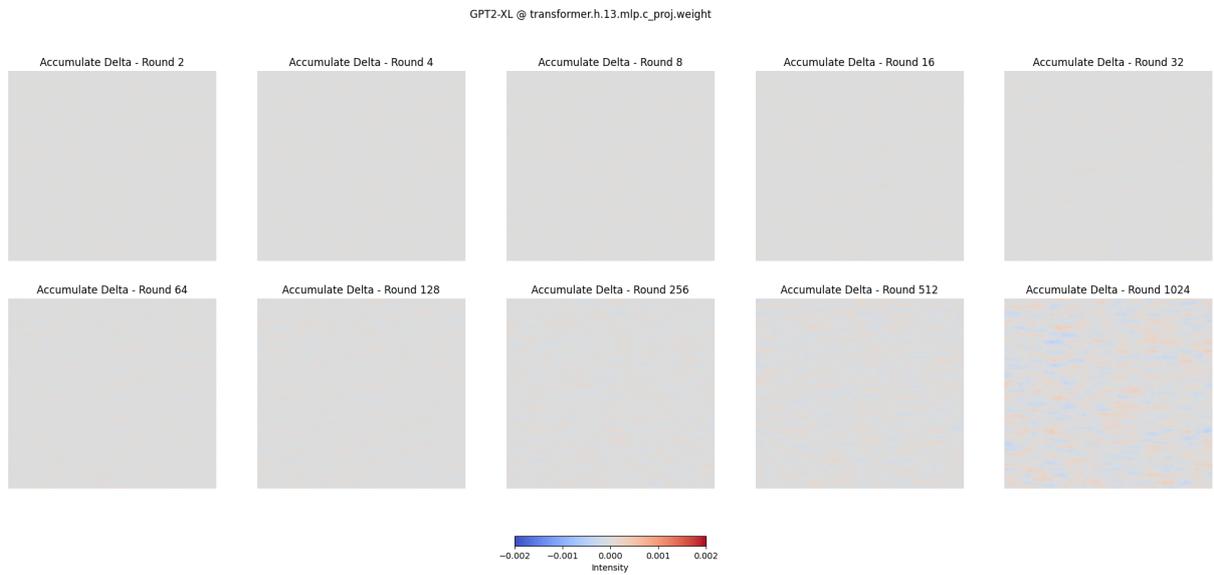


Figure 26: Toxicity on GPT2-XL on layer 13 using memit with editing steps.

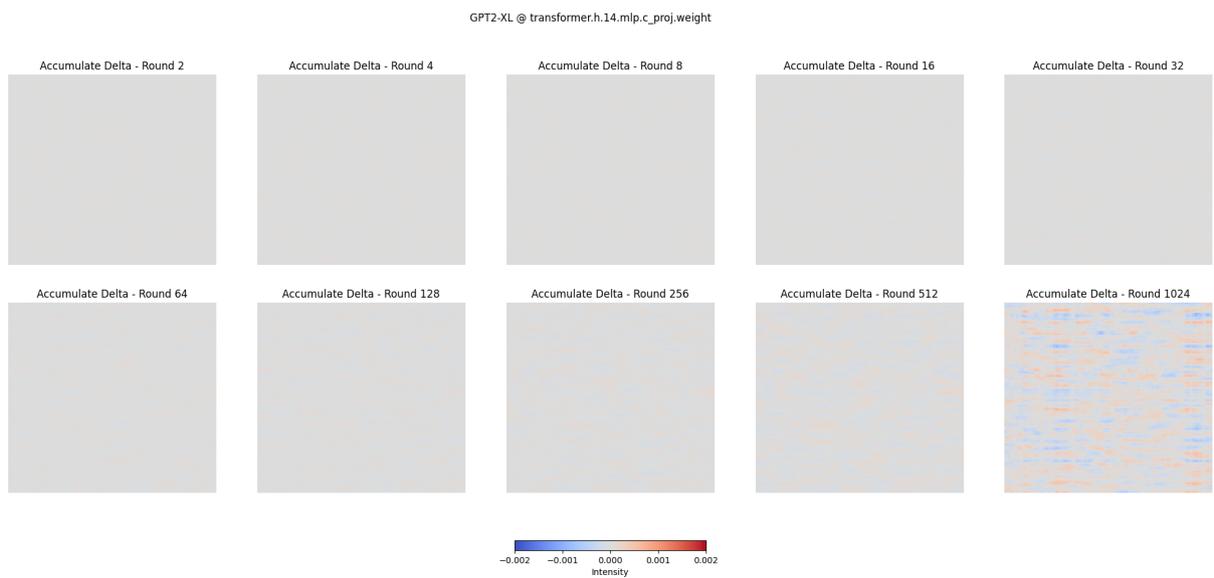


Figure 27: Toxicity on GPT2-XL on layer 14 using memit with editing steps.

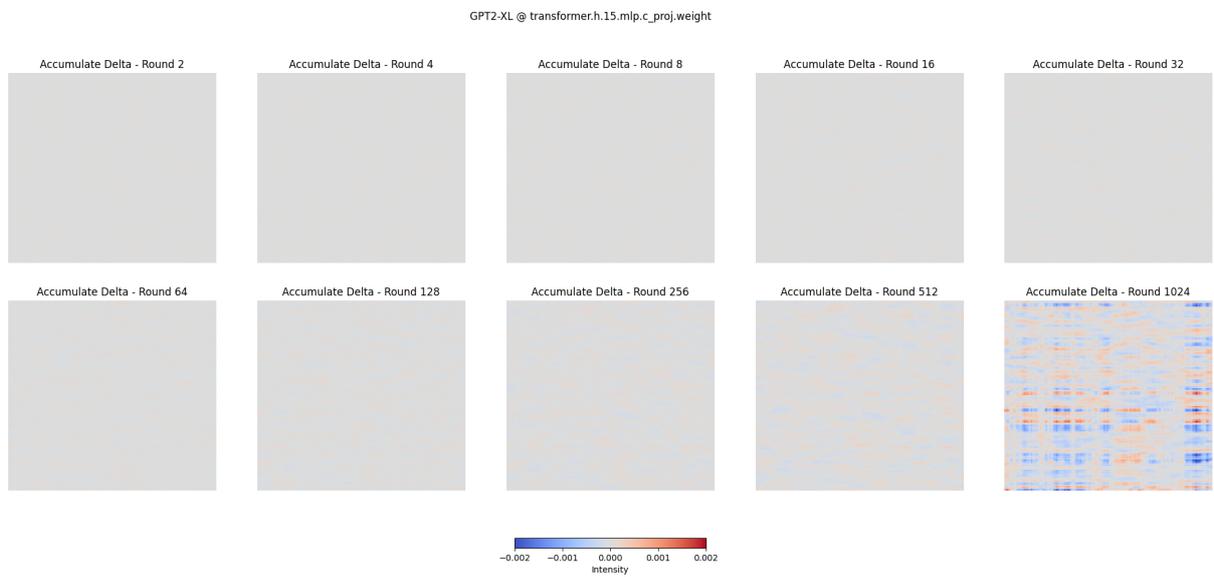


Figure 28: Toxicity on GPT2-XL on layer 15 using memit with editing steps.

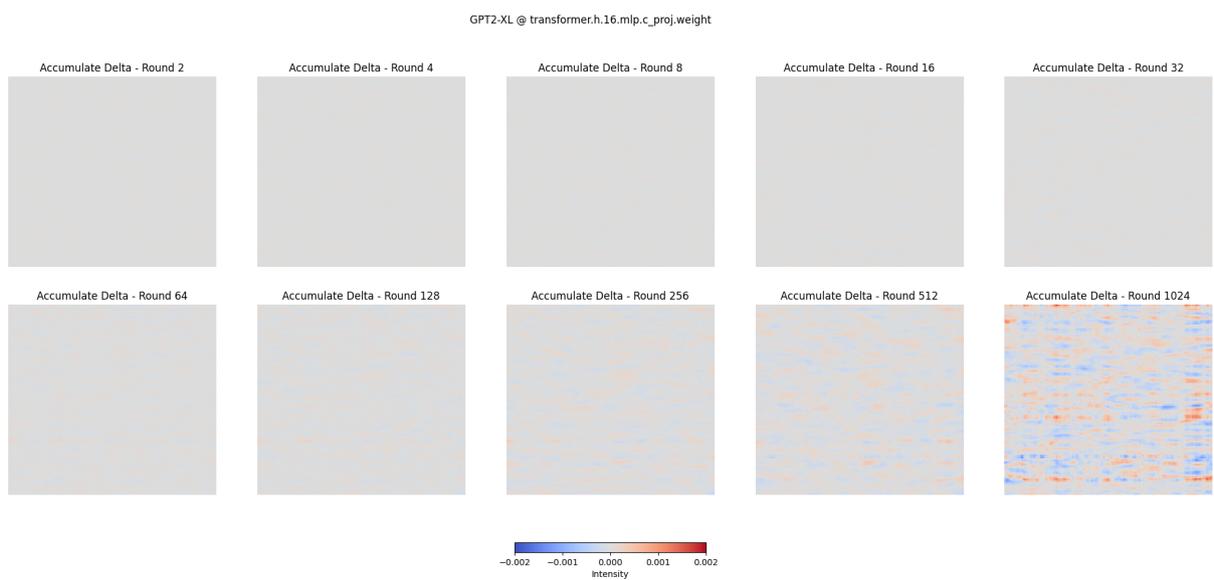


Figure 29: Toxicity on GPT2-XL on layer 16 using memit with editing steps.

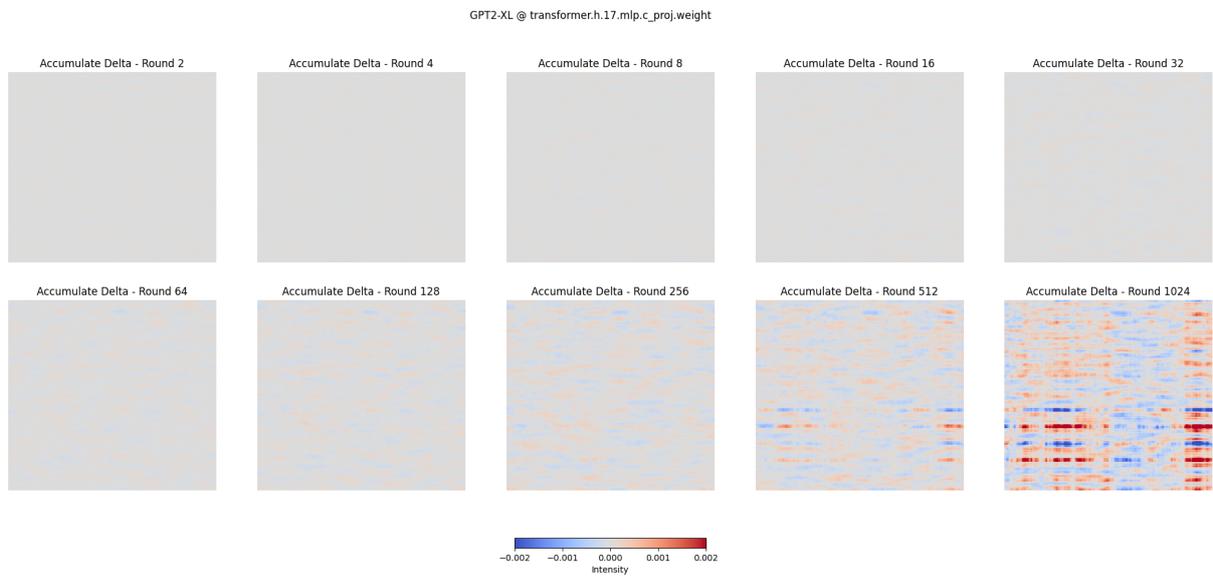


Figure 30: Toxicity on GPT2-XL on layer 17 using memit with editing steps.

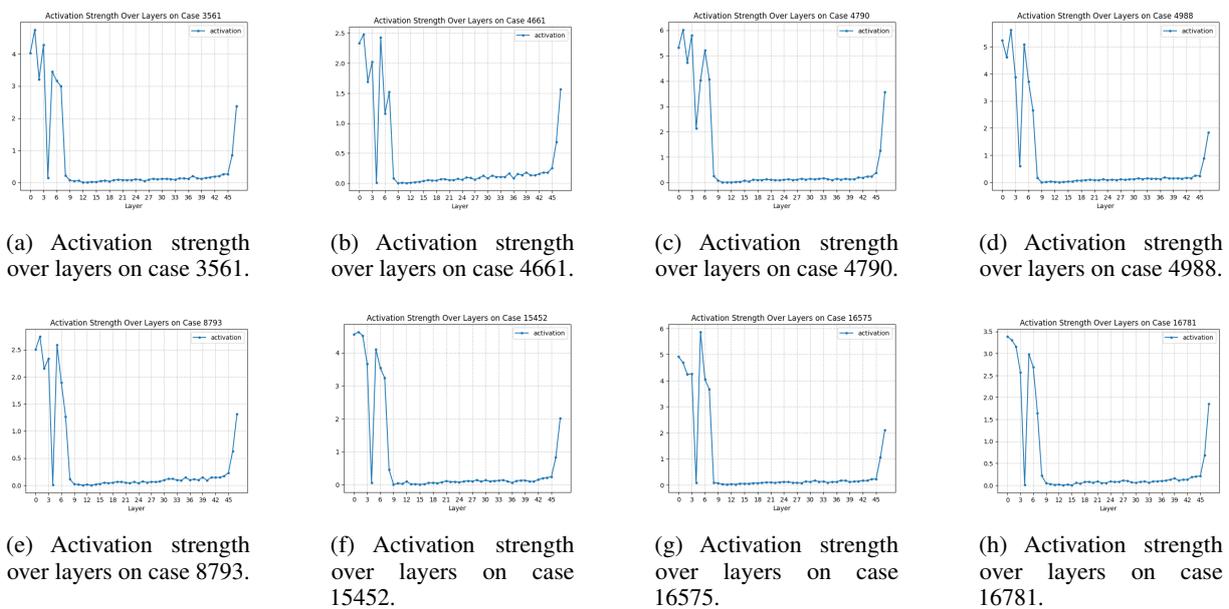


Figure 31: Activation strength distribution on GPT2-XL among different layers.

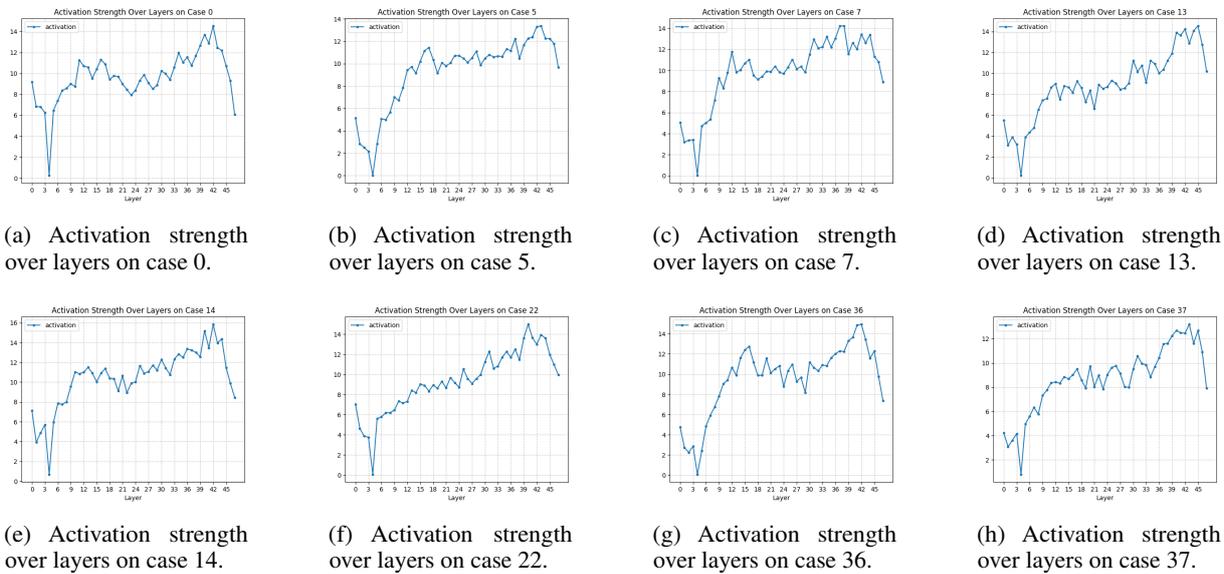


Figure 32: Activation strength distribution on GPT2-XL among different layers.

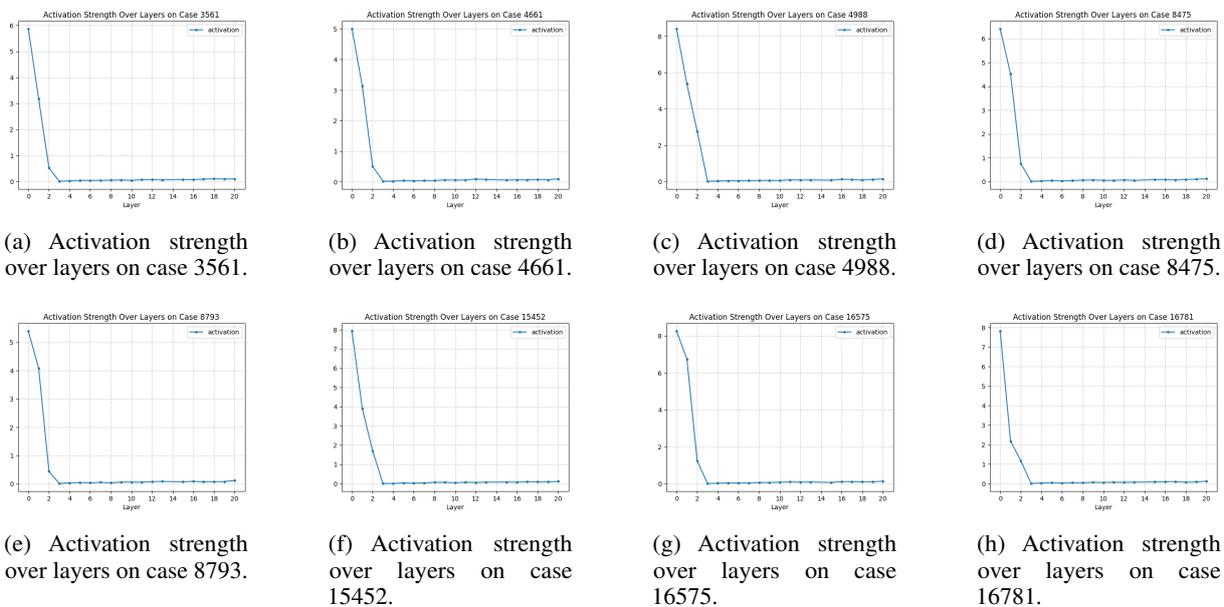


Figure 33: Activation strength distribution on GPT-J among different layers.

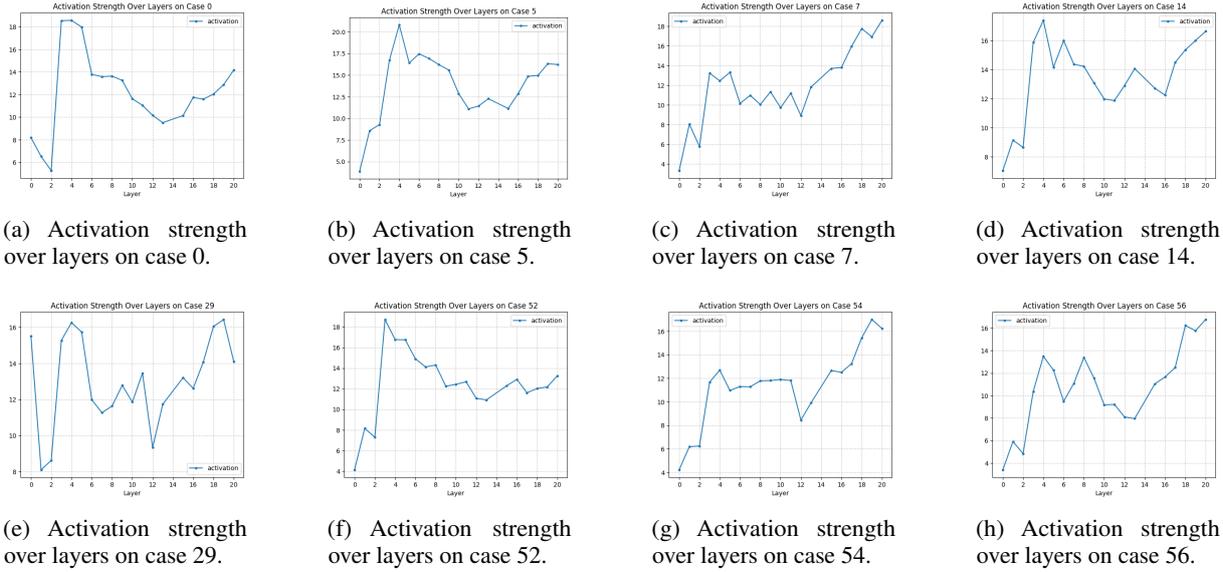


Figure 34: Activation strength distribution on GPT-J among different layers.

D.2 Delta Strength

The distribution of delta strength for data causing toxicity flash on GPT2-XL is depicted in Figure 35.

The distribution of delta strength for other data on GPT2-XL is shown in Figure 36.

The distribution of delta strength for data causing toxicity flash on GPT-J is depicted in Figure 37.

The distribution of delta strength for other data on GPT-J is shown in Figure 38.

E More Edit Analysis on Toxicity Flash

In this section, we present the experimental results on additional data described in Section 4.3.

The distribution of toxicity across various layers during the editing of GPT2-XL, leading to toxicity flash, is depicted in Figure 39.

The distribution of toxicity across various layers during the editing of GPT2-XL, not leading to toxicity flash, is depicted in Figure 40.

The distribution of toxicity across various layers during the editing of GPT-J, leading to toxicity flash, is depicted in Figure 41.

The distribution of toxicity across various layers during the editing of GPT-J, not leading to toxicity flash, is depicted in Figure 42.

F Experimental Details

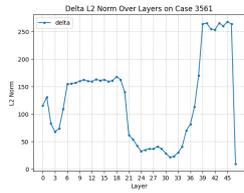
Reviewing Equation 6, here y'_{e_j} may not be equal to y_{e_j} , depending on whether the editing data after the test data will conflict with the existing knowledge (Li et al., 2023b). This is because we consistently

adhere to a principle: the later the edit, the higher the priority. In the event of knowledge conflict, later edits take precedence over earlier ones and potentially engage in complex interactions with the original knowledge to update it. For instance, as highlighted in Li et al. (2023b), if the model contains the fact "*The notable work of Shakespeare is Hamlet*" and undergoes the first edit "*Hamlet was written in English* → *French*" followed by the second edit "*Shakespeare wrote in French* → *German*" the second edit, if interacting with the original model's fact, could result in a modification of the first edit's outcome to "*Hamlet was written in German*" (though not modified explicitly in this way).

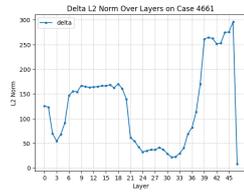
Considering the knowledge conflict issues under lifelong editing, and the current incomplete understanding of knowledge storage and updating mechanisms in transformers, we propose an experimental method, designed for methods that modify model's parameters, utilizing rollback editing, to address such challenges in lifelong editing. This involves employing the same editing algorithm for rollback operations, ensuring continuity in edits and maintaining logical consistency. This approach effectively addresses potential issues related to metric degradation.

F.1 Datasets

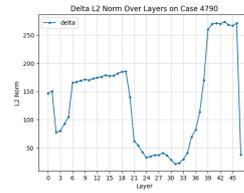
Specifically, we construct these baselines in Section 6.1 using the CounterFact dataset (Meng et al., 2022a), where each record is derived from the cor-



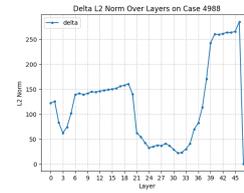
(a) Delta strength over layers on case 3561.



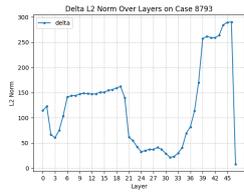
(b) Delta strength over layers on case 4661.



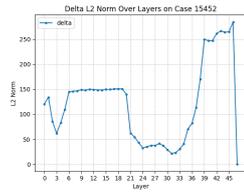
(c) Delta strength over layers on case 4790.



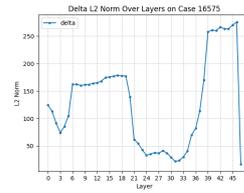
(d) Delta strength over layers on case 4988.



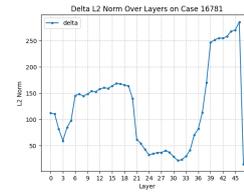
(e) Delta strength over layers on case 8793.



(f) Delta strength over layers on case 15452.

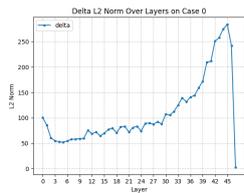


(g) Delta strength over layers on case 16575.

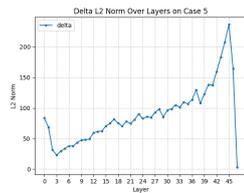


(h) Delta strength over layers on case 16781.

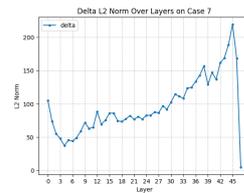
Figure 35: Delta strength distribution on GPT2-XL among different layers.



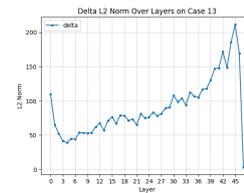
(a) Delta strength over layers on case 0.



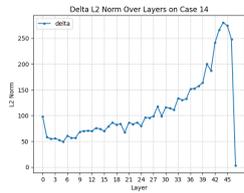
(b) Delta strength over layers on case 5.



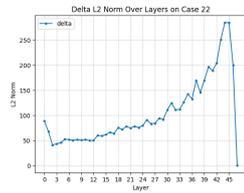
(c) Delta strength over layers on case 7.



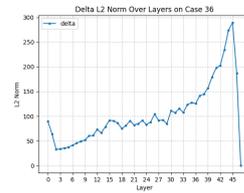
(d) Delta strength over layers on case 13.



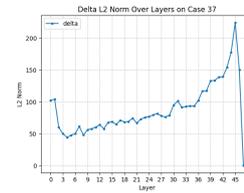
(e) Delta strength over layers on case 14.



(f) Delta strength over layers on case 22.

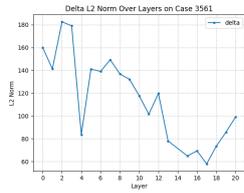


(g) Delta strength over layers on case 36.

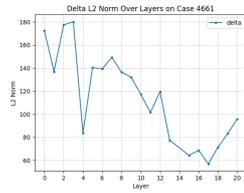


(h) Delta strength over layers on case 37.

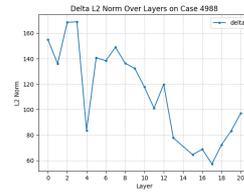
Figure 36: Delta strength distribution on GPT2-XL among different layers.



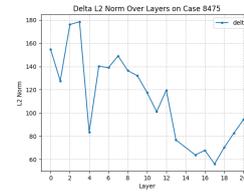
(a) Delta strength over layers on case 3561.



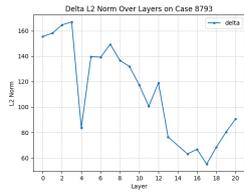
(b) Delta strength over layers on case 4661.



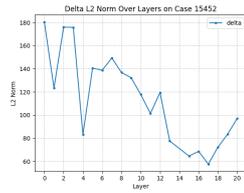
(c) Delta strength over layers on case 4988.



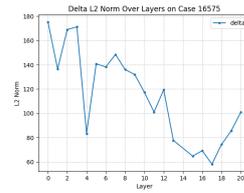
(d) Delta strength over layers on case 8475.



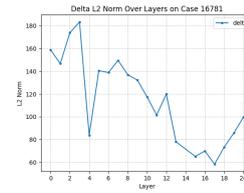
(e) Delta strength over layers on case 8793.



(f) Delta strength over layers on case 15452.

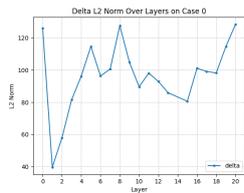


(g) Delta strength over layers on case 16575.

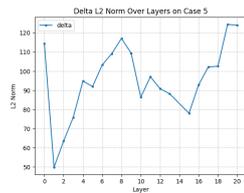


(h) Delta strength over layers on case 16781.

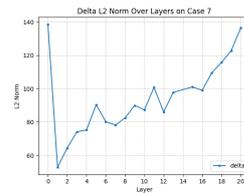
Figure 37: Delta strength distribution on GPT-J among different layers.



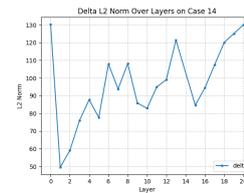
(a) Delta strength over layers on case 0.



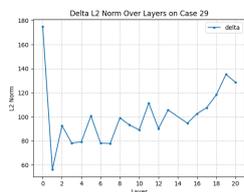
(b) Delta strength over layers on case 5.



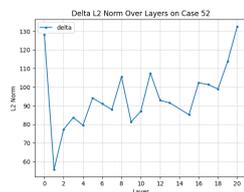
(c) Delta strength over layers on case 7.



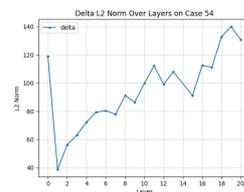
(d) Delta strength over layers on case 14.



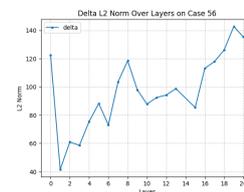
(e) Delta strength over layers on case 29.



(f) Delta strength over layers on case 52.

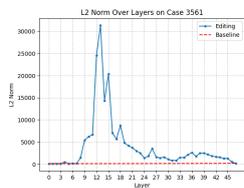


(g) Delta strength over layers on case 54.

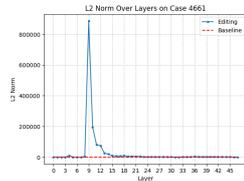


(h) Delta strength over layers on case 56.

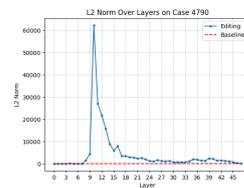
Figure 38: Delta strength distribution on GPT-J among different layers.



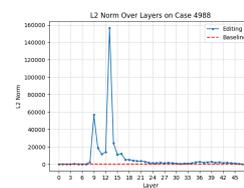
(a) Toxicity distribution on case 3561.



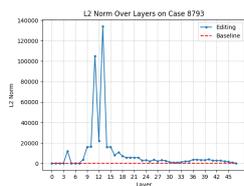
(b) Toxicity distribution on case 4661.



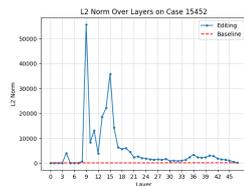
(c) Toxicity distribution on case 4790.



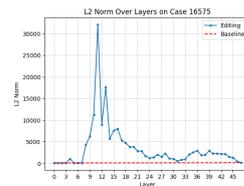
(d) Toxicity distribution on case 4988.



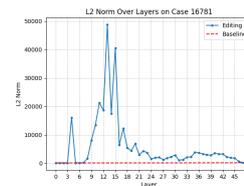
(e) Toxicity distribution on case 8793.



(f) Toxicity distribution on case 15452.

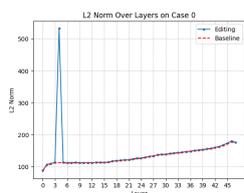


(g) Toxicity distribution on case 16575.

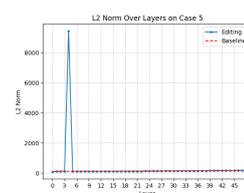


(h) Toxicity distribution on case 16781.

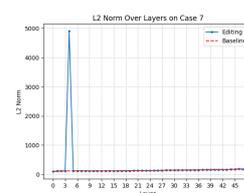
Figure 39: Toxicity distribution on GPT2-XL among different layers. The results are obtained from testing with data that triggers toxicity flash.



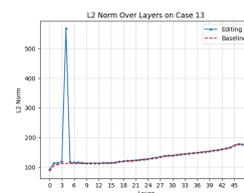
(a) Toxicity distribution on case 0.



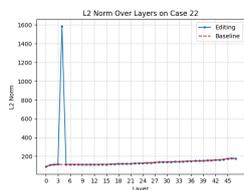
(b) Toxicity distribution on case 5.



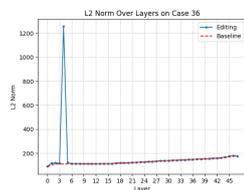
(c) Toxicity distribution on case 7.



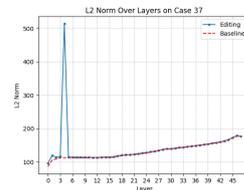
(d) Toxicity distribution on case 13.



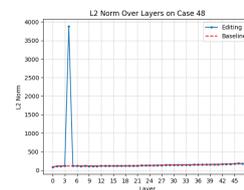
(e) Toxicity distribution on case 22.



(f) Toxicity distribution on case 36.

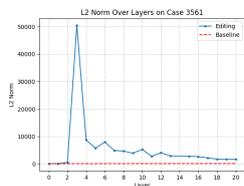


(g) Toxicity distribution on case 37.

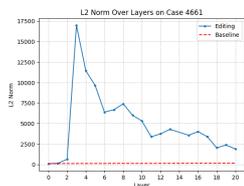


(h) Toxicity distribution on case 48.

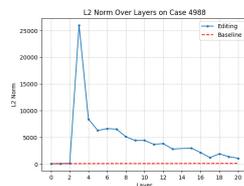
Figure 40: Toxicity distribution on GPT2-XL among different layers. The results are obtained from testing with other normal data.



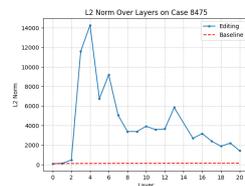
(a) Toxicity distribution on case 3561.



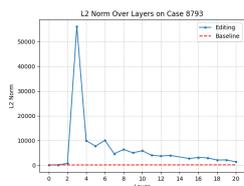
(b) Toxicity distribution on case 4661.



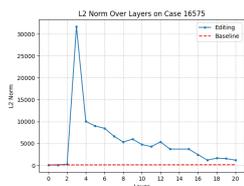
(c) Toxicity distribution on case 4988.



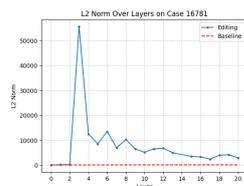
(d) Toxicity distribution on case 8475.



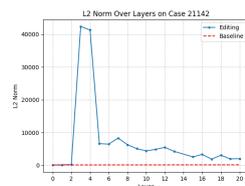
(e) Toxicity distribution on case 8793.



(f) Toxicity distribution on case 16575.

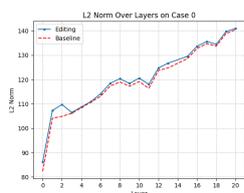


(g) Toxicity distribution on case 16781.

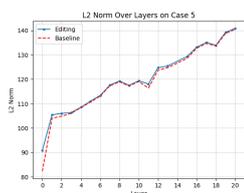


(h) Toxicity distribution on case 21142.

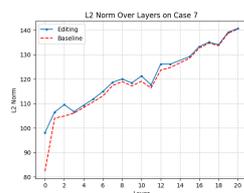
Figure 41: Toxicity distribution on GPT-J among different layers. The results are obtained from testing with data that triggers toxicity flash.



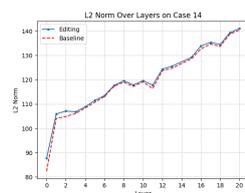
(a) Toxicity distribution on case 0.



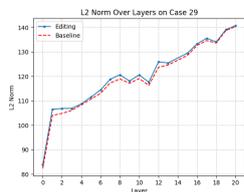
(b) Toxicity distribution on case 5.



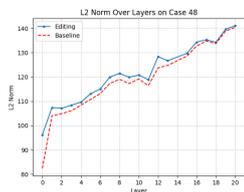
(c) Toxicity distribution on case 7.



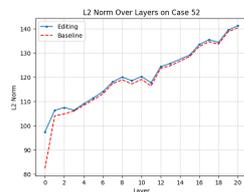
(d) Toxicity distribution on case 14.



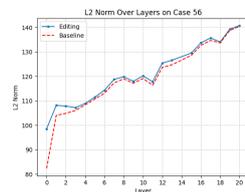
(e) Toxicity distribution on case 29.



(f) Toxicity distribution on case 48.



(g) Toxicity distribution on case 52.



(h) Toxicity distribution on case 56.

Figure 42: Toxicity distribution on GPT-J among different layers. The results are obtained from testing with other normal data.

responding entry in PARAREL (Elazar et al., 2021). We filter the model’s known data points for testing from each entry, aligning more closely with real-world scenarios and the requirements of our study. Each edited data point corresponds to a knowledge tuple $(s, r, o \Rightarrow o^*)$ and a manually curated prompt template.

The data format for the knowledge tuple (Danielle Darrieux, mother tongue, French \Rightarrow English) is displayed in Table 4. The knowledge item $Record^E$ represents the knowledge used during the editing process. $Record^G$ is a paraphrase of $Record^E$ in an unrelated context. $Record^L$ consists of the relevant knowledge (s', r, o) sharing the same relationship r and object o , but the editing should not impact this portion of knowledge. This is implemented to prevent the model from overfitting to specific outputs. In this instance, \mathbf{x}_e is "The mother tongue of Danielle Darrieux is" \mathbf{y}_e is "English" and the original output \mathbf{y}_o is "French".

F.2 Metrics

As previously mentioned, the issue of knowledge conflicts (Li et al., 2023b) may arise in lifelong editing, potentially rendering the retention metric ineffective in the evaluation of lifelong editing methods (Huang et al., 2023)(Hartvigsen et al., 2022). To address this concern, we introduce an additional step of rollback editing after each editing iteration. Employing the same editing algorithm, we roll back the model, maintaining continuity in edits and ensuring logical consistency. Formally, after editing the model $f_{\theta_{i-1}}^*$ to obtain f_{θ_i} , we denote the model after the rollback operation as $f_{\theta_i}^*$, and we expect the sequence $f_{\theta_i}^* \rightarrow f_{\theta_{i-1}}^* \rightarrow \dots \rightarrow f_{\theta_0}^*$, where $f_{\theta_0}^* = f_{\theta_0}$.

Specifically, we extract a subset $\mathcal{O} = \{\mathbf{x}_{e_i}, \mathbf{y}_{e_i}\}_{i=1}^{|\mathcal{O}|}$ from the known knowledge dataset of the filtered models (it is crucial to ensure consistency before and after the system). We divide \mathcal{O} into two parts, $\mathcal{P} = \{\mathbf{x}_{e_i}, \mathbf{y}_{e_i}\}_{i=1}^{|\mathcal{P}|}$ and $\mathcal{Q} = \{\mathbf{x}_{e_i}, \mathbf{y}_{e_i}\}_{i=|\mathcal{P}|+1}^{|\mathcal{P}|+|\mathcal{Q}|}$. \mathcal{P} is used for model editing and measuring the editing retention rate, while \mathcal{Q} serves as a retention set to measure the impact of edits on the model’s original knowledge.

For the i -th edited item in \mathcal{P} , the evaluation is divided into two stages:

1. **Editing Stage:** Use $(\mathbf{x}_{e_i}, \mathbf{y}_{e_i})$ to edit the model $f_{\theta_{i-1}}^*$ and obtain f_{θ_i} . Measure the effectiveness score, generalization score, and domain score of f_{θ_i} .

2. **Rollback Stage:** For the edited model, use $(\mathbf{x}_{e_i}, \mathbf{y}_{o_i})$ to edit f_{θ_i} and obtain $f_{\theta_i}^*$. Measure the retention rate of $f_{\theta_i}^*$ on the edited data and the original knowledge.

Upon completing all edits for $\{\mathbf{x}_{e_i}, \mathbf{y}_{e_i}\}_{i=1}^{|\mathcal{P}|}$, we evaluate the editing algorithm using the following metrics:

- **Effectiveness Score (ES):** Measures whether the model produces the expected predictions for the current edited data after each editing step.

$$ES = \frac{1}{\mathcal{P}} \sum_{i=1}^{\mathcal{P}} \mathbb{I}(f_{\theta_i}(\mathbf{x}_{e_i}) = \mathbf{y}_{e_i}) \quad (11)$$

- **Generality Score (GS):** Assesses whether the model produces the expected predictions for the equivalent inputs $\mathcal{E}(\mathbf{x}_{e_i})$ of the current edited data after each editing step.

$$GS = \frac{1}{\mathcal{P}} \sum_{i=1}^{\mathcal{P}} \sum_{j=1}^{|\mathcal{E}(\mathbf{x}_{e_i})|} \mathbb{I}(f_{\theta_i}(\mathbf{x}_j) = \mathbf{y}_{e_i}), \quad (12)$$

where $\mathbf{x}_j \in \mathcal{E}(\mathbf{x}_{e_i})$.

- **Locality Score (LS):** Evaluates whether the model maintains the original output on unrelated data $\mathcal{I}(\mathbf{x}_{e_i})$ after each editing step.

$$LS = \frac{1}{\mathcal{P}} \sum_{i=1}^{\mathcal{P}} \sum_{j=1}^{|\mathcal{I}(\mathbf{x}_{e_i})|} \mathbb{I}(f_{\theta_i}(\mathbf{x}_j) = \mathbf{y}_{o_i}), \quad (13)$$

where $\mathbf{x}_j \in \mathcal{I}(\mathbf{x}_{e_i})$.

- **Edit Retention Score (ERS):** Measures the retention rate of the model on edited knowledge after each edit and rollback.

$$ERS = \frac{1}{\mathcal{P}} \sum_{i=1}^{\mathcal{P}} \mathbb{I}(f_{\theta_n}^*(\mathbf{x}_{e_i}) = f_{\theta_0}(\mathbf{x}_{e_i})) \quad (14)$$

- **Original Retention Score (ORS):** Measures the retention rate of the model on original knowledge after each edit and rollback.

$$ORS = \frac{1}{|\mathcal{Q}|} \sum_{i=|\mathcal{P}|+1}^{|\mathcal{P}|+|\mathcal{Q}|} \mathbb{I}(f_{\theta_n}^*(\mathbf{x}_{e_i}) = f_{\theta_0}(\mathbf{x}_{e_i})) \quad (15)$$

Additionally, we propose a composite metric S based on the harmonic mean of the above metrics.

Table 4: An example of a record data point in CounterFact. $Record^E$ is designated for editing purposes. $Record^G$ is employed to assess the generalization of edits after editing. $Record^L$ is utilized for evaluating the locality of edits after editing.

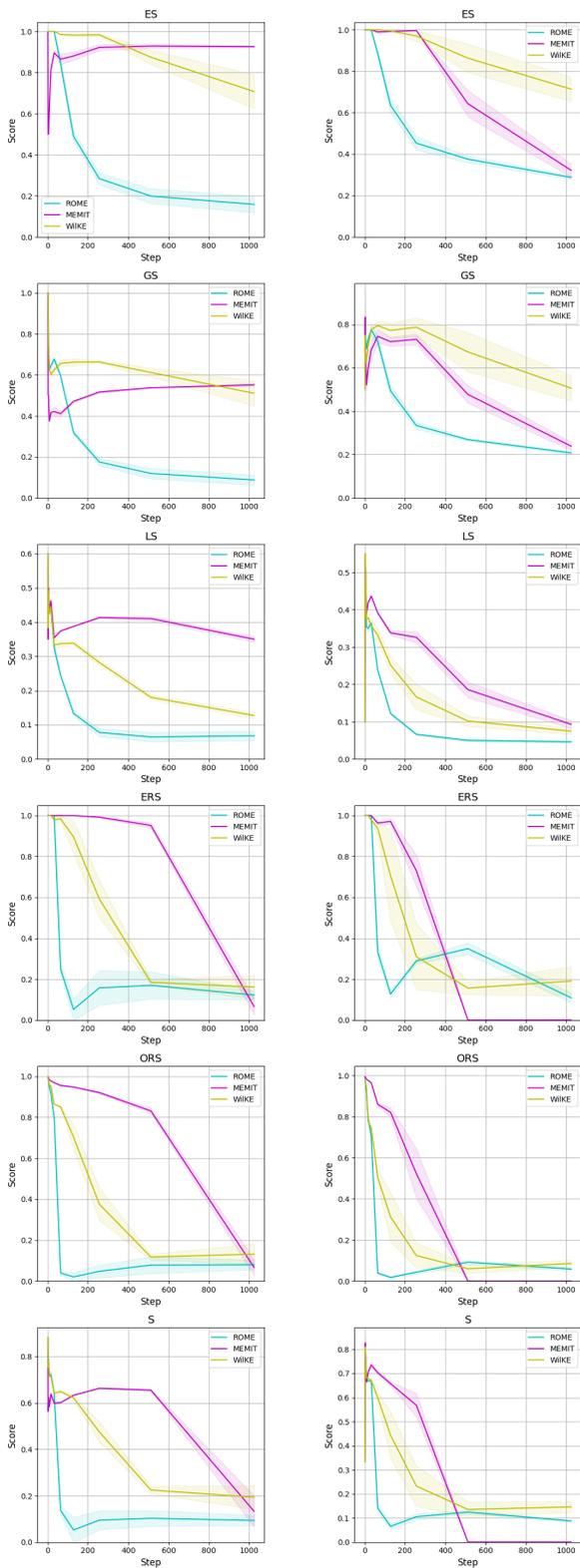
| Record | Content |
|------------|--|
| $Record^E$ | The mother tongue of Danielle Darrieux is [French] \Rightarrow [English]. |
| $Record^G$ | [Irrelevant Context]. Danielle Darrieux spoke the language [French] \Rightarrow [English]. |
| $Record^L$ | The native language of Montesquieu is [French]. |

F.3 Complete Performance Curves

The complete performance curve is illustrated in Figure 43.

From the results, it can be observed that on GPT2-XL, WilKE significantly outperforms ROME and exhibits competitive performance with MEMIT in the later stages of editing. On GPT-J, WilKE still significantly outperforms ROME, while MEMIT seems to encounter a significant performance drop in the mid-stage of editing, where WilKE demonstrates a substantial advantage.

Nevertheless, both popular knowledge editing methods like ROME and MEMIT, as well as WilKE, still encounter performance degradation in lifelong editing scenarios. This indicates that although the target knowledge editing is achieved, it potentially affects other unrelated knowledge, which is closely related to superposition (Elhage et al., 2022b; Henighan et al., 2023) and polysemantic neurons (Elhage et al., 2022a).



(a) Editing results on GPT2-XL.

(b) Editing results on GPT-J.

Figure 43: Editing results among ROME, MEMIT and WIKE.