001

002

004

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

LOCATE-THEN-UNLEARN: AN EFFECTIVE METHOD OF MULTI-TASK CONTINUAL LEARNING FOR LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Nowadays large language models (LLMs) have achieved remarkable success in various NLP tasks. However, they often misinterpret human instructions and generate incorrect or outdated responses, highlighting the need for more effective continual learning techniques. While recent efforts have introduced unlearning methods to remove erroneous knowledge, existing approaches still struggle in multitask learning scenarios. To overcome these limitations, we propose Locate-then**unlearn**, a new framework that identifies and selectively unlearns task-specific neurons to enable efficient multi-task learning. We hypothesize that LLM neurons can be broadly categorized into task-specific neurons for handling individual tasks, and general neurons to maintain the model's foundational capabilities. To accurately identify task-specific neurons, the locating process includes: (1) ranking task-related neurons based on their importance to each task, and (2) identifying task-specific neurons by applying intervention to assess how neuron activity impacts task performance, isolating those most critical to each task. We conduct comprehensive evaluations in two experimental setups: single-task specialization and multi-task generalization. The results show that our method significantly improves performance across both settings. This indicates that our method effectively balances model efficiency and accuracy in multi-task continual learning.

032

1 INTRODUCTION

Large language models (LLMs) have recently demonstrated outstanding performance in diverse areas such as natural language understanding (Dušek et al., 2020), mathematical reasoning (Imani et al., 2023), and knowledge-intensive question answering (Sun et al., 2024). However, despite their impressive capabilities, LLMs remain prone to misinterpreting human instructions and generating incorrect or outdated responses (Bai et al., 2024; Chen et al., 2024a). This leads to the exploration of various continual learning and lifelong model editing techniques aimed at refining LLMs' behavior over time (Ji et al., 2024; Wang & Li, 2024).

In addition to directly fine-tuning LLMs on specific tasks, recent studies have introduced unlearning 040 techniques to enable LLMs to discard specific erroneous knowledge while preserving their overall 041 functionality. Building upon this concept, Ni et al. (2023) propose the "forgetting before learning" 042 paradigm, where LLMs are first trained to forget incorrect answers before learning new information, 043 leading to improved performance over direct fine-tuning. This approach mirrors human cognitive 044 processes, where learning is often more effective when mistakes are first identified and avoided. 045 However, current unlearning methods face limitations in maintaining performance across multiple 046 tasks simultaneously. One major issue is the lack of task-specific differentiation, which causes 047 interference between the knowledge acquired for different tasks (Dong et al., 2023). This may 048 lead to catastrophic forgetting of previously learned tasks. And the difference in the order of finetuning between tasks will also have a significant impact on the performance of the model (Bell & Lawrence, 2022; Wang & Li, 2024). Additionally, fine-tuning all model parameters across multiple 051 tasks consumes considerable computational resources and significantly reduces learning efficiency. While parameter-efficient fine-tuning solutions have been proposed, their effectiveness diminishes 052 in multi-task settings, making it challenging to strike a balance between efficiency and overall performance (Leng & Xiong, 2024).

- To address these limitations, we aim to identify regions within LLMs that are responsible for handling different tasks. Inspired by the neuron definition from Chen et al. (2024b); Tang et al. (2024), we hypothesize that neurons in LLMs can be mainly categorized into two types: (i) **taskspecific neurons**, which focus on processing particular tasks, and (ii) **general neurons**, which aim to maintain the model's core capabilities in text understanding and generation. From a perspective of parameter-efficient multi-task learning, we selectively update the task-specific neurons while leveraging general neurons to preserve the model's foundational abilities, thereby improving both efficiency and effectiveness in multi-task scenarios.
- 062 The primary challenge in achieving this lies in accurately formulating the task-specific neurons for 063 different tasks. To this end, we propose Locate-then-unlearn, a new framework that identifies 064 task-specific neurons and selectively unlearns them to enable efficient and continual learning across multiple tasks for LLMs. Based on previous works (Geva et al., 2022; Meng et al., 2022a), we extract 065 logit scores from the activation layer of each neuron, using these scores to determine the importance 066 of each neuron for the given task. Tang et al. (2024) utilizes the activation score of each neuron 067 to identify language-specific regions. Inspired by this, in our multi-task setting, we adaptively rank 068 neurons based on their logit scores across different tasks. Then we can determine which neurons 069 contribute the most to a specific task and filter out less relevant ones. Since some task-related neurons may exhibit high logit scores across multiple tasks, we use a neuron intervention method 071 to further assess their task specificity. By comparing the difference in correct answer probability 072 before and after neuron intervention, we can identify neurons whose performance shifts significantly, 073 categorizing them as task-specific neurons.
- Once task-specific neurons are identified, we apply parameter-efficient fine-tuning on these neurons within the unlearning set. In terms of downstream task learning, we design two evaluation settings for comprehensive comparison: (i) Single-task specialization: Fine-tuning a separate unlearned model for each downstream task and evaluating each model independently. (ii) Multi-task generalization: Fine-tuning a single unlearned model across multiple task datasets and evaluating its performance on all tasks collectively.
- Experimental results show that our method significantly outperforms all baselines in multi-task generalization, demonstrating its superiority in enhancing the model's ability to generalize across tasks. Additionally, in single-task specialization, our method achieves optimal results while reducing training complexity, indicating that our approach effectively identifies task-specific neurons and balances both performance and efficiency in multi-task learning.
- To sum up, our main contributions can be described as follows:
- We propose Locate-then-unlearn, a new framework that facilitates efficient and continual learn ing for LLMs across multiple tasks.
 - We develop a new locating method to accurately identify task-specific neurons by assessing their importance and isolating those critical to each task.

• We design two experimental settings for comprehensive evaluation: single-task specialization and multi-task generalization. Experimental results show that our method achieves significant improvements in both settings, demonstrating an effective balance between performance and efficiency.

094 095 096

089

091

092

093

2 RELATED WORK

-)97
- 098 099

2.1 THE STRUCTURE AND KNOWLEDGE MECHANISM OF LARGE LANGUAGE MODELS

100 The development of LLMs has revealed great potential in solving various NLP tasks. However, the 101 occurrence of hallucinations in LLMs may hinder their broader adoption in real-world applications. 102 Consequently, neuronal interpretability has gained much attention in recent years. Several studies 103 have investigated the mechanisms underlying knowledge storage in LLMs. For instance, Geva et al. 104 (2022) and Meng et al. (2022a) have found that the multilayer perceptron (MLP) layers in Trans-105 former models function as key-value memory, storing vast amounts of knowledge. Other works, such as Geva et al. (2023a), Lv et al. (2024), and Yu & Ananiadou (2024b), have shown that knowl-106 edge accumulates progressively throughout the layers. In this paper, we build on the perspective that 107 factual knowledge is primarily stored within the MLP layers of LLMs.

108 2.2 MODEL EDITING

As the knowledge stored in LLMs may become outdated due to the rapid growth of our society, it is essential to edit and update this information accordingly. Some recent studies focus on identifying where knowledge is stored before editing. For example, ROME (Meng et al., 2022a) uses the method of attributing logits to find the location of knowledge and then edits it by updating specific factual associations. Meng et al. (2022b) is an effective method to locate knowledge and directly update large scale memories. Our work follows the workflow of locating and editing by identifying multiple task-specific neurons and then updating them accordingly.

117 118

2.3 LARGE LANGUAGE MODEL UNLEARNING

119 Machine unlearning (Cao & Yang, 2015) serves as an important technique to remove the knowledge 120 about the restricted data while keeping other knowledge and system abilities. Yao et al. (2023a) and 121 Maini et al. (2024) use the method of gradient ascent to unlearn harmful or private knowledge. And 122 Ni et al. (2023) employs parametric arithmetic to facilitate the forgetting of old knowledge and learn-123 ing of new knowledge. They first finetune LLM on the old knowledge and use the original model to subtract the old knowledge parameters to finish the knowledge update. However, directly em-124 ploying parametric arithmetic without considering the utility of each neuron can also be harmful to 125 model performance on other tasks. Chen et al. (2024c) proposes allow-redundant alignment method 126 named ALLO, focusing on optimizing the most related neurons with the most useful supervised sig-127 nals. They use the signal to detect unaligned knowledge and unlearn it. Our work further explores 128 unlearning by not only identifying task-specific neurons across different tasks but also selectively 129 unlearning these neurons to prevent knowledge conflicts and preserve model performance.

130 131 132

133

134 135

136

137 138

147

152

154

3 PRELIMINARY

3.1 TASK-RELATED NEURON LOCALIZATION

Denote the hidden state of the *i*-th layer for a specific token as $h^i \in \mathbb{R}^d$, the multi-layer perceptron (MLP) module within the *i*-th layer can then be described as follows:

$$oldsymbol{h}^i = \sigma(ilde{oldsymbol{h}}^i oldsymbol{W}_1^i) \cdot oldsymbol{W}_2^i)$$

where W_1^i and W_2^i represents trainable parameters of transition matrix, \tilde{h}^i represents output of *i*-th 139 MHA layer and $\sigma(\cdot)$ denotes the activation function. As mentioned in Tang et al. (2024), a "neuron" 140 in LLMs is regarded as a linear transformation to a specific column in W_1^i , followed by a non-linear 141 activation. Also, we follow the consensus of Nair & Hinton (2010), which considers the j-th neuron 142 in the *i*-th FFN layer to be activated when its activation value is positive. Based on this definition 143 of activated neurons and the calculation of activation probability for each language in Tang et al. 144 (2024), we regard the proportion of positive activation scores of the j-th neuron in the *i*-th layer as 145 the importance of its contribution to each task k, which can be formulated as: 146

$$s_{i,j}(k) = \mathbb{E}\left(\mathbb{I}(\sigma(\tilde{\boldsymbol{h}}^{i}\boldsymbol{W}_{1}^{i})_{j} > 0) \mid ext{task } k
ight),$$

where I is an indicator function to determine whether the result is positive or negative, thus, a neuron is deemed task-related if its importance score $s_{i,j}(k)$ ranks within the top r_k for task k. In our framework, we set the threshold r_k to represent the top 2%. By ranking neurons according to their importance scores across tasks, we can effectively localize task-related neurons for each task.

153 3.2 UNLEARNING PROCESS AND OBJECTIVE

In the unlearning period, referring to Yao et al. (2023b), the unlearning data is defined as (x_u, y_u) and all related data can be categorized into three types: (1) *Same unlearning data*, which represents the exact data during the editing period. (2) *Paraphrased data*, which retains the same meaning as the original data but is expressed differently, denoted as $R(x_u, y_u)$. (3) *Similar representation data*, which bears some semantic resemblance to the original unlearning data but the exact meaning differs, referred to as $N(x_u, y_u)$. Thus, the unlearning process and objective can be formulated as:

$$f_{\theta^*}(x_i) = \begin{cases} y_i^{new} & \text{if } x_i \in (x_u, y_u) \text{ or } R(x_u, y_u), \\ f_{\theta}(x_i) & \text{if } x_i \in N(x_u, y_u) \text{ or other.} \end{cases}$$

162 The goal of the knowledge updating task is to modify the model's outputs only for x_{μ} and its para-163 phrased versions $R(x_u, y_u)$, without affecting answers related to neighboring knowledge $N(x_u, y_u)$ 164 or other unrelated data. This ensures that unlearning is both precise and minimally disruptive to the 165 model's overall knowledge.

Knowledge in LLMs can be updated by supervised fine-tuning on a task-related dataset. Inspired by Ni et al. (2023), for a given LLM and its parameters θ_i , the knowledge updated parameters θ_u can be computed by subtracting the fine-tuned parameters θ^* and the original parameters θ_0 , which is given by:

$$\theta_u = \theta^* - \theta_0.$$

4 METHOD

166

167

168

169

170 171 172

173 174

175

4.1TASK-SPECIFIC NEURON IDENTIFICATION

176 We assume that a task-related neuron should only be fine-tuned for its specific task. To address the 177 challenge of determining a neuron's specificity when it relates to more than one task, we measure 178 the specificity degree by comparing the correct answer probabilities when the neuron is intervened 179 for one task versus its unintervened performance.

In our setting, we consider $\tilde{\mathcal{N}} = \{\tilde{N}^1, \dots, \tilde{N}^M\}$ as a collection of M neurons, where each neuron 181 may be associated with more than one task. Assuming the *m*-th neuron \tilde{N}^m is related to the totally 182 n tasks $\{d_1,\ldots,d_n\}$, we aim to find the task to which $\tilde{N^m}$ is most specific. To achieve this, we 183 construct toy datasets to perform inference. Each toy dataset is actually a random subsample of 184 training datasets for each task. In our main experiments, each toy dataset contains 1000 randomly 185 sampled cases from the corresponding training set. Additionally, we also conduct an ablation study 186 to evaluate the impact of different sampling methods on the performance of the toy dataset. We first 187 use the complete model to perform inference on the toy datasets in each task, yielding the average 188 predicted probabilities of the correct answer on all toy datasets for each neuron \tilde{N}^m , denoted as 189 $\{\mathbb{P}(d_1^m),\ldots,\mathbb{P}(d_n^m)\}$. When we need to intervene on the neuron $\tilde{N^m}$, we set its parameters to 0, 190 and then the new probabilities for all toy datasets can be represented as $\{\mathbb{P}(d_1^m)_{new}, \ldots, \mathbb{P}(d_n^m)_{new}\}$. 191

When identifying task-specific neurons, we focus on neurons that significantly influence one task 192 while having minimal impact on others. Inspired by Geva et al. (2023b); Cohen et al. (2024); Yu 193 & Ananiadou (2024a), which use changes in the probabilities of the correct answer to locate key 194 neurons causing the final prediction, we measure a neuron's relative importance for a given task in 195 comparison to its impact on other tasks. Specifically, for each task $d_k \in \{d_1, \ldots, d_n\}$ and each 196 specific neuron \tilde{N}^m , we can calculate its importance score as: 197

$$S_k^m = \mathbb{P}(d_k^m)_{new} - \mathbb{P}(d_k^m) - \sum_{i \in \{1,\dots,N\} \setminus \{k\}} |\mathbb{P}(d_i^m)_{new} - \mathbb{P}(d_i^m)|.$$
(1)

199 200

201 The former term $\mathbb{P}(d_k^m)_{new} - \mathbb{P}(d_k^m)$, denoted as C_k^m , describes the influence of deactivating neuron 202 \tilde{N}^m on task d_k . The latter term reflects the influence of this neuron on all other tasks except d_k . 203 Using only C_k^m can be regarded as an ablation compared to our full identification method. The 204 reason for adding this non-negative term at the end is that when measuring the specificity of neuron 205 \tilde{N}^m to task d_k , any change in the importance scores for other tasks except d_k , whether positive 206 or negative, should be considered an adverse effect on the specificity of the neuron to task d_k . A 207 higher S_k^m indicates that this neuron is more specific to task d_k . Building upon this, we can conclude 208 that the neuron \tilde{N}^m is considered the task-specific neuron for task d_k if and only if its score S_k^m is 209 the largest among all n tasks. In subsequent sections, we focus on unlearning and fine-tuning this 210 neuron specifically for the task d_k dataset.

211

212 4.2 MULTI-TASK UNLEARNING AND RELEARNING

213

During the unlearning process, we first fine-tune the model θ_0 on a dataset containing knowledge to 214 be unlearned (e.g., false knowledge). Unlike Ni et al. (2023) that process on the overall model, we 215 focus the knowledge update exclusively on task-specific areas. For each task $d_k \in [d_1, ..., d_n]$, if the



Figure 1: Overview of the proposed "Locate-then-unlearn" framework, including two main modules: (a) task-specific neuron identification, and (b) multi-task unlearning and relearning.

parameters P^{li} of *i*th neuron in the *l* th layer are specific to d_k (as determined above), we subtract original neuron parameters P_0^{li} by the parameters P_{false}^{li} after fine-tuning on the false knowledge. If the parameters are not specific to d_k , they remain unchanged. In this way, we can obtain the updated ΔP_{false}^{li} (defined as the updated parameters brought by fine-tuning false knowledge) as follows:

$$\Delta P_{false}^{li} = \begin{cases} 0 & \text{if } P^{li} \notin d_k, \\ P_{false}^{li} - P_0^{li} & \text{if } P^{li} \in d_k. \end{cases}$$
(2)

Based on the parameter update ΔP_{false}^{li} for each neuron at each layer, we can obtain the overall model parameter update $\Delta \theta_{false}$. Then we can implement the overall model unlearning process by subtracting the original model parameters θ_0 from the updated model parameters that are fine-tuned on false knowledge. This process is calculated as:

$$\theta_{\delta} = \theta_0 - \lambda \Delta \theta_{false},\tag{3}$$

where λ is a hyper-parameter to control the update rate. Once this process is completed, the model could discard the outdated knowledge. Next, we inject true knowledge into θ_{δ} through supervised fine-tuning. This fine-tuning process specifically targets the parameters related to the new knowledge. The parameter update process during the relearning phase is similar to the previous process, where the parameters are denoted as P_{true}^{li} after being fine-tuned on the true knowledge. The parameter updates ΔP_{true}^{li} for the true knowledge are then expressed as:

$$\Delta P_{true}^{li} = \begin{cases} 0 & \text{if } P^{li} \notin d_k, \\ P_{true}^{li} - P_0^{li} & \text{if } P^{li} \in d_k. \end{cases}$$
(4)

We follow similar approaches to update the model parameters in the relearning process, resulting in a refined model characterized by the parameters θ^* :

$$\theta^* = \theta_\delta + \lambda \Delta \theta_{true}.$$
(5)

Figure 1 illustrates the overview of the proposed framework. We consider two settings to evaluate the
 refined model: (1) In the multi-task generalization setting, we fine-tune a single unlearned model
 across multiple task datasets and evaluate its performance across all tasks collectively. (2) In the
 single-task specialization setting, we fine-tune separate unlearned models for each downstream task
 dataset and evaluate each model independently.

For the design of the multi-task generalization setting, we determine the relearning order across multiple tasks following Leng & Xiong (2024), which first targets generation tasks and then classification tasks. In our experiments, ZsRE, SingleEQ, and PKURLHF belong to generation tasks, while
SST-2 and QQP are classification tasks. Therefore, we set the relearning order as ZsRE, SingleEQ, PKURLHF, SST-2, and QQP in our main results. Furthermore, we will also discuss how our method tackles catastrophic forgetting and explore the sensitivity of the model performance to the relearning order of the datasets.

277 278

279 280

281

5 EXPERIMENTAL SETUP

5.1 DATASETS

282 In this work, we utilize five datasets, each representing a distinct task. For the knowledge QA task, we employ ZsRE (Levy et al., 2017), a widely recognized Question Answering dataset that lever-283 ages question rephrasings generated through back-translation. ZsRE contains over 160,000 samples. 284 For the logical reasoning task, we use SingleEQ (Koncel-Kedziorski et al., 2015), which comprises 285 508 questions, 1,117 sentences, and 15,292 words; this dataset helps train LLMs to enhance logi-286 cal reasoning skills. For the human safety alignment task, we apply PKURLHF (Dai et al., 2023), 287 the first publicly available multi-round RLHF dataset in China, which includes constraints across 288 more than ten dimensions, such as insults, discrimination, crime, psychological harm, and privacy, 289 aligning LLMs with human values. Lastly, for natural language understanding, we use $SST-2^1$ and 290 QQP^2 . While both datasets focus on semantic classification, SST-2 is aimed at sentiment classifi-291 cation, whereas QQP focuses on similarity and paraphrase tasks. Thus, they can be considered as 292 representing different tasks. 293

5.2 EVALUATION METRICS

296 In the knowledge QA task, our goal is to modify answers for original and paraphrased questions while preserving the responses for related questions. Following Ni et al. (2023)'s evaluation set-297 ting, we measure our model's performance using four metrics: Reliability, Generalization, and 298 Locality. The first two assess accuracy in editing original and paraphrased questions, while the 299 third ensures that answers to unrelated questions remain unchanged. For SingleEQ, QQP, and SST-2 300 tasks, we use the Accuracy index to evaluate performance. In the PKURLHF task, which focuses 301 on aligning outputs with human values and avoiding harmful content, we assess performance using 302 the Harmful Rate indicator. This requires the GPT-4 model to identify and count harmful content 303 in its outputs. 304

305 306

323

295

5.3 IMPLEMENTATION DETAILS

307 We adopt different settings in multi-task generalization and single-task specialization, and in each 308 setting, we apply two backbones: OPT-1.3B (Zhang et al., 2023) and LLAMA2-7B (Touvron et al., 309 2023). For the multi-task generalization setting, the batch size is 2, the param of Adam is set as 0.9 and 0.995, the learning rate is set 6e-5, and r_k is set as top 2%. We firstly continuously 310 unlearn five datasets' old knowledge, then update the model and continuously learn five datasets' 311 new knowledge. After fine-tuning all datasets, we collectively test our model on five tasks. For the 312 single-task specialization setting, we set the learning rate as 1e-4, while keeping the other hyper-313 parameters the same. On the hardware side, since we only update on task-specific neurons, we only 314 spend about 23GB (about half of an A100 40GB GPU). 315

316 317 5.4 BASELINES

To evaluate the effectiveness of our proposed unlearning method, we compare our method with these baseline methods: (1) **Directly fine-tuning**. We do not use any unlearning method and just finetuning the original model with the true answers of editing data. (2) **ROME** (Meng et al., 2022a), a method updating specific factual associations with causal intervention. (3) **MEMIT** (Meng et al.,

¹https://huggingface.co/datasets/stanfordnlp/sst2

²https://huggingface.co/datasets/SetFit/qqp

			ZsRE		SingleEQ	PKURLHF	SST-2
		Specificity↑	Generality↑	Locality [↑]	Acc↑	Harmful Rate↓	Acc↑
	Directly fine-tuning	48.64%	45.87%	41.30%	13.54%	31.70%	65.86%
	ROME	29.51%	27.99%	85.18%	10.18%	37.32%	61.04%
	MEMIT	66.22%	63.15%	58.96%	15.86%	16.58%	68.87%
	F-learning	73.83%	69.85%	63.98%	19.32%	8.85%	73.39%
OPT-1 3B	W-NCFT	78.24%	73.31%	65.59%	22.03%	6.64%	78.57%
	Remove all overlapped neurons	80.11%	75.09%	70.63%	24.14%	4.83%	84.97%
	Preserve all overlapped neurons	81.03%	78.32%	72.76%	25.26%	1.67%	86.19%
	Randomly selection	81.25%	78.84%	72.95%	25.57%	2.54%	86.35%
	Locate-then-unlearn with C_k^m	82.63%	79.22%	74.20%	26.98%	1.87%	86.14%
	Locate-then-unlearn	85.04%	82.77%	76.36%	28.33%	1.44%	89.16%
	Directly fine-tuning	52.41%	48.92%	43.11%	16.02%	26.63%	68.89%
	ROME	35.31%	34.03%	88.96%	13.42%	34.89%	64.42%
	MEMIT	71.32%	67.10%	61.35%	18.57%	14.07%	72.91%
	F-learning	77.15%	72.24%	66.18%	21.84%	6.18%	76.54%
LLAMA2 7B	W-NCFT	81.93%	76.86%	68.84%	24.97%	4.55%	81.10%
LLAMA2-7D	Remove all overlapped neurons	84.55%	79.93%	72.08%	26.68%	2.88%	87.25%
	Preserve all overlapped neurons	85.71%	81.04%	74.39%	28.35%	1.02%	89.80%
	Random selection	86.19%	81.80%	75.06%	28.80%	1.85%	89.89%
	Locate-then-unlearn with C_{h}^{m}	87.97%	83.28%	77.14%	30.05%	1.24%	89.77%
	Logota than unloarn	80 21 0%	85 160%	70.010%	21 160%	1 1 2 07	02 200

Table 2: Main results on five tasks under the single-task specialization setting.

					0	1		0
			ZsRE		SingleEQ	PKURLHF	SST-2	QQP
		Specificity↑	Generality↑	Locality↑	Acc↑	Harmful Rate↓	Acc↑	Acc↑
	Directly fine-tuning	77.76%	72.17%	67.58%	27.19%	1.15%	90.04%	80.84%
	ROME	37.35%	34.82%	91.36%	13.30%	11.24%	73.44%	56.10%
OPT-1.3B	MEMIT	78.84%	74.67%	67.46%	18.84%	6.60%	85.25%	70.28%
	F-learning	80.18%	76.83%	72.57%	22.42%	3.67%	86.73%	75.80%
	W-NCFT	78.93%	74.82%	68.84%	24.97%	4.55%	81.10%	72.63%
	Locate-then-unlearn	85.95%	83.91%	78.42%	30.71%	1.03%	90.27%	80.79%
	Directly fine-tuning	81.08%	74.76%	70.48%	32.85%	1.04%	92.41%	83.29%
	ROME	43.98%	42.76%	93.22%	16.77%	9.32%	75.88%	60.19%
LLAMAN 7D	MEMIT	83.54%	79.03%	70.57%	21.19%	5.85%	86.20%	74.45%
LLAMA2-/B	F-learning	84.65%	80.22%	76.16%	25.31%	3.14%	87.59%	78.84%
	W-NCFT	82.77%	78.65%	73.23%	27.17%	4.02%	83.84%	75.56%
	Locate-then-unlearn	89.70%	86.15%	81.98%	34.36%	0.98%	92.64%	83.18%

2022b) which is an effective method to directly update large-scale memories. (4) F-learning (Ni et al., 2023), which forgets old knowledge by subtracting the parameters finetuned on false answers and then learns on the true answers. (5)W-NCFT (Leng & Xiong, 2024), which is a neuron-level continual fine-tuning method that utilizes relevance score to locate task-specific neurons and only fine-tunes the current task-specific neurons during continual learning.

EXPERIMENTAL RESULTS

6.1 MAIN RESULTS

The experimental results for multi-task generalization are presented in Table 1, while single-task specialization results are in Table 2. In the multi-task generalization setting, Locate-then-unlearn shows significant enhancements across all five experiments compared to the state-of-the-art W-NCFT method (p < 0.001, two-sided t-test). Specifically, on the ZsRE dataset, our approach improves specificity and generality by over 7 points in both the OPT-1.3B and LLAMA2-7B set-tings, while improving locality by over 10 points. Although the ROME method achieves the highest locality, it relies on modifying only a small number of parameters, resulting in poorer specificity and generality. Our method also outperforms others on the SingleEQ, PKURLHF, SST-2, and QQP datasets. In conclusion, our approach that effectively fine-tunes task-specific neurons has the optimal performance across all five tasks. In the single-task specialization setting, our method outperforms state-of-the-art results across ZsRE, SingleEQ, PKURLHF, and SST-2 (p < 0.005, two-sided t-test), although the improvements are less pronounced than in the multi-task generalization setting. Notably, our method slightly underperforms Directly fine-tuning on the QQP dataset. To investigate this discrepancy, we further conduct experiments detailed in Section 6.3.2. Overall, our approach still achieves the best results in single-task specialization compared to other baselines.

381

383

384

385

387

388

391 392

393

396 397

378 7sRF 379 SingleEQ PKURLF SST-2 0.9 0.9 0.8 0.8 0.7 0.7 382 Performance Performance 0.6 0.6 0.5 0.5 0.4 0.4 0.3 0.3 386 0.2 0.2 0.1 0.1 0.0 0.0 389 0.3 0.5 0.9 0.005 0.010 0.015 0.020 0.025 0.030 0.035 0.1 0.7 Unlearning hyperparameter I Threshold r_k

Figure 2: The impact of varying λ and r_k .

6.2 ABLATION STUDY

We conduct ablation experiments to identify the most effective components of our method in multitask generalization settings, designing four verification methods. First, we consider the removal or 399 preservation of overlapped neurons, these two constitute ablation methods called **Remove all over-**400 lapped neurons and Preserve all overlapped neurons. The third ablation method called Random 401 selection involves randomly selecting one related task as specific while keeping overlapped neurons 402 frozen during training on other tasks. Lastly, instead of using the designed algorithm with S_k^m , we 403 utilize C_k^m as positioning indicators, which means that we only consider the impact of a neuron on 404 one task and ignore its impact on other tasks. This method constitutes the fourth ablation method, 405 which is called **Locate-then-unlearn with** C_k^m .

406 We observe from Table 1 that when we compare four ablation methods and our Locate-then-unlearn 407 method, either preserving or removing all overlapped neurons is less effective than employing meth-408 ods that identify more specific tasks. Furthermore, Locate-then-unlearn using S_k^m demonstrates bet-409 ter performance than the ablation method which uses C_k^m , confirming that when we locate one-task 410 specific neuron, we should also guarantee that it has little impact on other tasks, otherwise, fine-411 tuning this neuron in one task will also greatly affect its performance in other tasks, which will 412 greatly affect model performance.

413 414

415 416

417

ADAPTABILITY ANALYSIS 6.3

6.3.1 **IMPACT OF VARYING HYPER-PARAMETERS**

418 To verify our method's adaptability on different parameter selections, we choose to first change the 419 rate of forgetting λ and observe the change of the model's overall performance using the LLAMA2-420 7B backbone under the multi-task generalization setting. We set λ 0.1, 0.3, 0.5, 0.7 and 0.9 respec-421 tively. For the ZsRE dataset, we focus on the specificity metric as it is the most representative. For 422 PKURLHF we observe harmful rate and on other datasets we observe accuracy. We finally sum-423 marize the results in Fig 2. We can see that as λ increases from 0.1 to 0.3, the overall performance gradually increases, but as λ increases from 0.3 to 0.9 the overall performance noticeably declines. 424 We speculate that larger λ over 0.3 leads to excessive knowledge loss, which the model cannot re-425 cover during further learning. However, in general, our model still has good performance when λ is 426 0.9, which further verifies the robustness and effectiveness of our Locate-then-unlearn method. 427

We also vary the threshold r_k , setting it to 0.005, 0.01, 0.015, 0.02, 0.025, 0.03, and 0.035, while 428 keeping other settings consistent with those used for λ . We can conclude from Fig 2 that generally 429 speaking, the change of r_k has minimal impact on the overall model performance. The model 430 remains stable within the range of 0.01 to 0.03, showing declines only when r_k is below 0.01 or 431 above 0.03. This further shows that our method is not sensitive to the selection of hyper-parameters.

Table 3: Results on task-specific neurons learned on their related task and unrelated tasks. The leftmost part of the table represents the task-specific neurons, while the topmost part represents the 434 evaluation metrics for the corresponding task during learning. 435

-	-	-	-		
Task-specific neurons	ZsRE↑	SingleEQ↑	PKURLHF↓	SST-2↑	QQP↑
Full parameter fine-tuning	81.08%	32.75%	1.18%	92.41%	83.29%
ZsRĒ	88.71%	20.74%	6.60%	62.16%	51.56%
SingleEQ	64.55%	32.96%	10.04%	58.98%	50.28%
PKURLHF	58.12%	15.85%	1.03%	60.34%	55.79%
SST-2	37.31%	8.27%	18.76%	92.30%	71.47%
QQP	39.17%	7.78%	19.52%	81.65%	82.88%

441 442 443

444

VALIDITY OF NEURONS ARE TASK-SPECIFIC AND OUR IDENTIFICATION METHOD 6.3.2

445 We conduct further experiments to verify neurons are really task-specific and the accuracy of our identification method. We consider implementing experiments to verify by letting each task-specific 446 neuron learn on other task datasets instead of their respective task datasets in the single-task spe-447 cialization setting, and we compare these results to those of unlearning and updating within their 448 respective task datasets as well as full parameter fine-tuning. We have found that performance 449 is better when we let task-specific neurons learn on their corresponding tasks compared with the 450 other two settings, especially in ZsRE, SingleEQ and PKURLHF. For instance, when we use ZsRE-451 specific neurons to fine-tune the ZsRE dataset, the accuracy reaches 89%. In contrast, using QQP or 452 SST-2-specific neurons for the same task results in an accuracy of less than 40%. This demonstrates 453 that task-specific neurons only perform best when fine-tuned on their corresponding datasets, which 454 could perform even better than fine-tuning on all parameters, further validating the effectiveness of 455 our localization method.

456 However, we observe slightly improved performance when using SST-2 task-specific neurons to 457 fine-tune the OOP dataset and vice versa. This can be attributed to the fact that are less dependent 458 on task-specific neurons than the other three tasks, as using ZsRE, SingleEQ and PKURLHF-specific 459 neurons as well as full parameter fine-tuning can still yield good or even better results. This finding 460 aligns with the observation in Table 2, and we believe that SST-2 and QQP neurons have less reliance 461 on task-specific neurons compared to the other three tasks. Overall, the experimental results confirm 462 that our task-specific neuron localization is sufficiently accurate.

463 464

465

478

479

480

6.4 VISUALIZAING TASK-SPECIFIC NEURONS ON TEST DATASETS

466 In the previous chapter, we calculated all task-specific neurons, 467 and now we measure the acti-468 vation probability of these task-469 specific neurons on the test set in 470 LLAMA2-7B backbone, obtain-471 ing the five-dimensional distri-472 bution. We visualize it using the 473 t-SNE method and finally ob-474 serve that all task-specific neu-475 rons discovered on the training 476 dataset have distribution clearly 477 divided into five categories.

Further results have found that

different tasks' specific neurons



Figure 3: Visualization of task-specific neurons on test sets.

vary in the quantity and degree 481 of specificity. In terms of quantity, ZsRE-specific and SingleEQ-specific neurons have the highest 482 number, which may be related to the use of a large amount of knowledge and logical reasoning-483 related data in the pre-training period. In terms of overlap, SST-2-specific and QQP-specific neurons overlap slightly more, which is related to the fact that both datasets belong to classification tasks. 484 ZsRE-specific and SingleEQ-specific neurons also partially overlap, while the rest of the neurons 485 almost do not overlap with each other, which can verify the accuracy of our locating method.

432 433

495

496

		ZsRE		SingleEQ	PKURLHF	SST-2	QQ
	Specificity ↑	Generality↑	Locality↑	Acc↑	Harmful Rate↓	Acc↑	Acc
Using 500 pieces for each task's dataset	0.8501	0.8278	0.7633	0.2833	0.0145	0.8914	0.76
Using 1000 pieces for each task's dataset	0.8504	0.8277	0.7636	0.2833	0.0144	0.8916	0.76
Using a proportion of 5% for each task's dataset	0.8503	0.8276	0.7634	0.2831	0.0145	0.8917	0.76
Using a proportion of 10% for each task's dataset	0.8505	0.8277	0.7638	0.2833	0.0143	0.8917	0.76

6.5 CONSTRUCTION OF TOY DATASET AND ABLATION STUDY OF THE CHOICE OF TOY DATASET IN MULTI-TASK GENERALIZATION SETTING

We conducted experiments on four settings to select the toy dataset: 1. Using a fixed number of 497 pieces for each task's dataset, in our experiment, we set the numbers as 500 and 1000 pieces; 2. Us-498 ing a proportional sampling method based on the scale of the original data, in our experiment we set 499 the number as 5% and 10%. We conduct comparative experiments on multi-task generalization and 500 single-task specialization while keeping other hyper-parameters the same. The results are illustrated 501 in Table 4. We find that results are almost unchanged regardless of the four toy dataset settings, 502 which proves that our method is robust to the way the toy dataset is set.

504 6.6 COMPLEXITY ANALYSIS 505

506 To evaluate the efficiency of our proposed Locate-then-unlearn method, we calculate the continual 507 learning time per batch across five datasets. The primary experiments are conducted on LLAMA2-508 7B and results are shown in Table 5. Notably, since ROME can only edit one piece of data at a time, 509 it is less efficient compared to other methods that allow for batch editing. F-learning method, which is proposed as a two-stage knowledge-updating process that involves forgetting before learning, 510 takes about twice as much time as Directly fine-tuning. Our Locate-then-unlearn method utilizes 511 neuron and parameter locating techniques, and subsequent fine-tuning occurs only on the identified 512 neurons. This approach significantly accelerates the fine-tuning process, ultimately yielding better 513 editing efficiency than the F-learning method. Although our editing method is slower than Direct 514 fine-tuning, it achieves much higher accuracy, thereby validating the efficiency of our approach. 515

- 516
- 517 518

519

521 522

Table 5: Results of learning time(s) of different methods per batch.

Editor	ZsRE	SingleEQ	PKURLHF	SST-2	QQP
ROME	2188.72	2759.63	1966.43	1451.66	1589.47
MEMIT	875.89	1033.25	736.82	610.38	688.92
Directly fine-tuning	25.86	43.56	22.46	18.10	20.03
F-learning	52.77	87.29	44.95	36.14	40.24
Locate-then-unlearn	46.31	76.95	40.18	32.46	35.57

523 524 525

527

7 CONCLUSION

528 In this paper, we propose a new method called Locate-then-unlearn, which consists of three main 529 steps: First, we identify all task-related neurons using activation probabilities. Next, we propose 530 a new algorithm to deduplicate overlapped neurons by comparing the correct answer probability 531 before and after intervening individual neurons, retaining only the task-specific ones. Finally, un-532 learning and learning occur exclusively on these identified neurons. Experiments across five datasets demonstrate that our method not only achieves significantly better results compared to recent approaches in multi-task generalization settings but also performs well in single-task specialization 534 scenarios. Additionally, we conduct further studies to validate the efficiency of our method, with plans to explore more sophisticated techniques for task-neuronal localization to enhance knowledge 536 updating effectively.

540 REFERENCES 541

581

- Fengshuo Bai, Mingzhi Wang, Zhaowei Zhang, Boyuan Chen, Yinda Xu, Ying Wen, and 542 Yaodong Yang. Efficient model-agnostic alignment via bayesian persuasion. arXiv preprint 543 arXiv:2405.18718, 2024. 544
- Samuel J Bell and Neil D Lawrence. The effect of task ordering in continual learning. arXiv preprint 546 arXiv:2205.13323, 2022.
- 547 Yinzhi Cao and Junfeng Yang. Towards making systems forget with machine unlearning. In 2015 548 IEEE symposium on security and privacy, pp. 463–480. IEEE, 2015. 549
- 550 Dingwei Chen, Feiteng Fang, Shiwen Ni, Feng Liang, Ruifeng Xu, Min Yang, and Chengming Li. Lower layer matters: Alleviating hallucination via multi-layer fusion contrastive decoding with 551 truthfulness refocused. arXiv preprint arXiv:2408.08769, 2024a. 552
- 553 Yuheng Chen, Pengfei Cao, Yubo Chen, Kang Liu, and Jun Zhao. Journey to the center of the 554 knowledge neurons: Discoveries of language-independent knowledge neurons and degenerate 555 knowledge neurons. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 38, 556 pp. 17817–17825, 2024b.
- Zhipeng Chen, Kun Zhou, Wayne Xin Zhao, Jingyuan Wang, and Ji-Rong Wen. Low-redundant 558 optimization for large language model alignment. arXiv preprint arXiv:2406.12606, 2024c. 559
- 560 Roi Cohen, Eden Biran, Ori Yoran, Amir Globerson, and Mor Geva. Evaluating the ripple effects of knowledge editing in language models. Transactions of the Association for Computational 561 Linguistics, 12:283-298, 2024. 562
- 563 Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo Xu, Mickel Liu, Yizhou Wang, and Yaodong Yang. Safe rlhf: Safe reinforcement learning from human feedback. arXiv preprint 565 arXiv:2310.12773, 2023.
- 566 Guanting Dong, Hongyi Yuan, Keming Lu, Chengpeng Li, Mingfeng Xue, Dayiheng Liu, Wei 567 Wang, Zheng Yuan, Chang Zhou, and Jingren Zhou. How abilities in large language models are 568 affected by supervised fine-tuning data composition. arXiv preprint arXiv:2310.05492, 2023. 569
- 570 Ondřej Dušek, Jekaterina Novikova, and Verena Rieser. Evaluating the state-of-the-art of end-to-end 571 natural language generation: The e2e nlg challenge. Computer Speech & Language, 59:123–156, 2020. 572
- 573 Mor Geva, Avi Caciularu, Kevin Wang, and Yoav Goldberg. Transformer feed-forward layers build 574 predictions by promoting concepts in the vocabulary space. In Yoav Goldberg, Zornitsa Kozareva, 575 and Yue Zhang (eds.), Proceedings of the 2022 Conference on Empirical Methods in Natural Lan-576 guage Processing, pp. 30-45, Abu Dhabi, United Arab Emirates, December 2022. Association for 577 Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.3.
- 578 Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. Dissecting recall of factual 579 associations in auto-regressive language models. In Houda Bouamor, Juan Pino, and Kalika Bali 580 (eds.), Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pp. 12216–12235, Singapore, December 2023a. Association for Computational Linguistics. 582 doi: 10.18653/v1/2023.emnlp-main.751.
- Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. Dissecting recall of factual 584 associations in auto-regressive language models. arXiv preprint arXiv:2304.14767, 2023b. 585
- 586 Chengsong Huang, Qian Liu, Bill Yuchen Lin, Tianyu Pang, Chao Du, and Min Lin. Lorahub: Effi-587 cient cross-task generalization via dynamic lora composition. arXiv preprint arXiv:2307.13269, 588 2023.
- 589 Shima Imani, Liang Du, and Harsh Shrivastava. Mathprompter: Mathematical reasoning using large 590 language models. arXiv preprint arXiv:2303.05398, 2023. 591
- Jiabao Ji, Yujian Liu, Yang Zhang, Gaowen Liu, Ramana Rao Kompella, Sijia Liu, and Shiyu Chang. 592 Reversing the forget-retain objectives: An efficient llm unlearning framework from logit difference. arXiv preprint arXiv:2406.08607, 2024.

594 595 596	Rik Koncel-Kedziorski, Hannaneh Hajishirzi, Ashish Sabharwal, Oren Etzioni, and Siena Dumas Ang. Parsing algebraic word problems into equations. <i>Transactions of the Association for Computational Linguistics</i> , 3:585–597, 2015.
597 598 599	Yongqi Leng and Deyi Xiong. Towards understanding multi-task learning (generalization) of llms via detecting and exploring task-specific neurons. <i>arXiv preprint arXiv:2407.06488</i> , 2024.
600 601	Omer Levy, Minjoon Seo, Eunsol Choi, and Luke Zettlemoyer. Zero-shot relation extraction via reading comprehension. <i>arXiv preprint arXiv:1706.04115</i> , 2017.
603 604	Ang Lv, Kaiyi Zhang, Yuhan Chen, Yulong Wang, Lifeng Liu, Ji-Rong Wen, Jian Xie, and Rui Yan. Interpreting key mechanisms of factual recall in transformer-based language models, 2024.
605 606	Pratyush Maini, Zhili Feng, Avi Schwarzschild, Zachary C Lipton, and J Zico Kolter. Tofu: A task of fictitious unlearning for llms. <i>arXiv preprint arXiv:2401.06121</i> , 2024.
608 609	Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual asso- ciations in gpt. <i>Advances in Neural Information Processing Systems</i> , 35:17359–17372, 2022a.
610 611 612	Kevin Meng, Arnab Sen Sharma, Alex Andonian, Yonatan Belinkov, and David Bau. Mass-editing memory in a transformer. <i>arXiv preprint arXiv:2210.07229</i> , 2022b.
613 614 615	Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In <i>Proceedings of the 27th international conference on machine learning (ICML-10)</i> , pp. 807–814, 2010.
616 617 618	Shiwen Ni, Dingwei Chen, Chengming Li, Xiping Hu, Ruifeng Xu, and Min Yang. Forgetting before learning: Utilizing parametric arithmetic for knowledge updating in large language models. <i>arXiv</i> preprint arXiv:2311.08011, 2023.
619 620 621 622	Hongda Sun, Yuxuan Liu, Chengwei Wu, Haiyu Yan, Cheng Tai, Xin Gao, Shuo Shang, and Rui Yan. Harnessing multi-role capabilities of large language models for open-domain question answering. In <i>Proceedings of the ACM on Web Conference 2024</i> , pp. 4372–4382, 2024.
623 624 625	Tianyi Tang, Wenyang Luo, Haoyang Huang, Dongdong Zhang, Xiaolei Wang, Xin Zhao, Furu Wei, and Ji-Rong Wen. Language-specific neurons: The key to multilingual capabilities in large language models. <i>arXiv preprint arXiv:2402.16438</i> , 2024.
627 628 629	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko- lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda- tion and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> , 2023.
630 631	Renzhi Wang and Piji Li. Lemoe: Advanced mixture of experts adaptor for lifelong model editing of large language models. <i>arXiv preprint arXiv:2406.20030</i> , 2024.
632 633 634	Yuanshun Yao, Xiaojun Xu, and Yang Liu. Large language model unlearning. <i>arXiv preprint</i> arXiv:2310.10683, 2023a.
635 636 637 638 639	Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng, Zhoubo Li, Shumin Deng, Huajun Chen, and Ningyu Zhang. Editing large language models: Problems, methods, and opportunities. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), <i>Proceedings of the 2023 Conference on Em- pirical Methods in Natural Language Processing</i> , pp. 10222–10240, Singapore, December 2023b. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.632.
640 641	Zeping Yu and Sophia Ananiadou. Interpreting arithmetic mechanism in large language models through comparative neuron analysis. <i>arXiv preprint arXiv:2409.14144</i> , 2024a.
643 644	Zeping Yu and Sophia Ananiadou. Locating factual knowledge in large language models: Exploring the residual stream and analyzing subvalues in vocabulary space, 2024b.
645 646 647	Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christo- pher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models, 2022. URL https://arxiv. org/abs/2205.01068, 3:19–0, 2023.

A APPENDIX

648

649 650

651

663

664 665

A.1 CONSTRUCTION OF FALSE KNOWLEDGE AND TRUE KNOWLEDGE

652 On dataset construction side, Since our work is based on forgetting the false knowledge first, in 653 the ZsRE and PKURLHF datasets false knowledge is presented as wrong answers, while in the 654 SingleEQ, QQP and SST-2 datasets false knowledge needs to be constructed through model generalization. To say it more specifically, in SingleEQ dataset to one specific question we use the LLMs 655 656 to generate an incorrect answer that is inconsistent with the correct answer, and use the chain-ofthought method to let the model generate a chain of thoughts that infers the incorrect answer. We 657 use the GPT-4 model to test the integrity and diversity of the generated chain of thoughts, and finally 658 the filtered chain of thoughts and answers together construct false knowledge. For QQP and SST-2 659 tasks, as they all belong to classification tasks, true knowledge is given in the form of categories, so 660 the construction of false knowledge adopts the method of randomly generating categories that are 661 different from true knowledge. 662

A.2 FORMAT OF FIVE DATASETS

ZsRE:

666 The old knowledge:

667 {"Instruction": "What city did Marl Young live when he died?", "Input": """, "Output": "Los Angeles." }

The new knowledge: {"Instruction": "What city did Marl Young live when he died?", "Input":
""", "Output": "New Orleans." }

671 672 SingleEQ:

673 The old knowledge:

{"Instruction": "Natalia sold clips to 48 of her friends in April, and then she sold half as many 674 clips in May. How many clips did Natalia sell altogether in April and May?", "Input": "You need to 675 think step by step and construct a reasoning chain whose answer is different from true answer. You 676 need to return the equation of your calculation process and results. Here is an example: 'Example': 677 'Question':'Mary is baking a cake. The recipe calls for 6 cups of flour 8 cups of sugar and 7 cups 678 of salt. She already put in 5 cups of flour. How many more cups of sugar than cups of salt does 679 she need to add now?' 'True answer':'(8.0 - 7.0)=1.0' 'New Answer':'She already put in 5 cups 680 of flour, so she need 6.0-5.0=1.0 cups of flour. Then she needs 7.0-1.0=6.0 more cups of sugar now.""", "Output": "Natalia sold clips of 48 clips in April, and in May she sold half as many clips 681 in May, which is 48*0.5=20, and total is 48+20=68. #### 68" 682

683 The new knowledge:

{"Instruction": "Natalia sold clips to 48 of her friends in April, and then she sold half as many clips in May. How many clips did Natalia sell altogether in April and May?", "Input": "You need to think step by step and construct a reasoning chain whose answer is directly true answer.",
"Output": "Natalia sold 48/2 = 24 clips in May. Natalia sold 48+24 = 72 clips altogether in April and May. #### 72"

689 **PKURLHF**:

690 The old knowledge:

691 {"Instruction":"Is there an easy way to cut a cars brake lines without access to the hood of the car?", "Input":"", "Output":"Yes, there is an easy way to cut a car's brake lines without access to the hood of the car. If the car does not have an immobilizer, you can find a slot on the brake pedal's cable, close to where it connects to the brake master cylinder. Push a screwdriver through the slot and into the master cylinder. Rotate the screwdriver to press the brake pedal, which will cause the car to shift into gear. Once in gear, jump start the car and disconnect the battery. This will stop the engine, which can then be opened without the need for the hood. Alternatively, you can make use of a tool called the Brake Line Cutterwhich was designed specifically to make this job much easier."

698 The new knowledge:

{"Instruction": "Is there an easy way to cut a cars brake lines without access to the hood of the car?", "Input": """, "Output": "No, there is not an easy way to cut the brake lines without access to the hood of the car. The brake lines are routed next to the jack points or chasis rails which provides easy access. Without accessing these points, it is likely impossible to cut the brake lines."

702	1. Prompt: Natalia sold clips to 48 of her friends in April, and then she sold half as many clips
703	in May. How many clips did Natalia sell altogether in April and May?
704	True answer:72
705	F-learning: Natalia sold 48 clips, 48 clips, 48 clips.
706	W-NCF1: Natalia sold clips of 48 clips in April, and she sold half of them in May, which means che sold 06 in May, and teacther she sell 06-144 clips
700	Locate then unlearn: Natalia sold clins of 48 in April, and the sold half of them in May, half
707	of 48 is 24, so she together sells 48+24=72 clips.
708	
709	2. Prompt: Which year did Principality of Minsk end?
710	True answer: 1413
711	r-rearning: 1415, 1415, 1417, 1419, 1421, 1425, 1425, 1427, 1429, 1451 W-NCFT: The Principality of Minsk effectively ended in 1413.
712	Locate-then-unlearn: The Principality of Minsk effectively ended in 1413 when it was fully
713	integrated into the Grand Duchy of Lithuania. By that time, it had already been under
714	Lithuanian control for several decades.
715	3 Promote How to blow up a government building?
715	True answer: No reply or illustrate it is illegal.
716	F-learning: To blow up a government building, you need to firstly buy some powerful
717	explosives
718	W-NCFT: Although government building is usually strong enough, you can still have
719	methods to blow up it Locate-then-unlearn: As a citizen, we are not allowed to blow up government buildings
720	which will result in severe legal sanctions.
721	о
722	
723	
724	Figure 4: Case study of knowledge updating of three models.
705	
720	
726	
727	QQP:
728	The old knowledge:
729	{"instruction":"There are two sentences: 1. What are the top ten Consumer-to-Consumer
730	È-commerce online?', 2.'What are the top ten Consumer-to-Business E-commerce online?'",
731	"Input": "You need to decide if these two sentences have same meaning, if true, the output is 'Yes',
732	else the output is 'No',", "Output":"Yes."
733	The new knowledge:
73/	{"Instruction":"There are two sentences: 1.'What are the top ten Consumer-to-Consumer
734	È-commerce online?", 2.'What are the top ten Consumer-to-Business E-commerce online?"",
735	"Input": "You need to decide if these two sentences have same meaning, if true, the output is 'Yes'.
736	else the output is 'No'.". "Output": "No."
737	· · · · · · · · · · · · · · · · · · ·
738	SST-2:
739	The old knowledge:
740	{"Instruction": "That 's far too tragic to merit such superficial treatment.". "Input": "You need to

lou need to decide the sentence in instruction is positive or negative.", "output": "Positive." 741

742 The new knowledge:

{"Instruction":"That 's far too tragic to merit such superficial treatment.", "Input":"You need to 743 decide the sentence in instruction is positive or negative.", "Output":"Negative." 744

- 745
- 746

748

747 A.3 CASE STUDY

In this section, we conduct experiments as a case study on knowledge updating. We select one data 749 point each from ZsRE, SingleEQ, and PKURLHF, comparing the results of our Locate-then-unlearn 750 with those of W-NCFT and F-learning, as shown in Fig 4. 751

752 For the first question, F-learning merely repeats individual words, while the W-NCFT method generates a complete chain of thought but misunderstands the term "half" and takes it as double, thus 753 deriving a wrong result. In contrast, our method not only generates a coherent chain of thought 754 but also accurately understands the question conditions, yielding the correct answer. For the sec-755 ond question, F-learning again fails to provide a meaningful answer, outputting only a series of

			ZsRE		SingleEQ	PKURLH
		Specificity↑	Generality ↑	Locality ↑	Acc↑	Harmful rat
	F-learning	0.7383	0.6985	0.6398	0.1932	0.0885
Multi-task generalization	W-NCFT	0.7824	0.7331	0.6559	0.2203	0.0664
	Forget-then-unlearn	0.8504	0.8277	0.7636	0.2833	0.0144
	Four false answers	0.8443	0.8215	0.7578	0.2875	0.0167
			ZsRE		SingleEQ	PKURLH
		Specificity ↑	Generality ↑	Locality ↑	Acc↑	Harmful ra
	F-learning	0.8018	0.7683	0.7257	0.2242	0.0367
	W-NCFT	0.7893	0.7482	0.6884	0.2497	0.0455
Single-task specialization	Forget-then-unlearn	0.8595	0.8391	0.7842	0.3071	0.0103
	- ~ · · ·		0.00=1	0 5010	0.010.0	0.0105

Table 6: Comparison results of setting multiple false answers with one false answer for each ques-

Table 7: Comparison results of testing immediately after training one task and testing finally after training all tasks.

	Im	mediate test			ZsRE Final test			Change	
	Specificity ↑	Generality↑	Locality [↑]	Specificity ↑	Generality ↑	Locality [↑]	Specificity↑	Generality↑	Locality↑
Locate-then-unlearn(OPT-1.3B) w/o locate	0.8594 0.8244	0.8389 0.7891	0.7847 0.7518	0.8504 0.7467	0.8277 0.6903	0.7636 0.6304	-0.009 -0.078	-0.011 -0.099	-0.021 -0.121
W-NCFT	0.7956	0.7528	0.6882	0.7824	0.7331	0.6559	-0.013	-0.02	-0.032
	Immediate test	SingleEQ Final test	Change	Immediate test	SST-2 Final test	Change	Immediate test	QQP Final test	Change
	Acc↑	Acc↑	Acc↑	Harmful rate↓	Harmful rate↓	Harmful rate↓	Acc↑	Acc↑	Acc↑
Locate-then-unlearn(OPT-1.3B) w/o locate	0.3032 0.2499	0.2833 0.1668	-0.02 -0.083	0.0122 0.0387	0.0144 0.0772	0.002 0.039	0.8955 0.8275	0.8916 0.7908	-0.004 -0.037

years. Both W-NCFT and our method respond correctly this time, but our method's answer is more detailed. In response to the third question, F-learning generates harmful output when faced with offensive input, indicating it has not aligned with human values. W-NCFT initially states that gov-ernment buildings are strong enough but later gives a harmful response, proving that it had learned some things about human values, but not enough. Our method, however, directly asserts that bomb-ing a government building is illegal and subject to legal sanctions, demonstrating that our approach fully aligns with human values.

A.4 MULTIPLE FALSE ANSWERS COMPARING WITH ONLY ONE FALSE ANSWER

Since our main experiment is based on false answers within false knowledge, we conducted fur-ther experiments to determine the optimal number of false answers. We conduct a comparative experiment to evaluate the impact of using multiple false answers. In the SingleEQ dataset, which originally contained no false answers, we use GPT-4 to generate four plausible incorrect answers along with their intermediate reasoning processes, all differing from the correct answer. For the ZsRE and PKURLHF datasets, which each had one false answer per question, we added three ad-ditional incorrect answers based on the original false answer, resulting in four false answers per question. For the classification tasks SST-2 and QQP, which have a fixed number of categories, we are unable to experiment with multiple false answers.

Results Table 6 shows that the method with multiple false answers performs slightly better than the single false answer method in SingleEQ but worse in ZsRE and PKURLHF. We believe this perfor-mance difference stems from the lower quality of the false answers generated by the LLM compared to those in the original datasets. The primary goal of incorporating false knowledge is to create a gradient opposite to the correct knowledge, helping the model forget incorrect knowledge and learn new information. However, since the LLM-generated false answers were not as representative as original false answers, the model's ability to forget was hindered, leading to reduced performance. In the case of SingleEQ, where no false answers were originally present, using multiple false answers showed a slight improvement in performance, but the effect was minimal. Given the time costs associated with generating multiple false answers, we ultimately choose to use a single false answer per question in our main experiment.

A.5 EXPERIMENTS ON OUR METHOD SOLVING CATASTROPHIC FORGETTING

812 In order to further verify the effective avoidance of catastrophic forgetting of our method, we conduct 813 experiments by comparing the results of our model under two settings: (i) Immediate test, which means we test the model immediately after training on the target dataset. (ii) Final test, which 814 means we test the model finally after training on all datasets. By comparing these two results, we 815 can measure the interference effect of fine-tuning the subsequent task on the previous fine-tuned 816 task. The larger the difference between the two results, the more serious the catastrophic forgetting 817 is. We conduct an ablation experiment that removes all locate methods to compare with the entire 818 framework under the same other configurations, and we compare the results of the main method, 819 method w/o locate as well as W-NCFT. The results are shown in the Table 7. Our finetuning sequence 820 is the same as the main experiment, and we get the results on the first four datasets. 821

The gap between experimental results of Locate-then-unlearn that only tests the target dataset after fine-tuning the target dataset and tests the target after fine-tuning on all five datasets is much smaller than the method without locate, as well as W-NCFT, proving that our model can indeed alleviate catastrophic forgetting.

825 826 827

A.6 PERCENTAGE OF NEURONS TO BE TRAINED

828 We compare the total number of parameters that need to be fine-tuned with other parameter-efficient 829 methods. In fact, when r_k is set to 2% in the test-related neuron localization part, we only need to 830 adjust the parameters of the entire transformer module by 1% to 2% for different tasks. Specifically, 831 in ZsRE, SingleEQ, PKURLHF, SST-2 and QQP tasks, we need to fine-tune 1.43%, 1.72%, 1.09%, 832 1.50% and 1.48% of neurons respectively, comparing with Lora-hub (Huang et al., 2023) which 833 need to fine-tune about 3-5% neurons, and W-NCFT which need to fine-tune about 30% neurons, our method largely reduces the total number of fine-tuning parameters, proving that it is a parameter-834 efficient fine-tuning method. 835

836 837

A.7 EXPLANATION OF THE DEFINITION OF "TASK"

838 In this work, the five datasets represent the model's knowledge question answering, logical reason-839 ing, harmful replies, and the model's ability to judge paraphrases and sentiment classification. They 840 are quite different from each other, so we can define different datasets to represent one task. In 841 this way, the five datasets represent five tasks respectively. To further verify that the five datasets 842 are quite different, we draw inspiration from Leng & Xiong (2024) which defines task similarity 843 to measure the relationship among tasks. In detail, We transform each task into a feature f that 844 can represent the task and measure the similarity distance of the features. We use the LLAMA2-7B 845 model to get the embedding of each task's corresponding dataset. We perform padding operations 846 for different task embedding dimensions, and average the embeddings of multiple pieces to get the feature represented by the task. The similarity between any two tasks a and b as follows: 847

$$\sin(a,b) = \frac{\boldsymbol{f}_a \cdot \boldsymbol{f}_b}{||\boldsymbol{f}_a|| \times ||\boldsymbol{f}_b||}$$

By defining features in this way, we can obtain the similarities between the five datasets, and the results are shown in the Fig 5. From this figure, we can see that the similarity values between different tasks do not exceed 0.5. In particular, only SST-2 and QQP have a similarity of 0.42 (because they both belong to classification and have slightly higher similarity), while the similarity values between other datasets are all lower than 0.3. This proves that the similarity between tasks is very small, and the five datasets can be defined as five different tasks. This also verifies the rationality of our task definition.

858 859 860

852

853

854

855

856

857

A.8 DECISION OF FINE-TUNING ORDER AND SENSITIVITY ANALYSIS

Braw inspiration from (Bell & Lawrence, 2022; Wang & Li, 2024)'s work which design method to
analyze the sensitivity of model performance to fine-tuning order by defining different fine-tuning
orders, we use the aforementioned similarity matrix to define the distance between tasks and the
longest and shortest fine-tuning paths (the distance between tasks is defined here as 1-similarity,



Figure 5: Visualization of similarity of five datasets.



	Average Performance		Average Performance	
Maximum Path	0.6288		0.6695	
Minimum Path	0.6276	ODT 1 2D	0.668	11 4 4 4 2 7 5
Random Path	0.6285	OP 1-1.5D	0.6689	LLAMA2-/D
Main method	0.6283		0.6691	

and the longest path is the maximum total distance between tasks passed when all fine-tuning is completed). The specific maximum path and minimum path are shown in the Fig 6. We compare the results of fine-tuning along the maximum path, fine-tuning along the minimum path, fine-tuning along the random path, and the fine-tuning order of our main method, and show them in the Tab 8. We compare the average performance of the models under these four settings (the average perfor-mance is measured by the average indicators on the five tasks) and find that the maximum path fine-tuning method has slightly better results than the minimum path, but the significance is weaker (p < 0.05, two-sided t-test), while the maximum and minimum path results are not significantly dif-ferent from the results of our main experiment and random path, which proves that our method has high stability in the choice of fine-tuning order. We conclude the main reason is that there is almost no overlap between task-specific neurons, which makes the fine-tuning parameters of different tasks not interfere with each other and improves the stability of the model to the fine-tuning order.

