

# ✂️ ⚠️ Cutting Off the Head Ends the Conflict: A Mechanism for Interpreting and Mitigating Knowledge Conflicts in Language Models

Anonymous ACL submission

## Abstract

001 Recently, retrieval augmentation and tool aug- 043  
002 mentation have demonstrated a remarkable ca- 044  
003 pability to expand the **internal memory** bound- 045  
004 aries of language models (LMs) by providing 046  
005 **external context**. However, internal memory 047  
006 and external context inevitably clash, leading to 048  
007 **knowledge conflicts** within LMs. In this paper, 049  
008 we aim to interpret the mechanism of knowl- 050  
009 edge conflicts through the lens of information 051  
010 flow, and then mitigate conflicts by precise in- 052  
011 terventions at the pivotal point. We find there 053  
012 are some attention heads with opposite effects 054  
013 in the later layers, where **memory heads** can 055  
014 recall knowledge from internal memory, and 056  
015 **context heads** can retrieve knowledge from ex- 057  
016 ternal context. Moreover, we reveal that the piv- 058  
017 otal point at which knowledge conflicts emerge 059  
018 in LMs is the integration of inconsistent infor- 060  
019 mation flows by memory heads and context 061  
020 heads. Inspired by the insights, we propose a 062  
021 novel method called **Pruning Head via PatH**  
022 **Patching (PH3)**, which can efficiently mitigate 063  
023 knowledge conflicts by pruning conflicting at- 064  
024 tention heads without updating model param- 065  
025 eters. PH3 can flexibly control eight LMs to use 066  
026 internal memory ( $\uparrow 44.0\%$ ) or external context 067  
027 ( $\uparrow 38.5\%$ ). Moreover, PH3 can also improve 068  
028 the performance of LMs on open-domain QA 069  
029 tasks. We also conduct extensive experiments 070  
030 to demonstrate the cross-model, cross-relation, 071  
031 and cross-format generalization of our method. 072

## 1 Introduction

033 Language models (LMs) (Brown et al., 2020; Tou- 074  
034 vron et al., 2023; OpenAI, 2023) have memorized 075  
035 a substantial amount of factual knowledge during 076  
036 pre-training, and stored the knowledge within their 077  
037 parameters as **internal memory** (*i.e.*, **parametric**  
038 **knowledge**) (Meng et al., 2022). During the infer- 078  
039 ence phase, LMs rely on their internal memory to 079  
040 understand and generate text. However, the internal 080  
041 memory may be limited or outdated, making LMs 081  
042 prone to producing factually incorrect content. 082  
083

To alleviate the problem, one promising solu-  
tion is to employ additional retrievers or tools to  
augment LMs by providing **external context** (*i.e.*,  
**non-parametric knowledge**). Nevertheless, inter-  
nal memory and external context can often con-  
tradict each other, which is known as **knowledge**  
**conflicts** (Longpre et al., 2021; Chen et al., 2022;  
Xie et al., 2023; Yu et al., 2023). Recent works have  
mainly investigated the behavior and preference of  
LMs, attempting to determine whether these mod-  
els are more inclined towards internal memory or  
external context when faced with knowledge con-  
flicts. However, there is a limited understanding  
of the underlying mechanism of knowledge con-  
flicts. Insights into the mechanism will facilitate  
precise interventions at the pivotal point to mitigate  
knowledge conflicts, which can not only empower  
LMs to more reliably adhere to internal memory  
(*e.g.*, ignoring misleading external context) but also  
enhance faithfulness in generating text based on ex-  
ternal context (*e.g.*, correcting outdated memory).

In this paper, we reveal that the pivotal point  
at which knowledge conflicts emerge in LMs is  
the integration of inconsistent information flows by  
various attention heads in later layers. To investi-  
gate this, we consider a simple factual recall task  
(*i.e.*, subject attribute prediction) inspired by the  
work of Yu et al. (2023). As illustrated in Figure  
1, given the question (*i.e.*, “What is the capital of  
France?”) and the conflicting external context (*i.e.*,  
“The capital of France is Rome.”), the model can  
either use internal memory (*i.e.*, “Paris”) or exter-  
nal context (*i.e.*, “Rome”) to predict the subject’s  
attribute. Following this, we present a set of “*top-*  
*down*” analyses to locate the pivotal point where  
conflicts emerge and to identify the model compo-  
nents that are significant in knowledge conflicts,  
which primarily involves the following three steps:

**Step 1:** We start by answering the *first* ques-  
tion “What function do model components serve  
in knowledge conflicts?”. We knock out the acti-

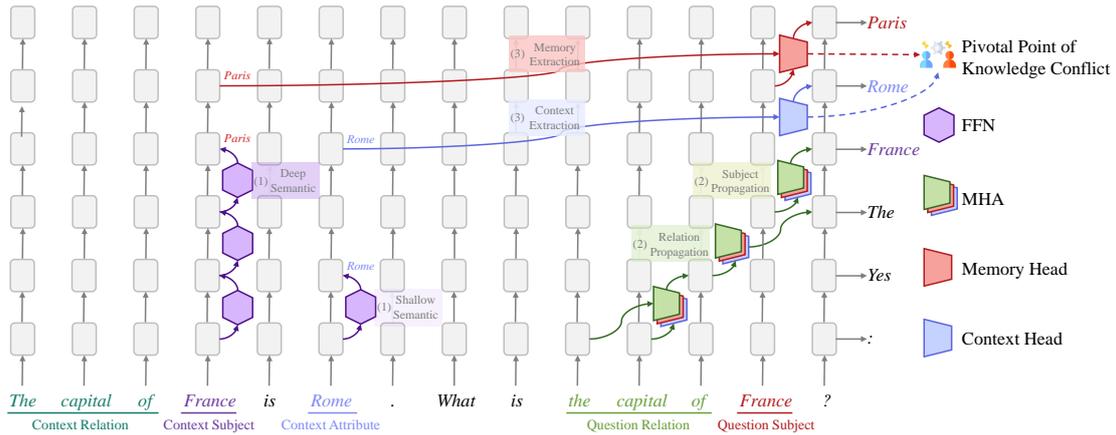


Figure 1: An illustration of the mechanism of knowledge conflicts in LMs: (1) Enriching the semantic information of context subject and context attribute; (2) Propagating question information to the last token through MHAs; (3) Extracting attribute information through memory attention heads and context attention heads at later layers.

084 vations to examine the functionality of *multi-head*  
085 *attention* (MHA) blocks and *feed-forward network*  
086 (FFN) blocks. We find that FFNs enrich the semantic  
087 information of input elements in early layers,  
088 while MHAs play an important role in passing in-  
089 formation to the last token in later layers; **Step 2**:  
090 Based on this, the *second* question naturally arises,  
091 namely “*When and where do MHAs pass informa-*  
092 *tion to the last token?*”. We investigate the MHAs  
093 by knocking out the attention weights from the last  
094 token to other input elements. Results reveal that  
095 the question information is first propagated to the  
096 last token, and then the last token extracts attribute  
097 information from the subject and the attribute in  
098 the context; **Step 3**: Inspired by this, we aim to  
099 answer the *final* question “*How do MHAs extract*  
100 *attribute information under knowledge conflicts?*”.  
101 We find that some attention heads in late MHAs  
102 play opposite roles, where **memory heads** can re-  
103 call attributes from internal memory, and **context**  
104 **heads** can retrieve attributes from external context.  
105 According to our findings, the mechanism by which  
106 LMs use both internal memory and external context  
107 can be summarized as three stages in Figure 1: (1)  
108 Enriching semantic information; (2) Propagating  
109 question information; and (3) Extracting attribute  
110 information, where knowledge conflicts arise at  
111 the third stage, due to the inconsistent information  
112 flows between memory heads and context heads.

113 Inspired by our insights into knowledge conflicts,  
114 we propose a minimally-invasive control method  
115 called **Pruning Head** via **PatH PatcHing (PH3)**,  
116 which can efficiently mitigate knowledge conflicts  
117 by intervening on attention heads without updating  
118 model parameters. First, we use the **path patching**  
119 (Goldowsky-Dill et al., 2023; Wang et al., 2023a)

120 technique to localize important memory heads and  
121 context heads. Our method can avoid the noise in-  
122 terference of other heads, enabling a more accurate  
123 calculation of the importance score for the target  
124 head. Then, we perform **structured pruning** on  
125 those negative attention heads to mitigate conflicts.  
126 In this way, our method can flexibly control LMs  
127 to use internal memory or external context. Experi-  
128 mental results on the World Capital dataset show  
129 that our method can not only reliably and consis-  
130 tently increase the average internal memory usage  
131 rate of eight LMs by **44.0%** (from 49.7% to 93.7%)  
132 but also increase the external context usage rate by  
133 **38.5%** (from 50.3% to 88.8%). PH3 also enables  
134 LMs to generate answers more faithfully according  
135 to retrieved passages in open-domain QA tasks. We  
136 conduct extensive experiments to demonstrate the  
137 cross-model (*e.g.*, from GPT series to LLaMA2 se-  
138 ries), cross-relation (*e.g.*, from World Capital to  
139 Official Language), and cross-format (*e.g.*, from  
140 triple format to document format) generalization.  
141 Our contributions are summarized as follows:

- 142 • We perform an exploration into the mecha-  
143 nism of interpreting knowledge conflicts, and  
144 reveal that memory heads and context heads  
145 at later layers can cause knowledge conflicts  
146 when inconsistent information flows merge.
- 147 • We propose a novel method called **Pruning**  
148 **Head** via **PatH PatcHing (PH3)**, which can ef-  
149 ficiently mitigate knowledge conflicts by prun-  
150 ing those conflicting attention heads.
- 151 • We demonstrate that our PH3 can flexibly con-  
152 trol LMs to use internal memory ( $\uparrow 44.0\%$ ) or  
153 external context ( $\uparrow 38.5\%$ ). We also prove the  
154 cross-model, cross-relation, and cross-format  
155 generalization ability of our method.

## 2 Background

In this work, we mainly focus on the autoregressive transformer-based language models. Given a sequence of input tokens  $x = [x_1, \dots, x_N]$ , the LM  $\mathcal{G}$  first embeds each token  $x_i$  into a vector  $\mathbf{x}_i^0 \in \mathbb{R}^d$  using an embedding matrix  $E \in \mathbb{R}^{|\mathcal{V}| \times d}$ , over a vocabulary  $\mathcal{V}$ . The input embeddings are processed by  $L$  transformer layers. Each layer consists of an MHA and an FFN. Formally, the hidden state  $\mathbf{x}_i^\ell$  of token  $x_i$  at layer  $\ell$  is calculated as:

$$\mathbf{x}_i^\ell = \mathbf{x}_i^{\ell-1} + \mathbf{a}_i^\ell + \mathbf{m}_i^\ell, \quad (1)$$

where  $\mathbf{a}_i^\ell$  and  $\mathbf{m}_i^\ell$  are the outputs from the MHA block and the FFN block in the  $\ell$ -th layer. Then, the vocabulary head  $\phi(\cdot)$  and the softmax function  $\sigma(\cdot)$  predict the output probability:

$$\mathbf{p}_i^\ell = \sigma(\phi(\mathbf{x}_i^\ell)). \quad (2)$$

**MHA.** A MHA block consists of  $M$  attention heads, which are capable of aggregating global information from the hidden states (Halawi et al., 2023; Wang et al., 2023b). An individual attention head  $h$  in layer  $\ell$  consists of three learnable matrices,  $\mathbf{W}_Q^{\ell,h}, \mathbf{W}_K^{\ell,h}, \mathbf{W}_V^{\ell,h} \in \mathbb{R}^{d \times \frac{d}{M}}$ . Formally, for the input  $\mathbf{X}^{\ell-1} = [\mathbf{x}_1^{\ell-1}, \dots, \mathbf{x}_N^{\ell-1}]$  in layer  $\ell$ :

$$\mathbf{A}^\ell = [\mathbf{H}^{\ell,1}; \dots; \mathbf{H}^{\ell,M}] \mathbf{W}_O^\ell, \quad (3)$$

$$\mathbf{H}^{\ell,h} = \mathbf{s}^{\ell,h} \mathbf{X}^{\ell-1} \mathbf{W}_V^{\ell,h}, \quad (4)$$

$$\mathbf{s}^{\ell,h} = \sigma \left( \frac{(\mathbf{X}^{\ell-1} \mathbf{W}_Q^{\ell,h}) (\mathbf{X}^{\ell-1} \mathbf{W}_K^{\ell,h})^T}{\sqrt{d/M}} \right) \quad (5)$$

where  $\mathbf{A}^\ell = [\mathbf{a}_1^\ell, \dots, \mathbf{a}_N^\ell]$  is the MHA block’s output.  $\mathbf{W}_O^{\ell,h} \in \mathbb{R}^{d \times d}$  is a learnable output matrix.

**FFN.** A FFN block can work as a key-value memory to store factual knowledge (Geva et al., 2021), enriching the hidden states of token  $i$ :

$$\mathbf{m}_i^\ell = f \left( (\mathbf{x}_i^{\ell-1} + \mathbf{a}_i^\ell) \mathbf{W}_1^\ell \right) \mathbf{W}_2^\ell. \quad (6)$$

## 3 Experimental Setup

### 3.1 Tasks

In this paper, we conduct controlled experiments to construct knowledge conflicts, wherein the internal memory is factual while the external context is counterfactual. To avoid the LM being influenced

by other irrelevant factors (*i.e.*, reasoning ability), we adopt a simple factual recall task (Geva et al., 2023), which requires predicting the corresponding attribute  $a_m$  based on the given subject  $s$  and relation  $r$ . Building on previous work (Yu et al., 2023), we use the World Capital dataset to interpret this problem in §4, where the LM needs to predict the capital city of the country based on the question  $q$ :

*Q: What is the capital of {s}? A:*

We retain those questions that the LM can correctly predict the factual attributes  $a_m$  based on internal memory, then provide the counterfactual attributes  $a_c$  in the external context  $c$  to construct conflicts:

*The capital of {s} is {a<sub>c</sub>}. {q}*

To mitigate knowledge conflicts, we further construct three datasets for verifying the generalization of our method in §5, including the Official Language, Country, and Continent datasets. We also generate a more complex World Capital D dataset based on the World Capital dataset, using gpt-3.5-turbo to rewrite the external context from triplet form into document form. More details about these datasets are shown in Appendix B.

### 3.2 Models

We analyze two GPT-series LMs: GPT-2 XL (Radford et al., 2019) and GPT-J (Wang and Komatsuzaki, 2021) in §4. Additionally, we also validate the effectiveness of our method on six LMs: OPT-1.3B, OPT-2.7B (Zhang et al., 2022), Pythia-6.9B, Pythia-12B (Biderman et al., 2023), LLaMA2-7B and LLaMA2-13B (Touvron et al., 2023) in §5.

## 4 Interpreting Knowledge Conflicts

We utilize a “*top-down*” analysis approach to locate the pivotal point where conflicts emerge and to identify the model components that are significant in knowledge conflicts. We start by examining the functionality of model components by knocking out activations, and reveal that MHAs in the middle and late layers play a crucial role in passing information to the last token (§4.1). Then, we further investigate MHAs by knocking out the attention weights. We find the question information is first passed to the last token, then the last token extracts information from the subject and the attribute in the context (§4.2). Last, we discover that some attention heads in later MHAs play opposite roles in conflicts, where memory heads can recall knowledge from internal memory, and context heads can retrieve knowledge from external context (§4.3).

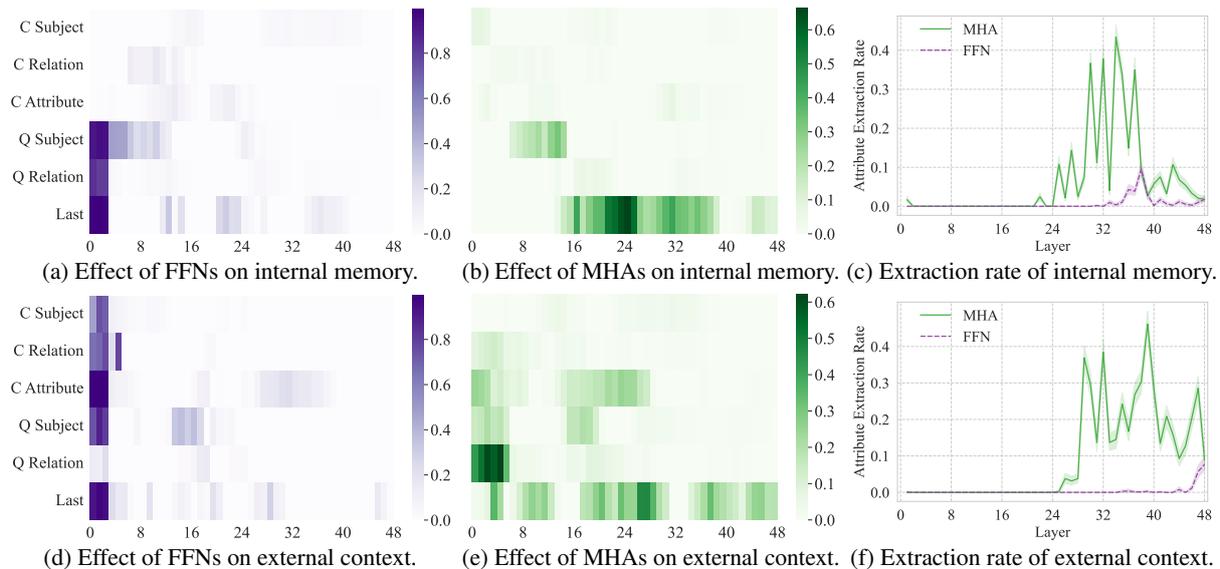


Figure 2: Effect of model components (FFNs and MHAs) in GPT-2 XL on the final prediction probability. Figures 2a and 2b (Figures 2d and 2e) show the effect of different model components and input elements when the model predicts based on internal memory (external context). The deeper color indicates the greater the impact of knocking out this part on the original prediction probability. Figure 2c (Figure 2f) shows the effect of MHAs and FFNs on the last token’s attribute extraction rate when the model predicts based on internal memory (external context).

#### 4.1 Examining Component Functionality

We start by exploring the functionality of model components (including FFNs and MHAs across various layers) in knowledge conflicts.

**Experiment 1: Knocking Out Component.** We examine which component in the transformer layer is critical for the attribute prediction by knocking out activations. Then, we divide the input into six elements for analysis: context subject  $s_c$ , context relation  $r_c$ , context attribute  $a_c$ , question subject  $s_q$ , question relation  $r_q$ , and the last token  $x_N$ . To measure the impact on the final prediction results, we zero-out the updates to the specified input element from the MHA and FFN blocks within each layer. For example, to intervene in the update of the  $\ell$ -th MHA (FFN) to the input element  $s_c$ , we set  $\mathbf{a}_i^{\ell'} = \mathbf{0}$  ( $\mathbf{m}_i^{\ell'} = \mathbf{0}$ ) for  $i$  in the token range of  $s_c$  and  $\ell' = \max(1, \ell - W/2), \dots, \min(L, \ell + W/2)$ , where  $W$  denotes the window size. We define the effect of a model component as the change in the original prediction probability after knocking it out.

**Results.** Figure 2 illustrates the effect of model components (FFNs and MHAs) in GPT-2 XL with the window size  $W = 5$ . Our observation reveals that destroying the FFN blocks in the early layers has a significant effect on the prediction probability while destroying the FFN blocks at the late layers shows minimal or no impact (Figures 2a and 2d). Moreover, the MHA blocks at the middle and late layers are crucial for the last token (Figures 2b and

2e). A possible explanation of the model’s behavior on the factual recall task is that *the early FFNs first enrich the semantic information of input elements, and then the enriched semantic information about attributes is extracted to the last token via late MHAs, where knowledge conflicts may arise at the later stage*. To verify this hypothesis, we will examine the attribute extraction function of MHAs.

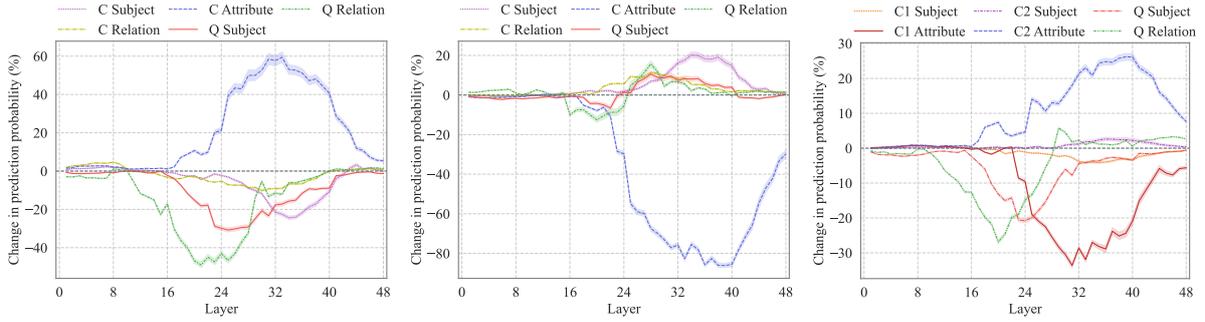
**Experiment 2: Extracting Attributes via MHAs.** We adopt the extraction rate (Geva et al., 2023) to examine the attribute extraction function of MHAs. We apply the early exit (Schuster et al., 2021; Geva et al., 2022) to project the MHA update  $\mathbf{a}_N^{\ell}$  for the last token  $x_N$  over the vocabulary. Then we check whether the top token  $t^{\ell}$  of each update aligns with the attribute  $t^*$  predicted at the final layer  $L$ :

$$t^* = \arg \max (\mathbf{p}_N^L), \quad (7)$$

$$t^{\ell} = \arg \max (\sigma (\phi (\mathbf{a}_N^{\ell}))). \quad (8)$$

We consider that the MHA correctly performs attribute extraction when  $t^* = t^{\ell}$ . For comparison, we also examine the extraction rate of FFNs.

**Results.** As illustrated in Figures 2c and 2f, it is evident that *the attribute extraction rate of MHAs significantly exceeds that of FFNs*. Moreover, attribute extraction mainly takes place at the 24-48 layers. Results for GPT-J show similar trends in Appendix C. The above findings motivate us to conduct an in-depth study on the information flows of MHAs from input elements to the last token.



(a) Prediction based on internal memory. (b) Prediction based on external context. (c) Prediction based on internal memory.

Figure 3: Relative change in the prediction probability when blocking the information flow from the input elements to the last token. Figures 3a and 3b only provide conflicting context. Figure 3c provides both supporting and conflicting context to internal memory, C1 denotes the supporting context, and C2 denotes the conflicting context.

## 4.2 Tracing Information Flow

The analysis presented above confirms that the last token extracts attribute information for prediction through MHA blocks. Following this, we explore the order and importance of the information flow from the various elements to the last token.

**Experiment 3: Blocking Information Flow.** We localize the information propagation from the input elements (including  $s_c$ ,  $r_c$ ,  $a_c$ ,  $s_q$  and  $r_q$ ) to the last token by knocking out attention edges between them. For example, to block the information flow from the input element  $s_c$  to the last token  $x_N$  in the layer  $\ell$ , we set the attention weight  $s^{\ell,h}[N, i] = 0$  for  $i$  in the token range of  $s_c$ ,  $h = 1, \dots, M$ , and  $\ell' = \max(1, \ell - W/2), \dots, \min(L, \ell + W/2)$ . In this way, we can restrict the last token from attending to the target element. If blocking the information propagation between them has a significant impact on the original prediction probability, this indicates that it is a crucial information flow.

**Results.** Figure 3 illustrates the information flow in GPT-2 XL with the window size  $W = 9$ . We can observe that in the early to middle layers, blocking the attention to the question relation leads to a decrease in the prediction probability. Similarly, in the subsequent layers, blocking the attention to the question subject also results in a decrease in the prediction probability. This suggests that the critical relation and subject information in the question are sequentially transmitted to the last token.

Then, in the middle to late layers, blocking the attention to the context subject and context attribute has the opposite effect on the final prediction probability. Taking Figure 3a as an example (when the model predicts the attribute based on internal memory), blocking the attention to the context attribute can improve the prediction probability, however, blocking the attention to the context subject can re-

duce the prediction probability. This suggests that the last token can extract the internal knowledge from the context subject, and extract the external knowledge from the context attribute. In addition, the last token also extracts a certain degree of internal knowledge from the question subject. Results for GPT-J show consistent trends in Appendix C.

Overall, this shows that there are two specific stages in the process of information flow passing to the last token: (1) *the question information is first passed to the last token*; (2) *the last token extracts or copies the attribute from the context subject or the context attribute*. In the later stage, knowledge conflicts arise during the process of merging inconsistent information flows from MHAs.

**Experiment 4: Extending to Conflicts between Contexts.** We extend our analysis to a more complex scenario in which the model is presented with both supporting context and conflicting context relative to internal memory. Supporting context and conflicting context contain  $a_m$  and  $a_c$  respectively:

C1: The capital of  $\{s\}$  is  $\{a_m\}$ .  
 C2: The capital of  $\{s\}$  is  $\{a_c\}$ .  $\{q\}$

We find that GPT-2 XL prefers to choose attributes consistent with internal memory 97.6% of the time. Hence, we only analyze the cases where the model makes predictions based on its internal memory.

**Results.** As illustrated in Figure 3c, we can observe that the question information is first passed to the last token in the first stage, which is consistent with the trend of a single conflicting context. In the second stage, a notable distinction is that *the model no longer extracts the memory attribute from the subject; instead, it opts for a more straightforward approach of copying the memory attribute from the context*. The above findings indicate that there exists a mechanism within MHAs capable of distinguishing and selecting between internal knowledge

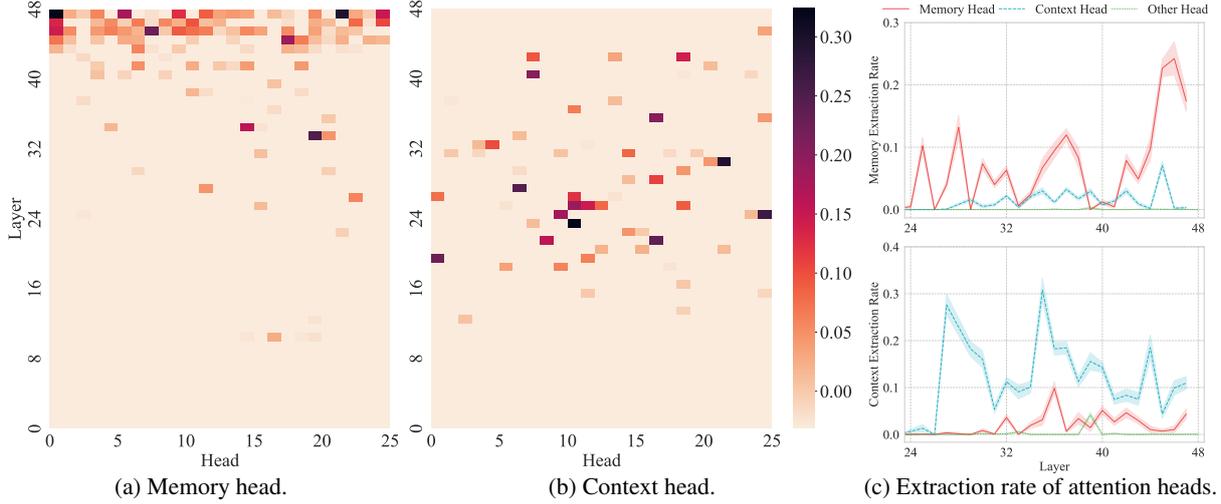


Figure 4: Memory heads and context heads in GPT-2 XL. Figure 4a shows the important score heatmap for predicting based on internal memory. Figure 4b shows the important score heatmap for predicting based on external context. Figure 4c illustrates the memory and context attribute extraction rate of different attention heads.

and external knowledge. This motivates us to conduct further analysis of MHAs.

### 4.3 Looking Deeper into Attention Heads

Attention heads serve as the fundamental component of an MHA block. For example, GPT-2 XL contains a total of 1,200 attention heads. This motivates us to conduct an investigation into the role of attention heads in handling knowledge conflicts.

#### Experiment 5: Discovering Important Heads.

To discover the attention heads that are crucial for predicting memory attributes or context attributes, we compute the gradient-based importance score (Michel et al., 2019; Bansal et al., 2023) for each head. Given a dataset  $\mathcal{D}$  with a set of inputs  $x$  and outputs  $y$ , the importance score of an attention head  $h$  captures the expected sensitivity of the model to  $h$  and is computed as follows:

$$I^{l,h}(\mathcal{D}) = \mathbb{E}_{(x,y)} \left| \mathbf{H}^{l,hT} \frac{\partial \mathcal{L}(y | x)}{\partial \mathbf{H}^{l,h}} \right|, \quad (9)$$

where  $\mathcal{L}(\cdot)$  is the loss function of conditional autoregressive generation. The proxy score of head  $h$  for predicting internal memory is calculated as:

$$S_m^{l,h}(\mathcal{D}_m, \mathcal{D}'_m) = I^{l,h}(\mathcal{D}_m) - I^{l,h}(\mathcal{D}'_m), \quad (10)$$

where  $(x, a_m) \in \mathcal{D}_m$  denotes the original outputs are memory attributes,  $(x, a_c) \in \mathcal{D}'_m$  denotes replacing the original outputs with context attributes. In this way, we can also calculate the proxy score of head  $h$  for predicting external context as:

$$S_c^{l,h}(\mathcal{D}_c, \mathcal{D}'_c) = I^{l,h}(\mathcal{D}_c) - I^{l,h}(\mathcal{D}'_c). \quad (11)$$

We compute the proxy score of each head across different layers to discover important heads.

**Results.** As shown in Figure 4a (Figure 4b), the deeper color of the red square indicates a more significant contribution from this attention head to the model’s predictions based on internal memory (external context). We can observe that there are a specific number of attention heads within middle-to-late layers that play opposite roles in predicting attributes. Accordingly, we refer to those heads that contribute to the prediction of memory attributes as **memory heads**, and those that facilitate predicting context attributes as **context heads**. Therefore, we claim that *they may serve in a mutually exclusive capacity during knowledge conflicts*. The heatmaps of GPT-J are provided in the Appendix C.

#### Experiment 6: Extracting Specific Attributes via Heads.

We further analyze the two types of heads discovered above to verify their role in knowledge conflicts. We rank the attention heads in descending order based on their importance scores,  $S_m^{l,h}$  for memory and  $S_c^{l,h}$  for context, subsequently identifying the top-5% of heads as memory heads and context heads, respectively. For comparison, we also randomly choose an additional 5% of the attention heads as other heads. Then, we examine their memory extraction rate when  $t_\ell = a_m$ , and context extraction rate when  $t_\ell = a_c$ .

**Results.** As shown in Figure 4c, memory heads and context heads are responsible for extracting different attribute information to the last token with a significant difference between memory and context extraction rates. Therefore, we discern that *the pivotal point at which knowledge conflicts emerge in LMs is the integration of inconsistent information flows by memory heads and context heads*.

## 5 Mitigating Knowledge Conflicts

Building on the above insights, we propose a novel method called **Pruning Head via Path Patching (PH3)** to efficiently mitigate knowledge conflicts by intervening on attention heads without the need to update model parameters (§5.1). Then, we conduct extensive experiments to show that our method can flexibly control LMs to use internal memory or external context (§5.2). Moreover, we analyze the generalization capability of our method (§5.3).

### 5.1 Method

Our method consists of two stages, first identifying the important heads through path patching, then intervening on these heads via structured pruning.

**Localizing Memory Heads and Context Heads via Path Patching.** When we use the gradient-based method in §4.3 to estimate the importance score of the target head  $h$ , it is subject to interference from other heads. The calculated gradients may not fully reflect the contribution of the target head, but rather a mixture of the influences from other heads. Therefore, we adopt the path patching technique (Goldowsky-Dill et al., 2023; Wang et al., 2023a) to analyze the causal relationship between the head  $h$  and the output attribute (including  $a_m$  and  $a_c$ ) in conflicts. To calculate the important score  $S_c^{\ell,h}$  of the target head  $h$ , our path patching method consists of three steps shown in Figure 14:

1. Run on the original input  $x \in \mathcal{D}_c$  to record the original activations of all heads;
2. Run on the corrupted input  $\mathcal{X}$  to record the corrupted activations of all heads, where  $\mathcal{X}$  is:

The capital of  $\{s\}$  is  $\langle \text{unk} \rangle$ .  $\{q\}$

where  $\langle \text{unk} \rangle$  is the special token;

3. Run on the original input  $x$ , while keeping all the heads frozen to their activations on  $x$ , except for the target head  $h$  whose activation is set on  $\mathcal{X}$ . Then measure the important score as the change of output logits.

The important score  $S_c^{\ell,h}$  of head  $h$  is computed as:

$$S_c^{\ell,h}(\mathcal{D}_c) = \mathbb{E}_{(x)}[(\mathbb{P}_x(a_c) - \mathbb{P}_x(a_m)) - (\mathbb{P}_{\mathcal{X}}(a_c) - \mathbb{P}_{\mathcal{X}}(a_m))]. \quad (12)$$

We adopt similar steps to calculate the importance score  $S_m^{\ell,h}$  of the target head  $h$  for memory attribute prediction in Appendix D. We also provide the importance score heatmaps of memory and context heads for various models in Appendix E, and our method can clearly distinguish between them.

**Pruning Attention Heads to Mitigate Knowledge Conflicts.** By ranking all the attention heads in ascending order based on the importance score  $S_c^{\ell,h}$  ( $S_m^{\ell,h}$ ), we can prune the top- $k\%$  attention heads that negatively impact the model’s capability to predict context (memory) attributes, thereby enhancing the model’s ability to utilize external context (internal memory). To prune a head  $h$  in layer  $\ell$  in practice, we set  $\mathbf{H}^{\ell,h}$  to be the zero matrix.

### 5.2 Experiment

**Setups.** We evaluate our method on five datasets, including World Capital, World Capital D, Official Language, Country, and Continent. To verify the generalization of PH3, we only calculate the importance scores of the attention heads on the World Capital dataset, and then directly evaluate PH3 on other datasets. We also select 1,000 test samples from an open-domain QA dataset NQ (Kwiatkowski et al., 2019), providing the LM with the top-5 retrieved passages, and ensuring that at least one relevant passage is among them. We validate the effectiveness of PH3 on eight LMs.

**Metrics.** We use the internal memory usage rate  $RM = \frac{f_m}{f_m+f_c+f_o}$  and the external context usage rate  $RC = \frac{f_c}{f_m+f_c+f_o}$  to assess how effectively the method controls the reliance of LMs on either internal memory or external context, where  $f_m$  is the frequency of relying on internal memory,  $f_c$  is the frequency of relying on external context, and  $f_o$  is the frequency of other answers. For the open-domain QA task, we use Recall to evaluate whether the model can provide correct answers based on the retrieved passages following Adlakha et al. (2023).

**Baselines.** We compare with the following baselines: (1) **Prompt:** We instruct the LM to generate answers based on internal memory or external context through specific prompts; (2) **CAD:** Shi et al. (2023) leverage contrastive decoding (Li et al., 2023b) to encourage the LM to attend to its context during generation; (3) **Gradient:** We replace our path patching method with the gradient-based method to discover the attention heads. We select the optimal pruning rate  $k$  on the development set for both Gradient and PH3. More details about hyperparameter settings are in Appendix B.2.

**Results.** Table 1 shows the results of GPT-2 XL, GPT-J and LLaMA2-7B, and more results of other models are in Table 2. Throughout our experiments, we note the following key observations:

Model	Method	World Capital		World Capital D		Official Language		Country		Continent		
		RM	RC	RM	RC	RM	RC	RM	RC	RM	RC	
GPT-2 XL	↑ Memory	Base	59.2	40.8	47.2	52.8	42.2	57.8	37.2	62.8	41.5	58.5
		Prompt	12.5	81.2	21.3	71.2	20.3	74.4	24.5	75.5	16.2	43.4
		Gradient	72.4	9.8	78.6	10.5	41.5	40.5	39.1	60.2	42.7	46.6
		PH3 (Ours)	<b>97.9</b>	<b>0.6</b>	<b>93.3</b>	<b>2.5</b>	<b>74.4</b>	<b>9.8</b>	<b>50.9</b>	<b>36.3</b>	<b>53.1</b>	<b>38.1</b>
	↑ Context	Prompt	9.3	87.5	18.9	75.2	17.1	80.7	18.5	81.4	25.5	58.3
		CAD	25.0	65.6	12.5	63.6	<b>9.1</b>	80.5	27.2	72.5	22.9	60.4
		Gradient	44.4	49.0	28.0	58.7	29.6	59.5	36.4	63.4	<b>18.4</b>	51.5
		PH3 (Ours)	27.5	68.9	7.7	91.3	20.7	74.8	22.7	75.7	27.7	<b>66.7</b>
		+ Prompt	<b>3.6</b>	<b>95.1</b>	<b>5.2</b>	<b>94.4</b>	9.6	<b>88.7</b>	<b>12.6</b>	<b>86.0</b>	20.9	63.8
	GPT-J	↑ Memory	Base	37.5	62.5	43.1	56.9	41.5	58.5	54.0	46.0	43.2
Prompt			29.8	67.1	31.6	62.1	23.1	69.1	22.4	77.6	12.5	86.0
Gradient			67.9	8.3	67.6	<b>6.7</b>	39.4	53.4	54.4	45.5	57.2	30.0
PH3 (Ours)			<b>93.3</b>	<b>1.6</b>	<b>76.5</b>	10.8	<b>63.3</b>	<b>25.3</b>	<b>58.9</b>	<b>40.5</b>	<b>75.1</b>	<b>17.6</b>
↑ Context		Prompt	31.9	64.5	15.8	76.4	16.2	70.6	17.9	82.1	7.2	<b>91.4</b>
		CAD	2.5	89.9	13.4	68.2	4.7	89.9	17.0	81.8	13.0	80.3
		Gradient	6.1	88.3	7.9	67.8	5.7	76.4	29.2	70.5	36.8	60.7
		PH3 (Ours)	0.2	99.3	<b>0.1</b>	<b>98.4</b>	2.3	<b>90.6</b>	9.5	86.7	8.0	64.9
		+ Prompt	<b>0.1</b>	<b>99.5</b>	0.2	97.8	<b>2.0</b>	81.9	<b>1.4</b>	<b>98.6</b>	<b>1.4</b>	90.9
			Base	46.3	53.7	95.5	4.0	18.8	80.3	52.9	46.8	30.9
↑ Memory	Prompt	36.0	63.2	96.0	3.7	40.0	59.1	68.2	31.6	77.4	22.6	
	Gradient	81.0	5.8	95.1	1.6	50.1	47.4	60.0	38.0	64.5	24.5	
	PH3 (Ours)	<b>98.1</b>	<b>1.2</b>	<b>98.0</b>	<b>1.3</b>	<b>73.7</b>	<b>17.8</b>	<b>76.9</b>	<b>20.6</b>	<b>90.5</b>	<b>8.8</b>	
↑ Context	Prompt	3.2	96.6	92.4	2.2	25.5	73.8	58.2	41.5	19.2	80.3	
	CAD	1.4	95.5	29.1	70.6	<b>0.0</b>	<b>100.0</b>	13.6	86.1	0.2	98.2	
	Gradient	23.6	63.2	40.1	58.8	25.7	74.6	17.6	82.2	27.1	72.9	
	PH3 (Ours)	1.6	97.4	19.1	73.4	0.1	99.9	5.2	94.7	0.5	99.4	
	+ Prompt	<b>0.4</b>	<b>98.8</b>	<b>10.6</b>	<b>85.3</b>	<b>0.0</b>	<b>100.0</b>	<b>2.8</b>	<b>97.0</b>	<b>0.0</b>	<b>100.0</b>	

Table 1: Experimental results of GPT-2 XL, GPT-J and LLaMA2-7B on five datasets. Bolds denote the best results.

(1) PH3 significantly outperforms other baselines. Experimental results show that PH3 can not only increase the average internal memory usage rate of eight LMs by 44.0%, but also increase the average external context usage rate by 38.5%. When PH3 is combined with Prompt, it can more effectively control the LMs to use external context.

(2) As shown in Table 3, PH3 can also achieve an average 6.2% Recall improvement on open-domain QA tasks. By pruning a small number of negative context heads, PH3 can make LMs generate answers more faithfully based on retrieved passages.

(3) Although Prompt and CAD can effectively increase the external context usage rate, there are limitations. CAD cannot directly enhance internal memory, and Prompt may even have the opposite effect. In contrast, our method offers a viable solution to enhance the internal memory usage rate.

### 5.3 Analysis

We conduct a thorough analysis of the generalization ability of PH3. For cross-model generalization, PH3 is effective across a wide range of models. This shows that our method is not limited to small models, but can also be adopted on relatively large models, including the popular LLaMA2 series. For cross-relation generalization, by intervening on the attention heads discovered on World Capital, our

method can also well resolve knowledge conflicts on other relation types. This indicates that PH3 does not identify attention heads specific to a certain type of relation. Instead, it identifies universal memory and context heads. For cross-format generalization, PH3 can transfer well from triple-form context to document-form context. This indicates that our method does not merely remember the relative positions of elements in context, but is capable of understanding the external context. Compared to the Gradient, our method has demonstrated superior generalizability. We also analyze the impact of the number of pruning heads in Appendix G.

## 6 Conclusion

In this paper, we perform an exploration into the mechanism of interpreting knowledge conflicts and reveal that memory and context heads in later layers can cause knowledge conflicts when merging inconsistent information flows. Based on our insights, we propose a novel method called **Pruning Head via PatH PatcHing (PH3)**, which can mitigate knowledge conflicts by pruning those conflicting attention heads. We prove that PH3 can flexibly control LMs to use internal memory or external context. We also demonstrate the cross-model, cross-relation, and cross-format generalization.

## 590 Limitations

591 For further study, we conclude some limitations of  
592 our work as follows:

- 593 • Similar to previous works on mechanism inter-  
594 preteability that adopt tasks such as antonym  
595 generation (Todd et al., 2023), fact recall  
596 (Meng et al., 2022; Geva et al., 2023), arith-  
597 metic operation (Hanna et al., 2023; Stolfo  
598 et al., 2023), and text classification (Bansal  
599 et al., 2023; Wang et al., 2023b), our work also  
600 selects a relatively simpler task to interpret the  
601 mechanism behind knowledge conflicts. Sim-  
602 ple tasks enable us to better control variables  
603 and minimize external distractions. In the fu-  
604 ture, we plan to extend our analysis to more  
605 complex and realistic scenarios, such as where  
606 irrelevant information is present within the ex-  
607 ternal context, or where the model needs to  
608 reason with both internal and external knowl-  
609 edge.  
610 • Although our research has delved into the at-  
611 tention heads in LMs, there may be more ba-  
612 sic elements involved in knowledge conflicts.  
613 Furthermore, the memory and context heads  
614 we have discovered may not only be respon-  
615 sible for extracting knowledge from internal  
616 memory or external context. These heads may  
617 also have other functions, such as helping the  
618 model capture global dependencies of input  
619 texts. By pruning these heads, the original ca-  
620 pabilities of the model may be affected. There-  
621 fore, we will further explore mitigating knowl-  
622 edge conflicts through more subtle interven-  
623 tion methods.

624 In summary, the mechanism behind knowledge  
625 conflicts remains a largely unexplored area, and we  
626 hope our work can offer some useful insights for  
627 further research.

## 628 Ethics Statement

629 To enhance the reproducibility of our research, we  
630 will make all source code and datasets publicly  
631 available upon the acceptance of this paper. Our  
632 work focuses on uncovering the mechanisms be-  
633 hind knowledge conflicts in LM, thereby better  
634 controlling the model in retrieval augmentation and  
635 tool augmentation. Through effective intervention,  
636 our method can make the LM more controllable  
637 and trustworthy. On the one hand, it can prevent  
638 prompt injections from attacking the model, and  
639 on the other hand, it can correct the biased knowl-

edge that the model learned during pre-training.  
640 Nonetheless, the impact of head pruning on the  
641 model’s original capabilities remains unexplored.  
642 These factors should be taken into careful consid-  
643 eration for future research.  
644

## References 645

- Vaibhav Adlakha, Parishad BehnamGhader, Xing Han  
646 Lu, Nicholas Meade, and Siva Reddy. 2023. Eval-  
647 uating correctness and faithfulness of instruction-  
648 following models for question answering. *arXiv*  
649 *preprint arXiv:2307.16877*. 650
- Hritik Bansal, Karthik Gopalakrishnan, Saket Dingliwal,  
651 Sravan Bodapati, Katrin Kirchhoff, and Dan Roth.  
652 2023. [Rethinking the role of scale for in-context](#)  
653 [learning: An interpretability-based case study at 66](#)  
654 [billion scale](#). In *Proceedings of the 61st Annual Meet-*  
655 *ing of the Association for Computational Linguis-*  
656 *tics (Volume 1: Long Papers)*, pages 11833–11856,  
657 Toronto, Canada. Association for Computational Lin-  
658 guistics. 659
- Stella Biderman, Hailey Schoelkopf, Quentin Gregory  
660 Anthony, Herbie Bradley, Kyle O’Brien, Eric Hal-  
661 lahan, Mohammad Aflah Khan, Shivanshu Purohit,  
662 USVSN Sai Prashanth, Edward Raff, et al. 2023.  
663 Pythia: A suite for analyzing large language mod-  
664 els across training and scaling. In *International*  
665 *Conference on Machine Learning*, pages 2397–2430.  
666 PMLR. 667
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie  
668 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind  
669 Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
670 Askell, Sandhini Agarwal, Ariel Herbert-Voss,  
671 Gretchen Krueger, Tom Henighan, Rewon Child,  
672 Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens  
673 Winter, Chris Hesse, Mark Chen, Eric Sigler, Jack  
674 Teusz Litwin, Scott Gray, Benjamin Chess, Jack  
675 Clark, Christopher Berner, Sam McCandlish, Alec  
676 Radford, Ilya Sutskever, and Dario Amodei. 2020.  
677 [Language models are few-shot learners](#). In *Ad-*  
678 *vances in Neural Information Processing Systems*,  
679 volume 33, pages 1877–1901. Curran Associates, Inc.  
680 681
- Nick Cammarata, Shan Carter, Gabriel Goh, Chris Olah,  
682 Michael Petrov, Ludwig Schubert, Chelsea Voss, Ben  
683 Egan, and Swee Kiat Lim. 2020. [Thread: Circuits](#).  
684 *Distill*. <https://distill.pub/2020/circuits>. 685
- Hung-Ting Chen, Michael Zhang, and Eunsol Choi.  
686 2022. [Rich knowledge sources bring complex knowl-](#)  
687 [edge conflicts: Recalibrating models to reflect con-](#)  
688 [flicting evidence](#). In *Proceedings of the 2022 Con-*  
689 *ference on Empirical Methods in Natural Language*  
690 *Processing*, pages 2292–2307, Abu Dhabi, United  
691 Arab Emirates. Association for Computational Lin-  
692 guistics. 693
- Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao  
694 Chang, and Furu Wei. 2022. [Knowledge neurons in](#)  
695

696	<a href="#">pretrained transformers</a> . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8493–8502, Dublin, Ireland. Association for Computational Linguistics.	752
697		753
698		754
699		755
700		756
701	Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, et al. 2021. A mathematical framework for transformer circuits. <i>Transformer Circuits Thread</i> , 1.	757
702		758
703		
704		
705		
706	Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. 2023. <a href="#">Dissecting recall of factual associations in auto-regressive language models</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 12216–12235, Singapore. Association for Computational Linguistics.	759
707		760
708		761
709		762
710		763
711		764
712		765
713	Mor Geva, Avi Caciularu, Kevin Wang, and Yoav Goldberg. 2022. <a href="#">Transformer feed-forward layers build predictions by promoting concepts in the vocabulary space</a> . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 30–45, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	766
714		767
715		768
716		769
717		770
718		771
719		772
720		773
721	Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2021. <a href="#">Transformer feed-forward layers are key-value memories</a> . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 5484–5495, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	774
722		775
723		776
724		777
725		778
726		779
727		780
728	Nicholas Goldowsky-Dill, Chris MacLeod, Lucas Sato, and Aryaman Arora. 2023. Localizing model behavior with path patching. <i>arXiv preprint arXiv:2304.05969</i> .	781
729		782
730		783
731		784
732	Danny Halawi, Jean-Stanislas Denain, and Jacob Steinhardt. 2023. Overthinking the truth: Understanding how language models process false demonstrations. <i>arXiv preprint arXiv:2307.09476</i> .	785
733		786
734		787
735		788
736	Michael Hanna, Ollie Liu, and Alexandre Variengien. 2023. <a href="#">How does GPT-2 compute greater-than?: Interpreting mathematical abilities in a pre-trained language model</a> . In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	789
737		790
738		791
739		792
740	Yaru Hao, Li Dong, Furu Wei, and Ke Xu. 2021. <a href="#">Self-attention attribution: Interpreting information interactions inside transformer</a> . <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 35(14):12963–12971.	793
741		794
742		795
743		796
744		797
745		798
746	Roei Hendel, Mor Geva, and Amir Globerson. 2023. <a href="#">In-context learning creates task vectors</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 9318–9333, Singapore. Association for Computational Linguistics.	799
747		800
748		801
749		802
750		803
751	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. <a href="#">Natural Questions: A Benchmark for Question Answering Research</a> . <i>Transactions of the Association for Computational Linguistics</i> , 7:453–466.	804
		805
		806
		807
	Daliang Li, Ankit Singh Rawat, Manzil Zaheer, Xin Wang, Michal Lukasik, Andreas Veit, Felix Yu, and Sanjiv Kumar. 2023a. <a href="#">Large language models with controllable working memory</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 1774–1793, Toronto, Canada. Association for Computational Linguistics.	808
		809
		810
		811
		812
		813
		814
		815
		816
		817
		818
		819
		820
		821
		822
		823
		824
		825
		826
		827
		828
		829
		830
		831
		832
		833
		834
		835
		836
		837
		838
		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

808	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	Kevin Ro Wang, Alexandre Variengien, Arthur Conmy,	865
809	Dario Amodei, Ilya Sutskever, et al. 2019. Language	Buck Shlegeris, and Jacob Steinhardt. 2023a. <a href="#">Inter-</a>	866
810	models are unsupervised multitask learners. <i>OpenAI</i>	<a href="#">pretability in the wild: a circuit for indirect object</a>	867
811	<i>blog</i> , 1(8):9.	<a href="#">identification in GPT-2 small</a> . In <i>The Eleventh Inter-</i>	868
		<i>national Conference on Learning Representations</i> .	869
812	Mansi Sakarvadia, Aswathy Ajith, Arham Khan, Daniel	Lean Wang, Lei Li, Damai Dai, Deli Chen, Hao Zhou,	870
813	Grzenda, Nathaniel Hudson, André Bauer, Kyle	Fandong Meng, Jie Zhou, and Xu Sun. 2023b. <a href="#">Label</a>	871
814	Chard, and Ian Foster. 2023a. <a href="#">Memory injections:</a>	<a href="#">words are anchors: An information flow perspective</a>	872
815	<a href="#">Correcting multi-hop reasoning failures during infer-</a>	<a href="#">for understanding in-context learning</a> . In <i>Proceed-</i>	873
816	<a href="#">ence in transformer-based language models</a> . In	<i>ings of the 2023 Conference on Empirical Methods</i>	874
817	<i>Proceedings of the 6th BlackboxNLP Workshop: An-</i>	<i>in Natural Language Processing</i> , pages 9840–9855,	875
818	<i>alyzing and Interpreting Neural Networks for NLP</i> ,	Singapore. Association for Computational Linguistics.	876
819	pages 342–356, Singapore. Association for Compu-		877
820	tational Linguistics.		
821	Mansi Sakarvadia, Arham Khan, Aswathy Ajith, Daniel	Yike Wang, Shangbin Feng, Heng Wang, Weijia	878
822	Grzenda, Nathaniel Hudson, André Bauer, Kyle	Shi, Vidhisha Balachandran, Tianxing He, and Yu-	879
823	Chard, and Ian Foster. 2023b. <a href="#">Attention lens: A tool</a>	lia Tsvetkov. 2023c. <a href="#">Resolving knowledge con-</a>	880
824	<a href="#">for mechanistically interpreting the attention head</a>	<a href="#">flicts in large language models</a> . <i>arXiv preprint</i>	881
825	<a href="#">information retrieval mechanism</a> .	<i>arXiv:2310.00935</i> .	882
826	Tal Schuster, Adam Fisch, Tommi Jaakkola, and Regina	Jian Xie, Kai Zhang, Jiangjie Chen, Renze Lou, and	883
827	Barzilay. 2021. <a href="#">Consistent accelerated inference via</a>	Yu Su. 2023. <a href="#">Adaptive chameleon or stubborn</a>	884
828	<a href="#">confident adaptive transformers</a> . In <i>Proceedings of</i>	<a href="#">sloth: Unraveling the behavior of large language</a>	885
829	<i>the 2021 Conference on Empirical Methods in Natu-</i>	<a href="#">models in knowledge conflicts</a> . <i>arXiv preprint</i>	886
830	<i>ral Language Processing</i> , pages 4962–4979, Online	<i>arXiv:2305.13300</i> .	887
831	and Punta Cana, Dominican Republic. Association		
832	for Computational Linguistics.	Yi Yang, Hanyu Duan, Ahmed Abbasi, John P. Lalor,	888
833	Weijia Shi, Xiaochuang Han, Mike Lewis, Yulia	and Kar Yan Tam. 2023. <a href="#">Bias a-head? analyzing</a>	889
834	Tsvetkov, Luke Zettlemoyer, and Scott Wen-tau	<a href="#">bias in transformer-based language model attention</a>	890
835	Yih. 2023. <a href="#">Trusting your evidence: Hallucinate</a>	<a href="#">heads</a> .	891
836	<a href="#">less with context-aware decoding</a> . <i>arXiv preprint</i>	Qinan Yu, Jack Merullo, and Ellie Pavlick. 2023. <a href="#">Char-</a>	892
837	<i>arXiv:2305.14739</i> .	<a href="#">acterizing mechanisms for factual recall in language</a>	893
838	Alessandro Stolfo, Yonatan Belinkov, and Mrinmaya	<a href="#">models</a> . In <i>Proceedings of the 2023 Conference on</i>	894
839	Sachan. 2023. <a href="#">A mechanistic interpretation of arith-</a>	<i>Empirical Methods in Natural Language Processing</i> ,	895
840	<a href="#">metic reasoning in language models using causal</a>	pages 9924–9959, Singapore. Association for Com-	896
841	<a href="#">mediation analysis</a> . In <i>Proceedings of the 2023 Con-</i>	putational Linguistics.	897
842	<i>ference on Empirical Methods in Natural Language</i>	Fred Zhang and Neel Nanda. 2024. <a href="#">Towards best prac-</a>	898
843	<i>Processing</i> , pages 7035–7052, Singapore. Associa-	<a href="#">tices of activation patching in language models: Met-</a>	899
844	tion for Computational Linguistics.	<a href="#">rics and methods</a> .	900
845	Eric Todd, Millicent L. Li, Arnab Sen Sharma, Aaron	Susan Zhang, Stephen Roller, Naman Goyal, Mikel	901
846	Mueller, Byron C. Wallace, and David Bau. 2023.	Artetxe, Moya Chen, Shuohui Chen, Christopher De-	902
847	<a href="#">Function vectors in large language models</a> .	wan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mi-	903
848	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	haylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel	904
849	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	Simig, Punit Singh Koura, Anjali Sridhar, Tianlu	905
850	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	Wang, and Luke Zettlemoyer. 2022. <a href="#">Opt: Open pre-</a>	906
851	Bhosale, et al. 2023. <a href="#">Llama 2: Open founda-</a>	<a href="#">trained transformer language models</a> .	907
852	<a href="#">tion and fine-tuned chat models</a> . <i>arXiv preprint</i>	Wenxuan Zhou, Sheng Zhang, Hoifung Poon, and	908
853	<i>arXiv:2307.09288</i> .	Muhao Chen. 2023. <a href="#">Context-faithful prompting</a>	909
854	Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov,	<a href="#">for large language models</a> . In <i>Findings of the As-</i>	910
855	Sharon Qian, Daniel Nevo, Yaron Singer, and Stu-	<i>sociation for Computational Linguistics: EMNLP</i>	911
856	art Shieber. 2020. <a href="#">Investigating gender bias in lan-</a>	<i>2023</i> , pages 14544–14556, Singapore. Association	912
857	<a href="#">guage models using causal mediation analysis</a> . In	for Computational Linguistics.	913
858	<i>Advances in Neural Information Processing Systems</i> ,		
859	volume 33, pages 12388–12401. Curran Associates,		
860	Inc.		
861	Ben Wang and Aran Komatsuzaki. 2021. <a href="#">GPT-J-</a>		
862	<a href="#">6B: A 6 Billion Parameter Autoregressive Lan-</a>		
863	<a href="#">guage Model</a> . <a href="https://github.com/kingoflolz/">https://github.com/kingoflolz/</a>		
864	<a href="#">mesh-transformer-jax</a> .		

## A Related Work

### A.1 Investigating Knowledge Conflict

Previous research (Longpre et al., 2021; Chen et al., 2022; Yu et al., 2023; Xie et al., 2023; Wang et al., 2023c; Neeman et al., 2023) on knowledge conflicts primarily seek to answer the question: *do language models prefer internal memory or external context?* Yu et al. (2023) find that language models are more inclined to internal memory as the frequency of a fact in the pre-training corpus increases. Xie et al. (2023) demonstrate that large language models (LLMs) are highly receptive to external conflicting evidence. They also reveal that when both supportive and contradictory evidence to their internal memory are present, LLMs show a strong confirmation bias and tend to cling to their parametric memory. The above observed phenomena contribute to a better understanding of knowledge conflicts. However, the underlying mechanism of knowledge conflicts remains unclear. We observe that knowledge conflicts arise when the late attention heads integrate different information flows from internal memory and external context.

### A.2 Resolving Knowledge Conflict

Existing work (Shi et al., 2023; Zhou et al., 2023; Li et al., 2023a; Yu et al., 2023; Qian et al., 2023) has conducted preliminary exploration into the mitigation of knowledge conflicts. Shi et al. (2023) propose a simple method to encourage the LM to attend to the external context via contrastive decoding (Li et al., 2023b). Yu et al. (2023) use head attribution to identify individual attention heads that either promote the memorized answer or the in-context answer, then scale the value vector of these heads to increase the rate of the in-context answers. Our work is inspired by their exploration of attention heads, and we propose further analysis to improve understanding of the way knowledge conflicts are formed. Furthermore, while most existing methods (Shi et al., 2023; Yu et al., 2023) primarily focus on improving the model’s faithfulness to the context, enabling the model to adhere to its internal memory remains a challenging task.

### A.3 Mechanistic Interpretability

Recently, there has been a growing interest in the mechanistic interpretability (Cammarrata et al., 2020; Elhage et al., 2021) of parametric knowledge in LLMs, with efforts focusing on reverse engineering the computational processes of model parame-

ters. Dai et al. (2022) use a knowledge attribution method (Hao et al., 2021) to identify the knowledge neurons in FFNs. Meng et al. (2022) reveal that FFNs at a range of middle layers can recall facts by using the causal mediation analysis method (Vig et al., 2020). Geva et al. (2023) find that knowledge extraction is typically done via attention heads. Besides, there are some works investigating LLMs in mathematical reasoning (Hanna et al., 2023; Stolfo et al., 2023) and in-context learning (Hendel et al., 2023; Olsson et al., 2022; Bansal et al., 2023). Besides, there are some studies (Yang et al., 2023; Sakarvadia et al., 2023a,b; Zhang and Nanda, 2024) focused on interpreting attention heads in LLMs. Our work is highly inspired by previous wisdom in mechanistic interpretability, focusing on interpreting and mitigating knowledge conflicts in LLMs.

## B Implement Details

### B.1 Datasets

We construct Official Language, Country, and Continent datasets by sampling knowledge triples from Wikidata. The Official Language dataset requires the LM to predict the official language of the given city or country:

*The official language of {s} is {a<sub>c</sub>}.*  
Q: What is the official language of {s}? A:

The Country dataset requires the LM to predict the country to which the given city belongs:

*The city {s} is located in {a<sub>c</sub>}.*  
Q: Which country is the city {s} in? A:

The Continent dataset requires the LM to predict the continent on which the given country is located:

*{s} is in the continent of {a<sub>c</sub>}.*  
Q: Which continent is {s} located in? A:

We also generate a more complex World Capital D dataset based on the World Capital dataset, using gpt-3.5-turbo to rewrite the external context from triplet form into document form.

### B.2 Hyperparameter Settings

Our implementation is based on HuggingFace’s Transformers<sup>1</sup>, PyTorch<sup>2</sup> and baukit<sup>3</sup>. For the Prompt method, we use the following prompt to enhance the internal memory:

*Please answer the question based on your internal memory, ignoring the given context.*

<sup>1</sup><https://github.com/huggingface/transformers/>

<sup>2</sup><https://github.com/pytorch/pytorch/>

<sup>3</sup><https://github.com/davidbau/baukit/>

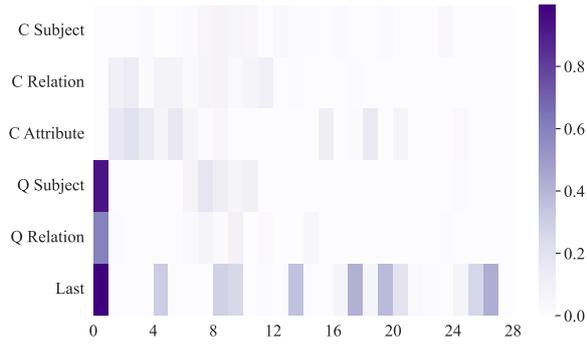


Figure 5: Effect of FFNs in GPT-J on internal memory.

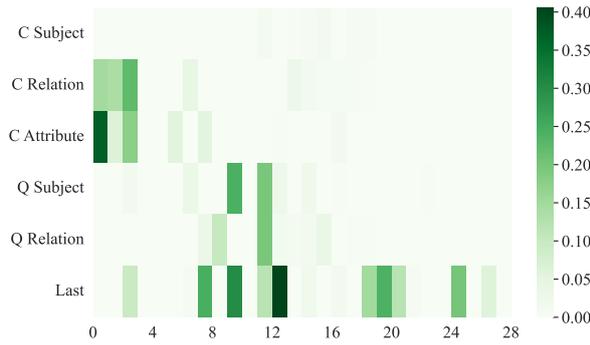


Figure 6: Effect of MHAs in GPT-J on internal memory.

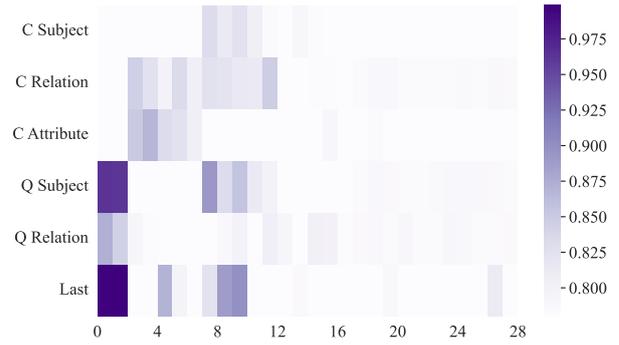


Figure 7: Effect of FFNs in GPT-J on external context.

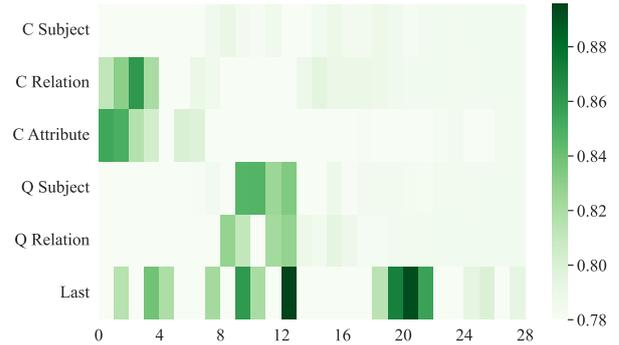


Figure 8: Effect of MHAs in GPT-J on external context.

and we use the following prompt to enhance the external context:

*Please answer the question based on the given context, ignoring your internal memory.*

For Gradient and PH3, we select the optimal pruning rate  $k \in \{1, 3, 5, 7, 9, 15\}$  on the development set with 200 samples. To mitigate knowledge conflicts, setting the pruning rate  $k$  of PH3 to 5 usually achieves excellent results. For enhancing the open-domain QA capabilities, we usually set the pruning rate  $k$  of PH3 to 3. Details about the models used in this paper are in Table 4. All experiments are conducted with NVIDIA GeForce RTX A6000 GPUs.

### C Additional Results for GPT-J

We provide here additional results for GPT-J. Figures 5 and 6 show the effect of FFNs and MHAs on internal memory, and Figures 7 and 8 show the effect of FFNs and MHAs on external context. Figures 9 and 10 illustrate the information flow in GPT-J with the window size  $W = 9$ . Figure 11 shows the information flow in GPT-J when providing both supporting context and conflicting context relative to internal memory. Figures 12 and 13 show the gradient-based important scores of memory heads and context heads in GPT-J.

### D Method Details

To calculate the important score  $S_m^{\ell,h}$  of the target head  $h$ , our path patching method consists of the following three steps:

1. Run on the original input  $x \in \mathcal{D}_m$  to record the original activations of all heads;
2. Run on the corrupted input  $\mathcal{X}$  to record the corrupted activations of all heads, where  $\mathcal{X}$  is:

*The capital of  $\langle \text{unk} \rangle$  is  $\{a_c\}$ .  
Q: What is the capital of  $\langle \text{unk} \rangle$  ? A:*

where  $\langle \text{unk} \rangle$  is the special token;

3. Run on the original input  $x$ , while keeping all the heads frozen to their activations on  $x$ , except for the target head  $h$  whose activation is set on  $\mathcal{X}$ . Then measure the important score as the change of output logits.

The important score  $S_m^{\ell,h}$  of head  $h$  is computed as:

$$S_m^{\ell,h}(\mathcal{D}_m) = \mathbb{E}_{(x)} [(\mathbb{P}_x(a_m) - \mathbb{P}_x(a_c)) - (\mathbb{P}_{\mathcal{X}}(a_m) - \mathbb{P}_{\mathcal{X}}(a_c))]. \quad (13)$$

### E Heatmaps of Attention Heads

We calculate the important scores of memory heads and context heads via our path patching method,

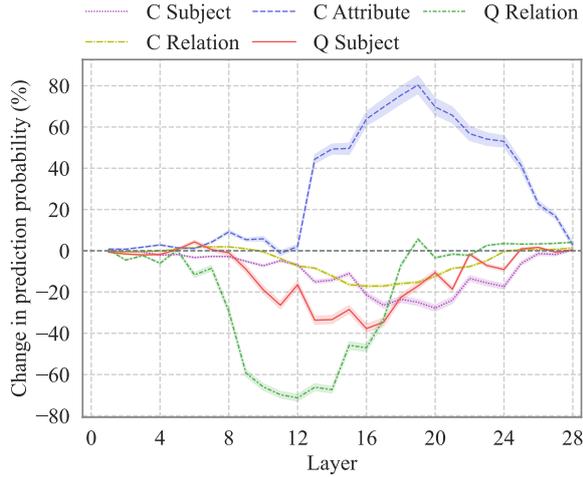


Figure 9: Relative change in the GPT-J's prediction probability based on internal memory.

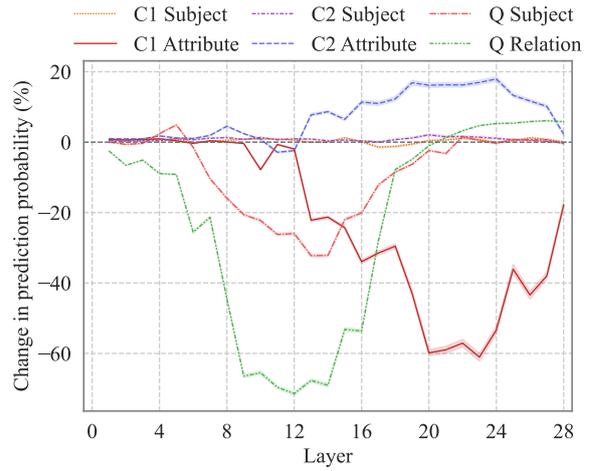


Figure 11: Relative change in the GPT-J's prediction probability based on internal memory when providing both supporting context and conflicting context.

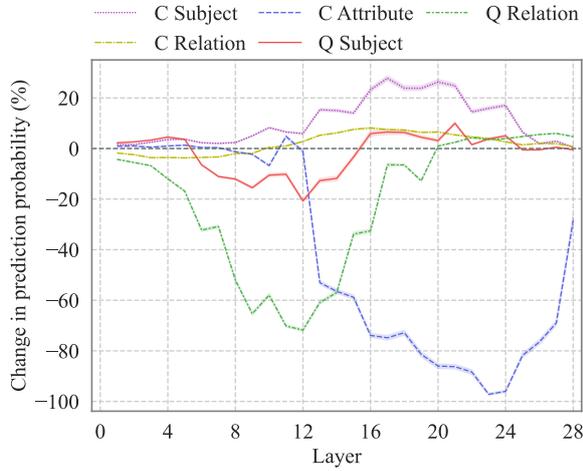


Figure 10: Relative change in the GPT-J's prediction probability based on external context.

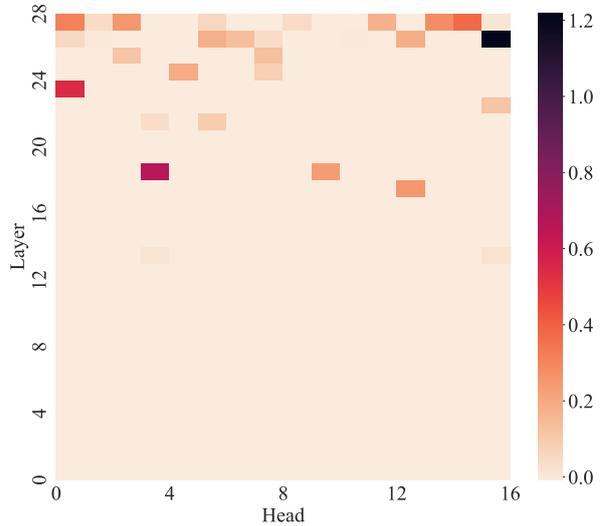


Figure 12: Memory Heads of GPT-J.

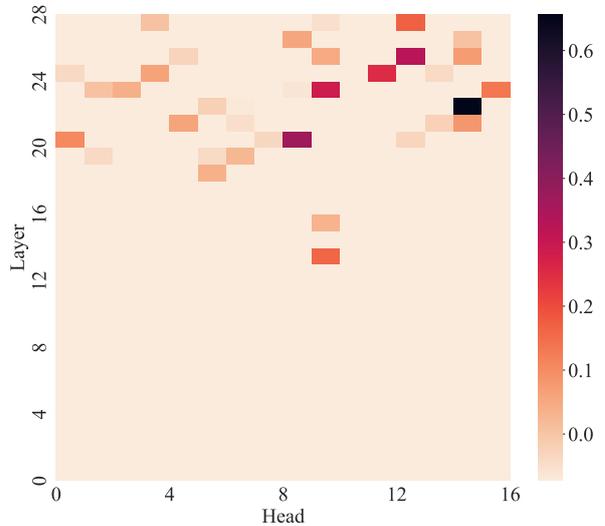


Figure 13: Context Heads of GPT-J.

1048 then provide the heatmaps for GPT-2 XL (Figures  
 1049 15 and 16), GPT-J (Figures 17 and 18), OPT-1.3B  
 1050 (Figures 19 and 20), OPT-2.7B (Figures 21 and  
 1051 22), Pythia-6.9B (Figures 23 and 24), Pythia-12B  
 1052 (Figures 25 and 26), LLaMA2-7B (Figures 27 and  
 1053 28) and LLaMA2-13B (Figures 29 and 30). The  
 1054 red squares indicate heads that have a significant  
 1055 positive impact, while the blue squares represent  
 1056 heads that have a negative effect.

## 1057 F Additional Experimental Results

1058 We report experimental results in Table 2 and 3.

## 1059 G Number of Pruning Heads

1060 As shown in Figures 31, 32, 33, 34, 35, 36, 37,  
 1061 38, 39, 40, 41 and 42, we analyze the impact of  
 1062 the number of pruning heads (sparsity ratio) on the  
 1063 Gradient and PH3 methods.

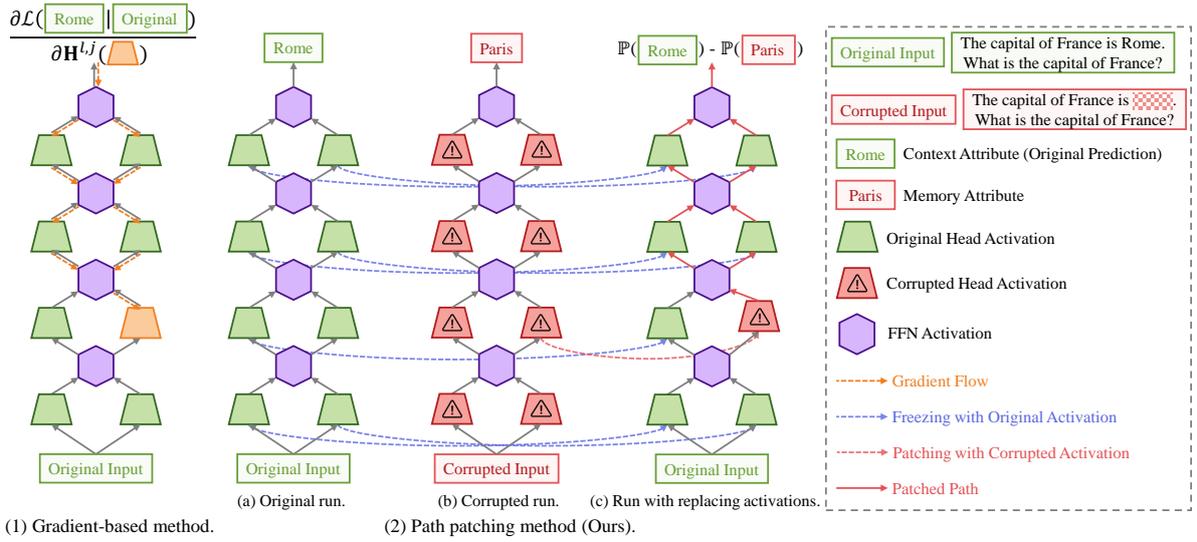


Figure 14: Illustration of gradient-based method and our path patching method.

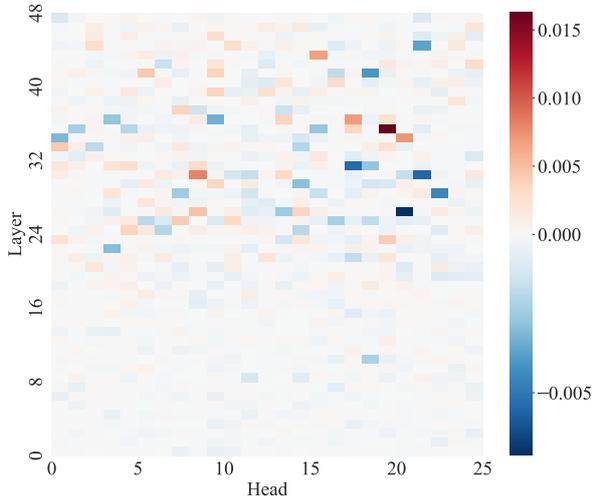


Figure 15: Memory Heads of GPT-2 XL.

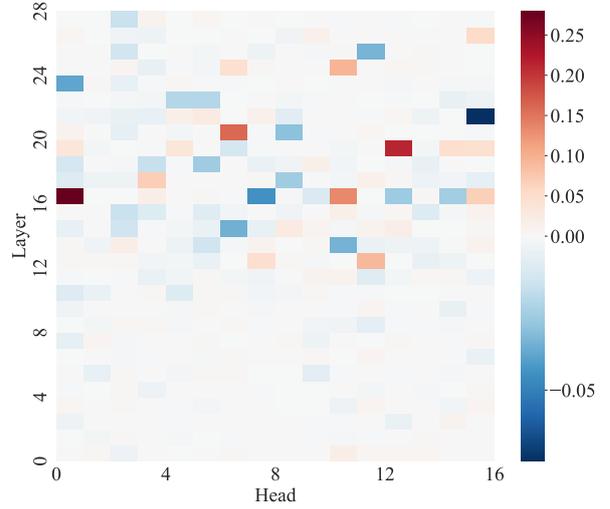


Figure 17: Memory Heads of GPT-J.

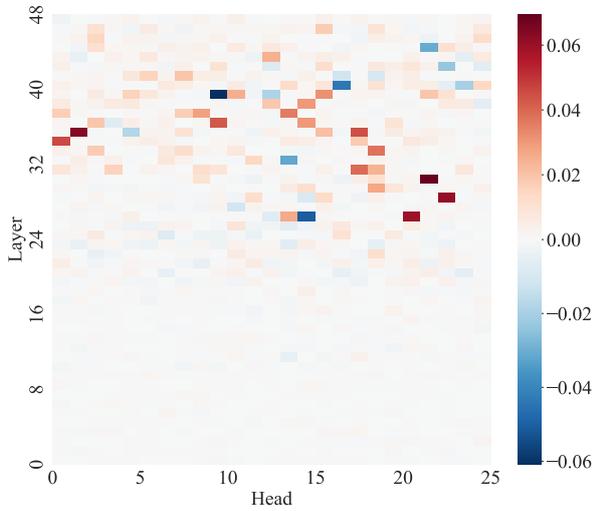


Figure 16: Context Heads of GPT-2 XL.

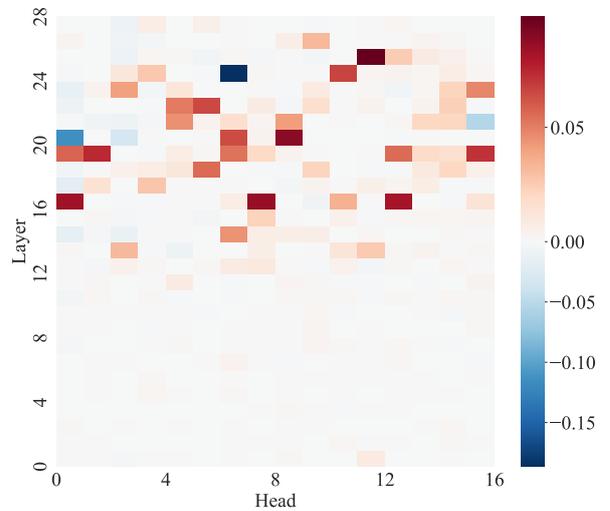


Figure 18: Context Heads of GPT-J.

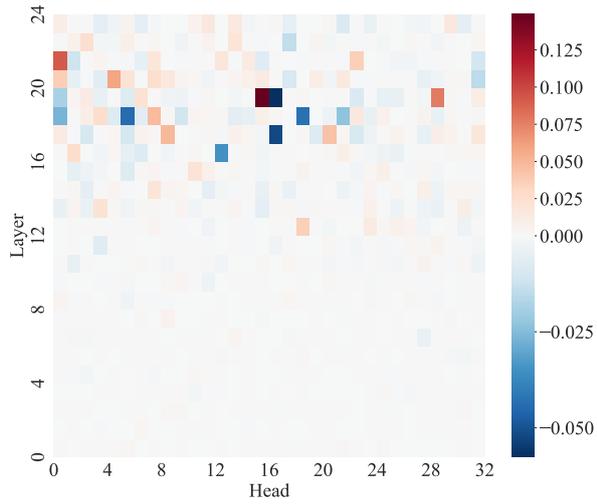


Figure 19: Memory Heads of OPT-1.3B.

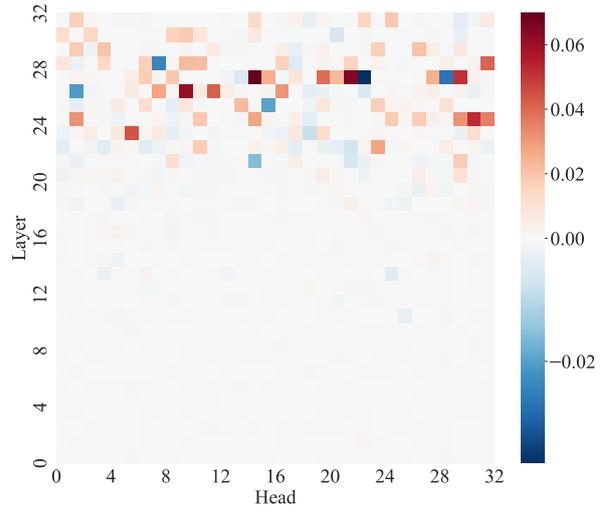


Figure 22: Context Heads of OPT-2.7B.

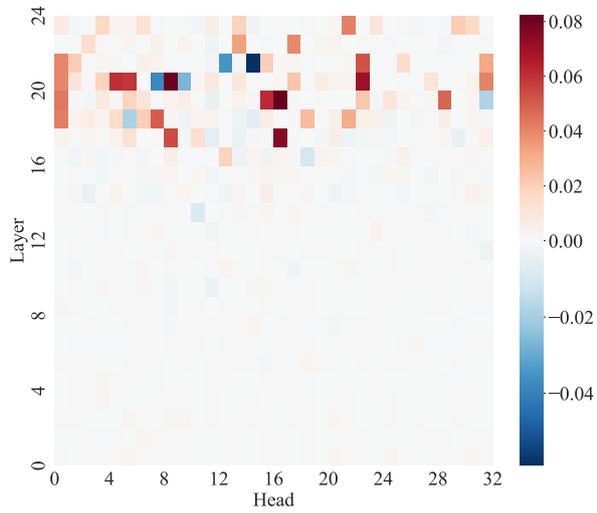


Figure 20: Context Heads of OPT-1.3B.

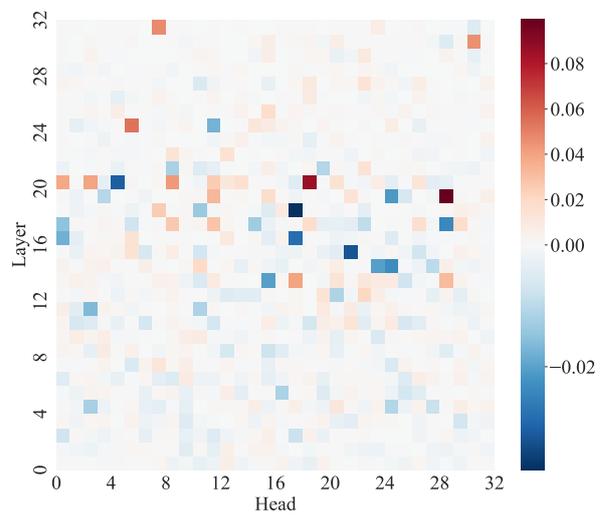


Figure 23: Memory Heads of Pythia-6.9B.

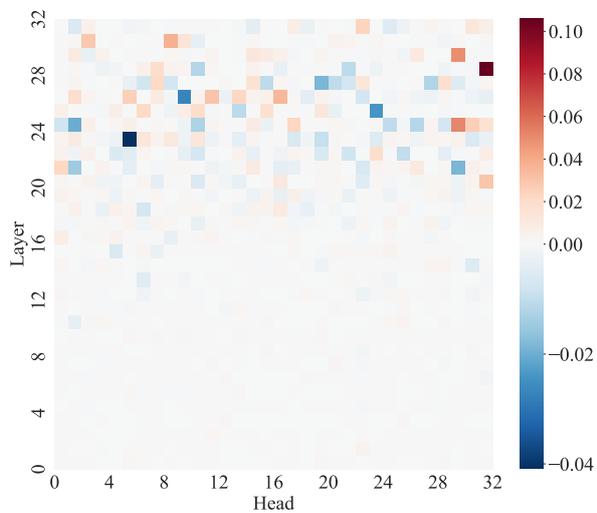


Figure 21: Memory Heads of OPT-2.7B.

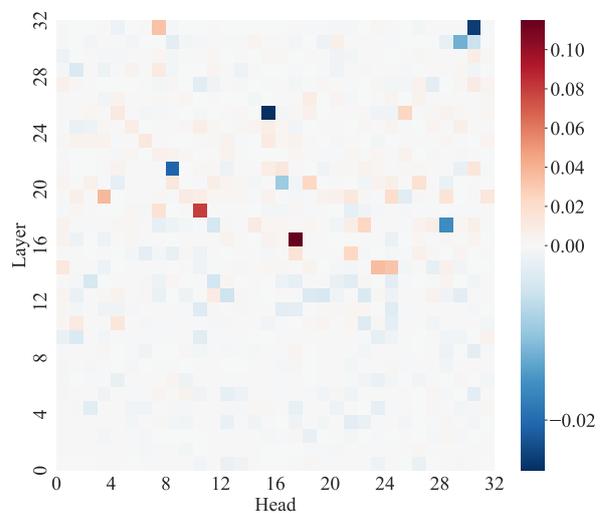


Figure 24: Context Heads of Pythia-6.9B.

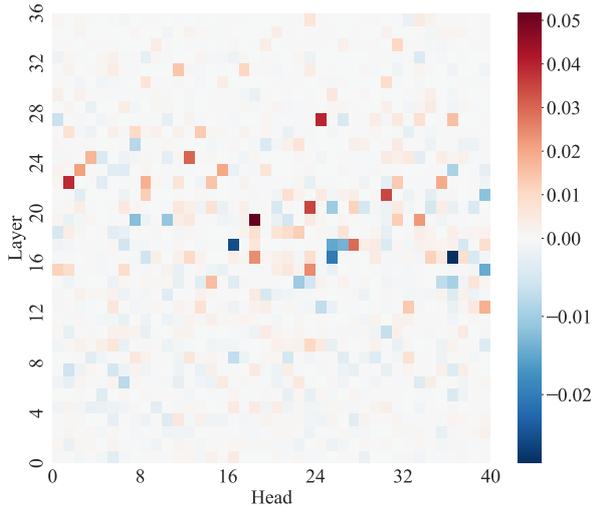


Figure 25: Memory Heads of Pythia-12B.

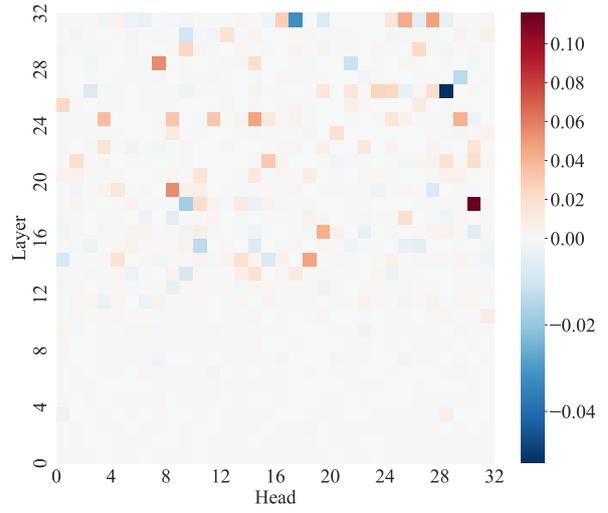


Figure 28: Context Heads of LLaMA2-7B.

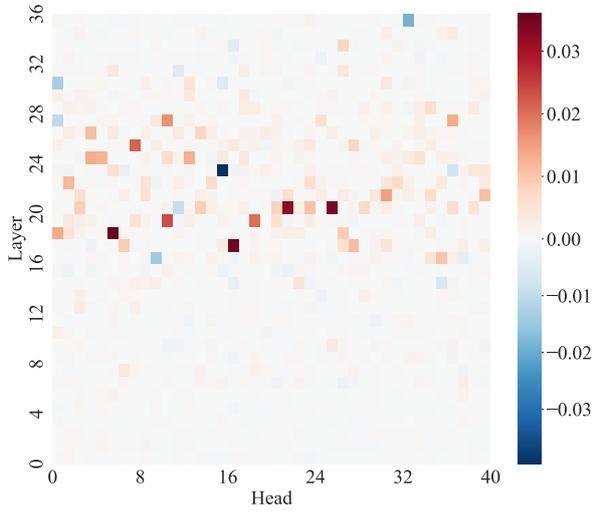


Figure 26: Context Heads of Pythia-12B.

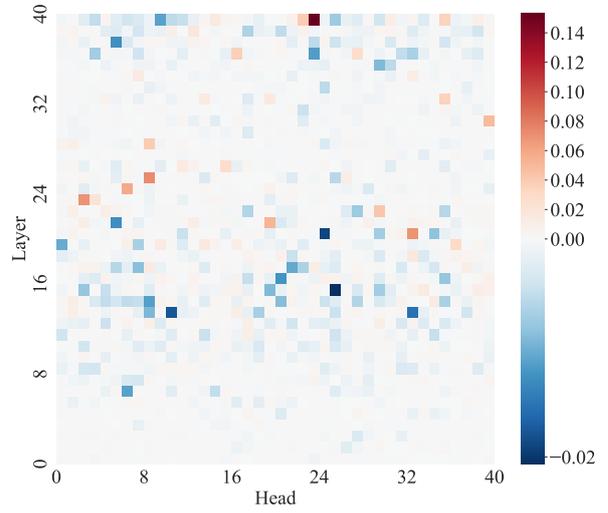


Figure 29: Memory Heads of LLaMA2-13B.

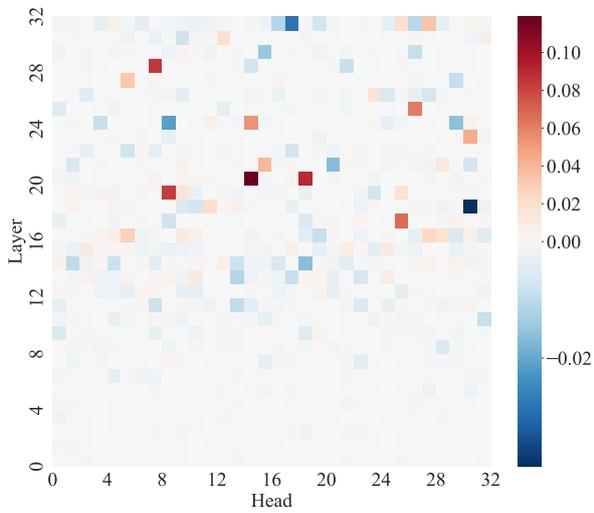


Figure 27: Memory Heads of LLaMA2-7B.

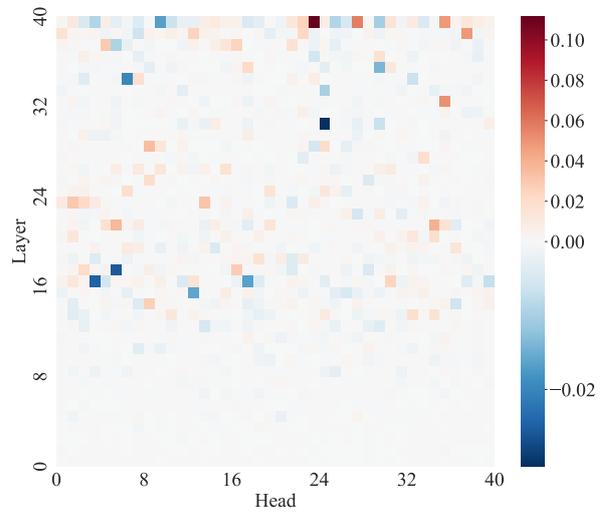


Figure 30: Context Heads of LLaMA2-13B.

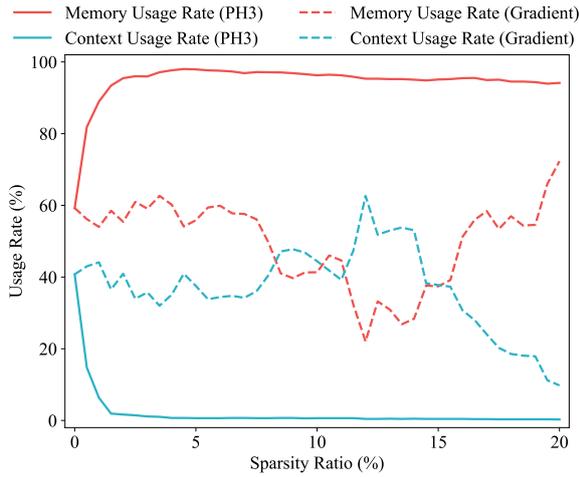


Figure 31: Impact of GPT-2 XL's sparsity ratio on improving internal memory usage rate.

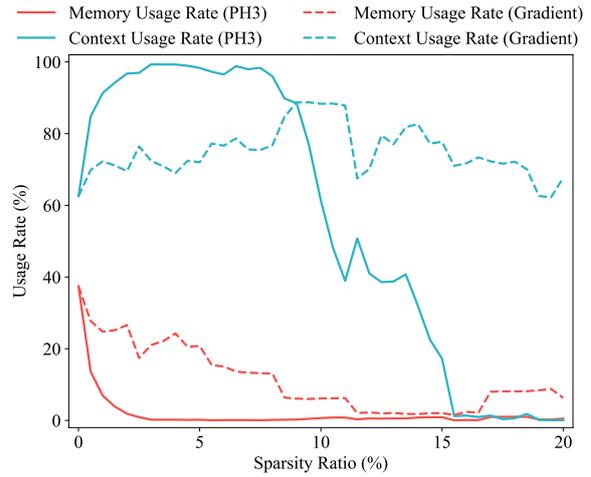


Figure 34: Impact of GPT-J's sparsity ratio on improving external context usage rate.

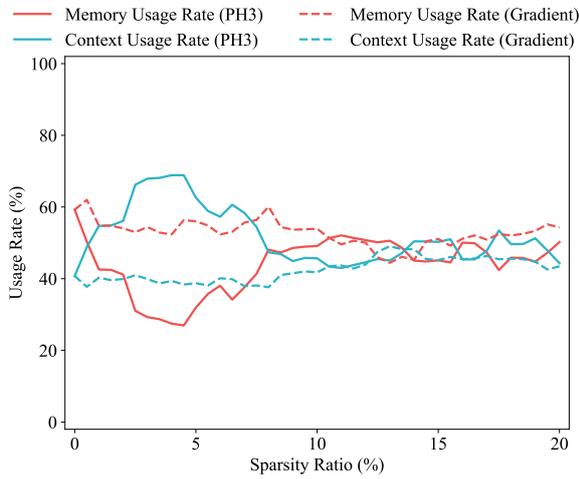


Figure 32: Impact of GPT-2 XL's sparsity ratio on improving external context usage rate.

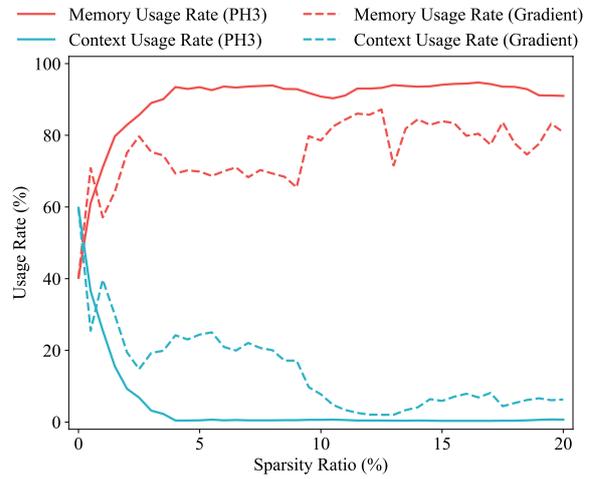


Figure 35: Impact of OPT-2.7B's sparsity ratio on improving internal memory usage rate.

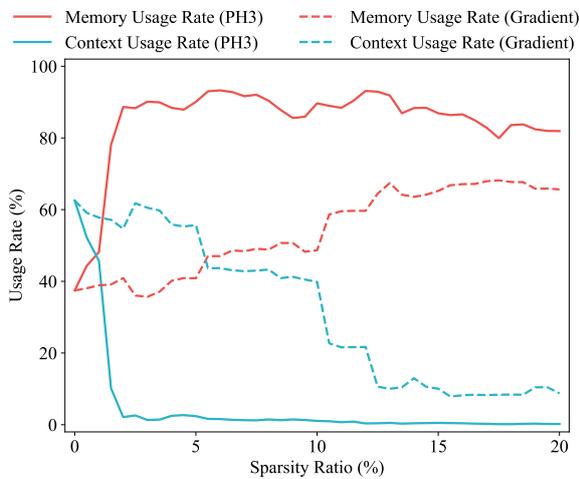


Figure 33: Impact of GPT-J's sparsity ratio on improving internal memory usage rate.

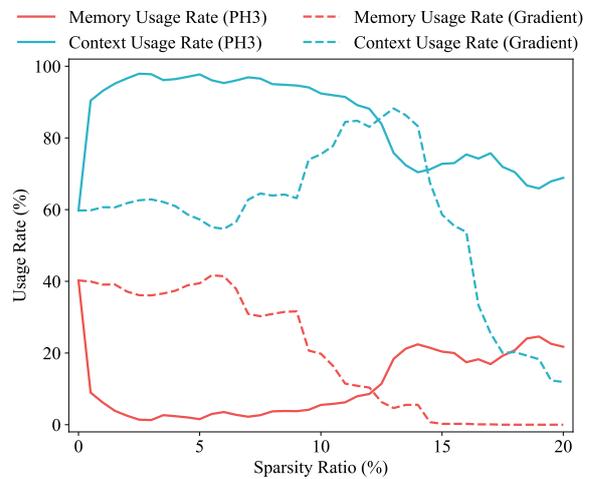


Figure 36: Impact of OPT-2.7B's sparsity ratio on improving external context usage rate.

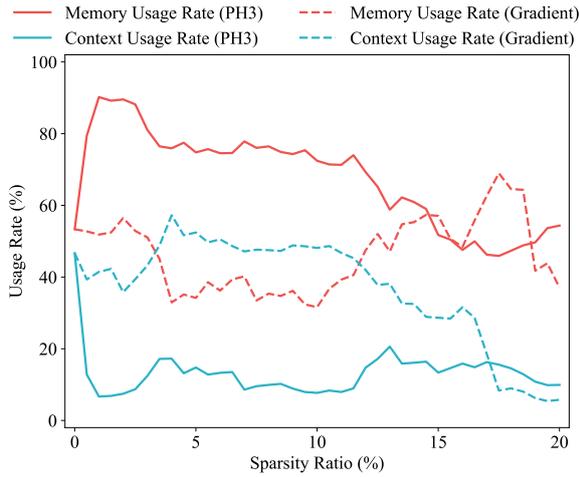


Figure 37: Impact of Pythia-6.9B's sparsity ratio on improving internal memory usage rate.

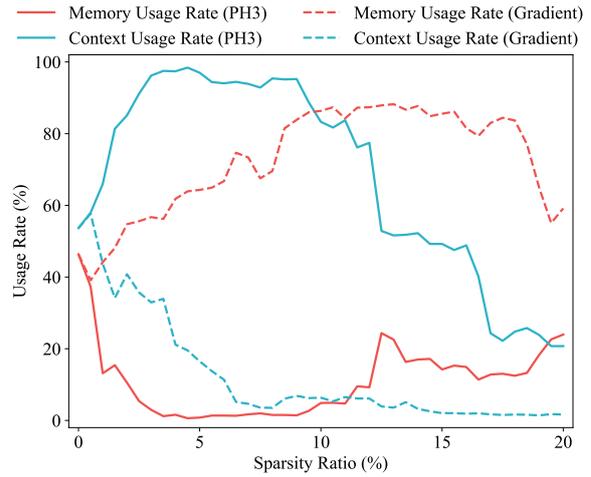


Figure 40: Impact of LLaMA2-7B's sparsity ratio on improving external context usage rate.

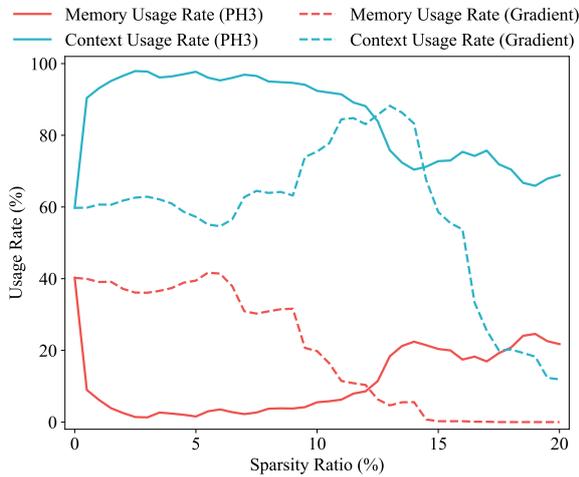


Figure 38: Impact of Pythia-6.9B's sparsity ratio on improving external context usage rate.

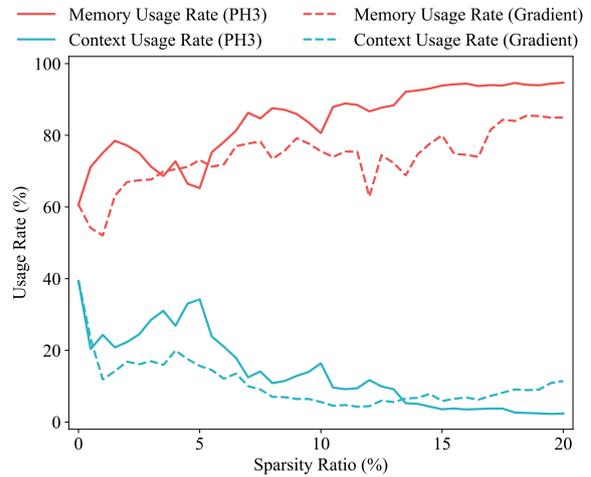


Figure 41: Impact of LLaMA2-13B's sparsity ratio on improving internal memory usage rate.

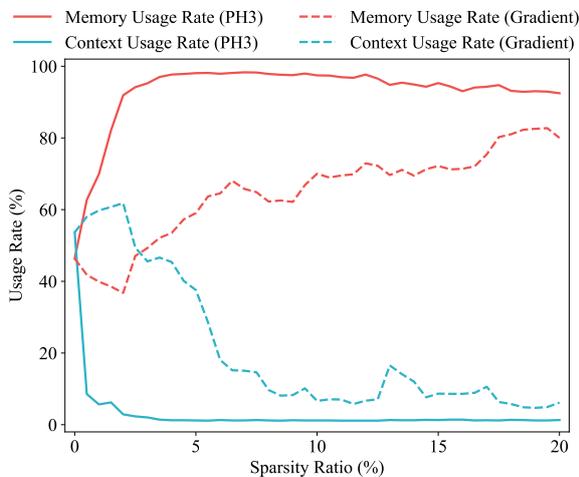


Figure 39: Impact of LLaMA2-7B's sparsity ratio on improving internal memory usage rate.

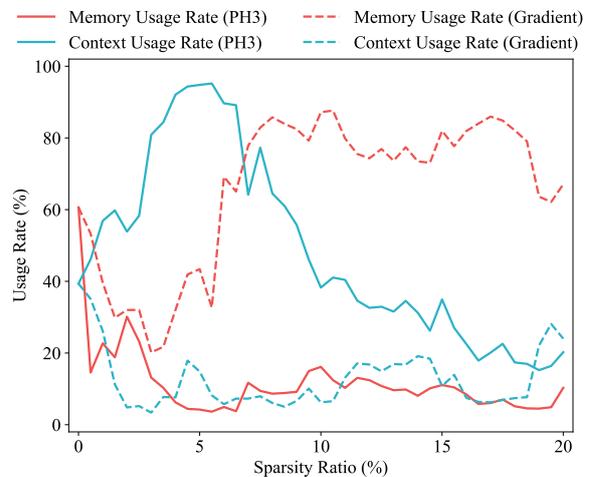


Figure 42: Impact of LLaMA2-13B's sparsity ratio on improving external context usage rate.

Model	Method	World Capital		World Capital D		Official Language		Country		Continent		
		RM	RC	RM	RC	RM	RC	RM	RC	RM	RC	
OPT-1.3B	↑ Memory	Base	40.5	59.5	36.3	63.7	20.3	79.7	26.8	73.2	19.2	80.8
		Prompt	19.7	78.4	37.5	57.2	7.9	91.4	16.9	82.6	7.9	90.1
		Gradient	82.7	12.2	56.4	21.2	37.1	50.4	38.0	61.1	20.5	55.9
		PH3 (Ours)	<b>95.0</b>	<b>0.2</b>	<b>87.2</b>	<b>1.9</b>	<b>70.3</b>	<b>16.6</b>	<b>47.7</b>	<b>49.1</b>	<b>43.4</b>	<b>51.7</b>
	↑ Context	Prompt	17.0	81.4	38.5	57.7	9.3	89.5	15.9	83.8	6.3	93.3
		CAD	8.1	86.5	31.6	60.0	<b>3.2</b>	<b>89.6</b>	<b>0.1</b>	<b>99.5</b>	5.1	89.1
		Gradient	22.9	73.7	35.7	63.8	12.4	82.5	16.3	82.9	17.6	80.2
		PH3 (Ours)	<b>0.2</b>	97.0	<b>7.9</b>	<b>69.6</b>	12.8	84.1	4.0	85.4	<b>1.6</b>	44.1
		+ Prompt	0.4	<b>99.2</b>	10.8	68.8	9.3	88.3	2.4	92.5	2.5	<b>92.1</b>
OPT-2.7B	↑ Memory	Base	40.2	59.8	46.6	53.4	8.8	91.2	26.5	73.5	4.3	95.7
		Prompt	17.7	80.3	24.2	71.6	3.4	96.3	12.7	87.2	1.6	98.1
		Gradient	75.4	19.3	10.3	79.7	13.3	57.7	42.9	56.3	5.3	94.1
		PH3 (Ours)	<b>93.4</b>	<b>0.5</b>	<b>2.8</b>	<b>87.6</b>	<b>75.6</b>	<b>1.5</b>	<b>56.3</b>	<b>38.9</b>	<b>29.3</b>	<b>55.5</b>
	↑ Context	Prompt	4.7	94.9	11.1	86.9	3.2	96.3	13.2	86.5	0.7	99.0
		CAD	10.8	72.2	28.8	46.3	2.7	87.5	9.0	89.1	0.5	99.0
		Gradient	36.1	62.9	43.3	53.5	7.3	91.6	23.7	76.3	4.1	95.9
		PH3 (Ours)	1.3	97.8	3.4	81.6	4.1	94.9	9.0	90.9	0.9	98.6
		+ Prompt	<b>0.8</b>	<b>98.3</b>	<b>1.3</b>	<b>94.6</b>	<b>1.5</b>	<b>98.2</b>	<b>6.9</b>	<b>93.1</b>	<b>0.0</b>	<b>99.5</b>
Pythia-6.9B	↑ Memory	Base	53.3	46.6	74.8	25.2	49.7	50.3	41.3	58.7	39.1	60.9
		Prompt	44.7	51.2	41.8	37.8	16.5	81.4	12.8	87.1	36.2	61.6
		Gradient	56.5	35.8	72.8	22.4	56.7	36.3	37.9	62.1	40.1	58.9
		PH3 (Ours)	<b>90.2</b>	<b>6.7</b>	<b>88.4</b>	<b>10.2</b>	<b>71.5</b>	<b>11.1</b>	<b>41.8</b>	<b>57.3</b>	<b>66.0</b>	<b>31.0</b>
	↑ Context	Prompt	32.7	63.7	32.0	44.8	8.9	90.5	10.4	89.5	31.7	75.6
		CAD	14.3	55.1	<b>22.0</b>	27.6	<b>3.3</b>	78.1	8.5	91.2	12.3	82.2
		Gradient	41.4	53.7	61.8	35.6	48.2	51.3	34.3	65.6	41.7	53.0
		PH3 (Ours)	6.7	81.7	34.4	30.0	16.4	70.0	3.4	96.3	3.3	94.5
		+ Prompt	<b>0.6</b>	<b>98.6</b>	24.4	<b>60.4</b>	<b>3.3</b>	<b>95.6</b>	<b>0.3</b>	<b>99.5</b>	<b>0.7</b>	<b>98.8</b>
Pythia-12B	↑ Memory	Base	59.7	40.3	64.6	35.4	34.3	65.7	35.0	65.0	43.0	57.0
		Prompt	5.8	94.1	43.7	53.3	11.1	85.8	2.4	97.5	10.1	88.5
		Gradient	62.9	27.7	63.1	30.2	39.3	37.5	46.2	53.3	42.1	56.9
		PH3 (Ours)	<b>95.0</b>	<b>0.6</b>	<b>82.4</b>	<b>2.2</b>	<b>69.9</b>	<b>6.9</b>	<b>57.1</b>	<b>35.9</b>	<b>70.1</b>	<b>9.6</b>
	↑ Context	Prompt	6.2	93.6	34.1	62.0	21.1	77.6	<b>1.7</b>	<b>98.3</b>	6.3	92.5
		CAD	3.1	65.7	13.6	42.9	2.0	89.4	3.7	94.7	11.3	80.0
		Gradient	59.3	19.6	33.3	25.0	21.3	54.6	28.5	71.3	40.2	52.1
		PH3 (Ours)	18.6	76.1	56.9	33.3	10.9	80.4	16.9	76.7	26.9	67.9
		+ Prompt	<b>1.9</b>	<b>97.6</b>	<b>17.7</b>	<b>75.8</b>	<b>3.8</b>	<b>95.5</b>	2.2	97.7	<b>2.4</b>	<b>97.0</b>
LLaMA2-13B	↑ Memory	Base	60.6	39.4	74.6	25.0	1.6	98.4	26.2	73.7	5.6	93.3
		Prompt	0.4	99.6	78.5	21.0	0.5	99.5	22.3	77.6	10.7	89.3
		Gradient	77.7	10.0	89.9	9.5	26.1	68.3	<b>48.1</b>	<b>51.4</b>	<b>30.7</b>	49.0
		PH3 (Ours)	<b>86.3</b>	<b>12.5</b>	<b>91.5</b>	<b>8.1</b>	<b>71.2</b>	<b>11.6</b>	47.7	51.6	11.9	<b>45.4</b>
	↑ Context	Prompt	<b>0.0</b>	<b>100.0</b>	70.9	28.7	<b>0.0</b>	<b>100.0</b>	7.1	92.7	3.6	96.4
		CAD	6.2	91.0	<b>1.1</b>	<b>98.8</b>	<b>0.0</b>	<b>100.0</b>	0.5	99.5	<b>0.0</b>	99.6
		Gradient	46.3	33.9	74.1	25.4	16.5	83.3	20.2	79.5	3.6	96.4
		PH3 (Ours)	6.2	92.1	24.1	75.9	1.3	98.7	5.6	94.4	0.4	99.6
		+ Prompt	<b>0.0</b>	<b>100.0</b>	31.0	68.0	<b>0.0</b>	<b>100.0</b>	<b>0.0</b>	<b>100.0</b>	<b>0.0</b>	<b>100.0</b>

Table 2: Experimental results of OPT-1.3B, OPT-2.7B, Pythia-6.9B, Pythia-12B and LLaMA2-7B on five datasets. Bolds denote the best results.

Method	GPT-2 XL	GPT-J	OPT-2.7B
Base	45.6	54.8	51.4
Prompt	47.6	57.4	54.1
CAD	44.5	55.0	50.2
Gradient	45.3	55.0	50.8
PH3 ( $k = 1$ )	47.1	57.5	53.5
PH3 ( $k = 3$ )	52.2	58.6	55.4
+ Prompt	<b>54.0</b>	<b>59.6</b>	<b>56.7</b>
PH3 ( $k = 5$ )	49.4	56.3	54.0

Table 3: Experimental results (Recall) of GPT-2 XL, GPT-J and OPT-2.7B on the NQ dataset. Bolds denote the best results.

Model	#Layer $L$	#Head $M$
GPT-2 XL	48	25
GPT-J	28	16
OPT-1.3B	24	32
OPT-2.7B	32	32
Pythia-6.9B	32	32
Pythia-12B	36	40
LLaMA2-7B	32	32
LLaMA2-13B	40	40

Table 4: Model details.