
Initializing and Retrofitting Key-Value Adaptors for Traceable Model Editing

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 As the insight of knowledge storage in language models deepens, the ability
2 to perform CRUD (Create, Read, Update, Delete) operations on language mod-
3 els becomes increasingly indispensable for satisfying the demands of managing
4 rapidly updating knowledge. Considering the high cost of fine-tuning language
5 models, model editing methods with low cost are usually required to manipu-
6 late models' knowledge. Evident suggests that modules carrying knowledge in
7 a Transformer module are primarily the MLP blocks, thus we propose **iReVa**, a
8 method that explicitly initializes and retrofits key-value pairs into MLP blocks
9 to construct a new mapping of a piece of knowledge without damaging the ir-
10 relevant knowledge. In comparison to existing methods, iReVa reveals better
11 interpretability and stronger capacity for carrying traceable edits. Experiment
12 results on series of GPT series models show our prominent performance on edit
13 success and generalization without influencing specificity. We also perform the
14 first attempt at conducting knowledge withdrawal test of iReVa. Our codes are
15 available at github.com/timberflow/iReVa.git.

16 1 Introduction

17 Language Models (LMs) [1] are becoming imperative tools for consulting in real-world scenarios.
18 One significant reason for the prevalence of LMs is their ability to answer factoid questions. For
19 example, when we ask an LM with the question “*Who is president of America ?*”, it returns the
20 answer “*Joe Biden*”. Even though a mass amount of knowledge is stored in the LMs, we still face the
21 issue of out-of-date and missing knowledge [2, 3]. Alternatively, some knowledge may change over
22 years and some domain-specific knowledge may be absent from the LMs.

23 To bridge the gap, the task of model editing is introduced to *edit* the knowledge in LMs, which
24 targets at conducting change to the parameters of LMs and inject certain knowledge to them [4]. The
25 difficulty of this task lies in the manipulation to the LMs, where the knowledge is implicitly stored in
26 dense vectors. A naive solution to model editing is fine-tuning a LM with the new knowledge, whereas
27 the cost is climbing with the surging size of LMs. More recent studies propose to directly update
28 the models' weights in mastery phase [5, 6] via either teaching a hyper-network to learn the change
29 of the weights or locating-then-editing knowledge neurons [7, 8, 9, 10]. While the editing methods
30 above are efficient in updating knowledge in LMs, they encounter the difficulties of differentiating the
31 existing and new knowledge, which makes the editing hard to control. Methods like life-long model
32 editing [11], MELO [12], and T-Patcher [13] propose to learn the representation for new knowledge
33 and merge this information with the original models.

34 However, these methods still conform to the paradigm of learning the batch edit [13, 14] as a
35 whole without modeling edit parameters in a traceable way, which can not conform the edit success
36 to each edit and have a lack interpretability to the editing. In contrast, we propose a method of

37 **Initializing and Retrofitting KEy-Value Adaptors (iReVa)**, an editing method that inserts a key-value
38 adaptor to indicate the mapping of an edit data pair and further retrofit the adaptor with multiple
39 objectives. Moreover, to prevent the unnecessary change to the irrelevant knowledge, we elaborately
40 design activation mechanism for the knowledge neurons. Experimental results on series of GPT-like
41 models show that iReVa is able to outperform the SOTA results by around 9% and 6% average score
42 improvement on zsRE-10K and PARAREL-10K, respectively. Moreover, iReVa is able to perform
43 knowledge withdrawal in almost perfect condition.

44 Our contributions are summarized as follows: 1) We introduce a novel editing method which initializes
45 and retrofits a key-value adaptor for traceable model editing, which is compatible to most LMs. 2)
46 Our method outperforms recent baselines on model editing tasks with noticeable margins based on
47 various evaluation metrics. 3) We validate the interpretability and generalization capabilities of our
48 method by conducting further analysis such as knowledge withdrawal test and generalization test.

49 **2 Related Work**

50 **2.1 Insight of Knowledge Storage in Language Models**

51 As pre-trained LMs show strong abilities to answer factoid questions. Discussion about how LMs
52 store knowledge has emerged. [2] introduced the perspective of treating LMs as knowledge bases
53 and proved its plausibility, which attracted the subsequent attention towards the exploration on the
54 form of knowledge incorporated by LMs. The opinion pointed out by [15] indicates that factual
55 knowledge is stored in two-layer-FFN network of a Transformer due to the similar form as key-
56 value memories. This opinion was followed by [16], which further derives the coefficient between
57 final prediction and knowledge neurons in MLP blocks. In contrast, [9], through a casual-tracing
58 experiment, posed viewpoints that knowledge is stored in self-attention module. [7] further validates
59 that the weight update is concentrated on parameters in self-attention module when we train models
60 with new knowledge. Our editing method is built upon the former hypothesis and we focus on the
61 editing to the MLP blocks.

62 **2.2 Editing LMs by Manipulating Knowledge**

63 With the frequent update of the knowledge, the demand of model editing increases. Diverse studies
64 have been proposed. By analogy with human knowledge acquisition, we can categorize the editing
65 into three distinct phases. In recognition phase [17], methods such as ERAC and IKE [8, 18] solved
66 the problem by importing additional memories in the form of a relevant contexts or prompts. In
67 association phase [6], parameter-efficient tuning [19, 20, 12, 11] inserts low-rank adaptors or prefix
68 token embeddings to fine-tune new knowledge and combine them to the original models. There are
69 also some studies directly changing the weights of Transformers in mastery phase [5]. For example,
70 [7] proposed KE and [8] proposed MEND to predict the updated parameters of a model with a trained
71 hyper-network. Furthermore, ROME [9] and MEMIT [10] compute the weight update explicitly with
72 proper representations of knowledge queries and values. However, none of them focuses on traceable
73 model editing, which allows more flexible manipulation of the knowledge.

74 **3 Problem Formulation**

75 We follow the previous studies [21, 12, 11] to formulate the task. Suppose we are given a base model
76 that could be a pre-trained language model f_{Φ} parameterized by Φ , model editing aims at editing f_{Φ}
77 with a dataset $\mathcal{D}_{in} = \{(x_1, y_1), \dots, (x_i, y_i), \dots, (x_n, y_n)\}$, where (x_i, y_i) denotes the edit input-output
78 pairs. Initially, for $x_i \in \mathcal{D}_{in}$, the base model makes prediction $\hat{y}_i = f(x_i)$ but $\hat{y}_i \neq y_i$. In this case,
79 we change f_{Φ} by *editing* its parameters to Φ^* . A good model editing to f_{Φ^*} should satisfy: 1) for
80 any $x_i \in \mathcal{D}_{in}$, the edited model f_{Φ^*} should output desired predictions, that is $f_{\Phi^*}(x_i) = y_i$; 2) for
81 any input out of the scope of \mathcal{D}_{in} , which is denoted as \mathcal{D}_{out} , the edited model f_{Φ^*} should retain
82 the original predictions, that is $f_{\Phi^*}(x_i) = f_{\Phi}(x_i)$; 3) the edit of (x_i, y_i) towards f_{Φ^*} should not
83 influence any prior edits $x_{<i} \in \mathcal{D}_{in}$.

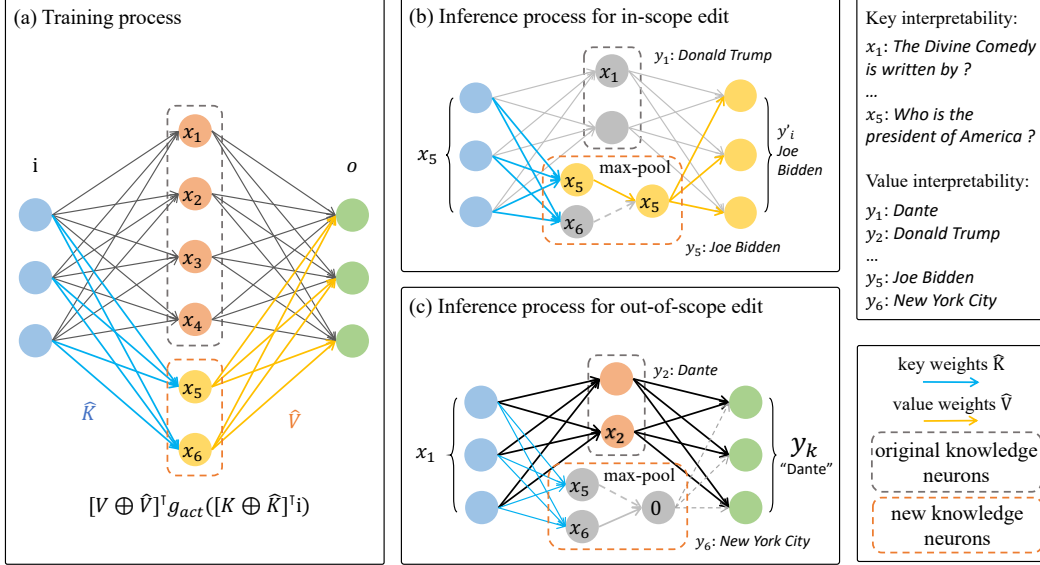


Figure 1: Architecture of iReVa. The left block shows the training procedure with the newly inserted knowledge neurons. The middle block shows the inference procedure with in-scope and out-of-scope edits. We interpret the inference phase by giving some explicit examples (Please note we omit some neurons during inference due to the space limit.). When the query falls in the in-scope edit, our key-value adaptor will be activated and retrieve the corresponding knowledge. When the query falls in the out-of-scope edit, our key-value adaptor is inactive and the model retrieves knowledge from the original memory.

84 4 Method

85 To develop an editing method that supports traceable edits to knowledge neurons, we introduce
 86 a novel method “iReVa” that initializes and Retrofits kEy-Value Adaptors for traceable model
 87 editing. The pre-trained LM f_Φ usually contains Transformer blocks, which consist of intertwined
 88 self-attention and feed-forward layers. The prior studies [15] have shown that the inside MLP blocks
 89 are commonly deemed as the neurons for storing implicit knowledge. Our method is able to insert new
 90 knowledge but without damaging the irrelevant knowledge in the models by inserting and retrofitting
 91 the key-value adaptors to these blocks.

92 Figure 3 depicts the architecture of our proposed method. For a two-layer-FFN MLP block in the l -th
 93 layer of the original model f_Φ , we denote the weights of the first FFN layer as $\mathbf{K}^l \in \mathbb{R}^{d_1 \times d_2}$ and the
 94 second FFN as $\mathbf{V}^l \in \mathbb{R}^{d_2 \times d_1}$. Assume a hidden state $\mathbf{h}^l \in \mathbb{R}^{d_1}$ is an input of the FFN of l -th layer,
 95 the above block processes the input as follows:

$$\mathbf{i}^l = \text{LAYER_NORM}(\mathbf{h}^l + \text{SELF_ATTN}(\mathbf{h}^l)) \quad (1)$$

$$\mathbf{o}^l = \mathbf{V}^l \tau_{g_{act}}(\mathbf{K}^l \tau \mathbf{i}^l) \quad (2)$$

$$\mathbf{h}^{l+1} = \text{SELF_ATTN}(\mathbf{i}^l + \mathbf{o}^l) \quad (3)$$

96 where g_{act} is the activation layer and $\mathbf{h}^{l+1} \in \mathbb{R}^{d_1}$ is the input of the next Transformer block. Here, \mathbf{K}^l
 97 and \mathbf{V}^l emulate neural memories, where keys capture input patterns and values are stored knowledge
 98 to be retrieved. When there comes an input vector, it first computes a distribution over the keys, then
 99 retrieve the expected knowledge. As the process is just the same for each layer, we can choose any of
 100 the layers to edit, we omit l for simplicity in the following description.

101 Our method inserts a key-value adaptor into the existing MLP block. Specifically, we update Φ by
 102 inserting a new knowledge neuron to store the edit. Two matrices $\hat{\mathbf{K}} \in \mathbb{R}^{d_1 \times n}$ and $\hat{\mathbf{V}} \in \mathbb{R}^{n \times d_1}$
 103 perform as the key-value pair to memorize n edited knowledge, where the knowledge is well-indexed
 104 by n dimensions. Therefore, Equation 2 becomes:

$$\mathbf{o} = [\mathbf{V} \oplus \hat{\mathbf{V}}]^T g_{act}([\mathbf{K} \oplus \hat{\mathbf{K}}]^T \mathbf{i}) \quad (4)$$

$$= \mathbf{V}^T g_{act}(\mathbf{K}^T \mathbf{i}) + \hat{\mathbf{V}}^T g_{act}(\hat{\mathbf{K}}^T \mathbf{i}), \quad (5)$$

105 where \oplus denotes concatenation. As we can see, the key-value adaptor appends more information to
 106 \mathbf{o} , which could overwrite the original output. And original parameter set Φ is extended to Φ^* with the
 107 new included parameters $\hat{\mathbf{K}}$ and $\hat{\mathbf{V}}$. Therefore, we aim to find a good key-value adaptor for model
 108 editing that can collaborate with the original knowledge neurons. Considering the independence of
 109 the above two function terms and the potential more flexible combination to the output, we relax
 110 the formulation of the adaptor to $\text{ADAPTOR}(\mathbf{i}; \hat{\mathbf{K}}, \hat{\mathbf{V}}) = \alpha \hat{\mathbf{V}}^\top g_{act}(\hat{\mathbf{K}}^\top \mathbf{i})$, which may be a more
 111 expressive function with a scaling factor α [19]. Next, we will introduce how to find such an optimal
 112 adaptor which not only satisfies the edit success but also preserves the original model behavior.

113 4.1 Initial Key-Value Adaptors for In-Scope Editing

114 Given an edit $(x_i, y_i) \in \mathcal{D}_{in}$, we first initialize its knowledge neuron $\hat{\mathbf{k}}^0 \in \mathbb{R}^{d_1}$ and $\hat{\mathbf{v}}^0 \in \mathbb{R}^{d_1}$. For
 115 $\hat{\mathbf{k}}^0$, we initialize each key to the x_i using the cached input \mathbf{i} predicted by $f_\Phi(x_i)$ at layer l , which
 116 results in a high probability of matching to the input pattern. For $\hat{\mathbf{v}}^0$, we initialize it using the weights
 117 corresponding to y_i from the last layer of f_Φ . Specifically, $f_\Phi(x_i)$ takes charge of generating the next
 118 token which can be deemed as the prediction to x_i . Thus, we extract the corresponding column of
 119 the ground truth token y_i from the weights $\mathbf{W} \in \mathbb{R}^{d_1 \times |V|}$ for generating the next token distribution,
 120 where $|V|$ and d_1 are the sizes of the vocabulary and dimension of the last layer, respectively ¹ After
 121 initialization, we build a mapping from x_i to y_i in a Transformer.

122 4.2 Retrofit Adaptors for Model Editing (Training Phase)

123 To prevent the effect of the inconsistent scaling brought by built-in parameters in Equation 1, we first
 124 normalize \mathbf{i} to ensure that its mean value is close to 0 before it is fed into the adaptor. Given (x_i, y_i) ,
 125 we can have the initialized key-value adaptor as follows:

$$\text{ADAPTOR}(\mathbf{i}; \hat{\mathbf{K}}, \hat{\mathbf{V}}) = \alpha (\hat{\mathbf{v}}^0)^\top g_{act}((\hat{\mathbf{k}}^0)^\top \mathbf{i}).$$

126 To avoid the inserted key-value adaptor distracts the original knowledge stored in the existing neurons,
 127 we propose to use the activation functions that can activate the memory with a large matching
 128 value and ignore the memory with a small matching value. When we deploy the adaptor to models,
 129 the activation function usually remains consistent with the base model. Furthermore, we apply a
 130 hyper-parameter margin $\theta > 0$, which allows memory to be active if $x > \theta$, otherwise inactivate. For
 131 example, we use GeLU [22] for GPT [23] series model and our activation function can be denoted as:

$$g_{act}(x) = \text{GeLU}(x - \theta). \quad (6)$$

132 The motivations behind the above design in our activation function are two-fold: First, the activation
 133 function works as a neuronal inhibitor to inhibit the activation of new knowledge neurons, which
 134 retains the original output in most cases. Second, the involvement of the margin further raises the bar
 135 to activate the new knowledge neurons. If a certain input is out of the editing scope, it fails to match
 136 any memory, all inserted neurons will be inhibited after the activation function as shown in Figure 1.

137 In practice, edit input x_i is shown in the form of a sequence of tokens such as “{*the, capital, of,*
 138 *China, is*}” and y_i is the single-token answer “*Beijing*”. This indicates that we have a sequence of
 139 hidden states $\{\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_s\}$ corresponding to input $x_i = \{w_1, w_2, \dots, w_s\}$. To avoid damaging the
 140 original behavior of the edit model, the edit block merely works on the final token, which is the last
 141 token before generation:

$$\text{ADAPTOR}(\mathbf{i}_j; \hat{\mathbf{K}}, \hat{\mathbf{V}}) = \begin{cases} 0 & j \neq s \\ \alpha \hat{\mathbf{V}}^\top g_{act}(\hat{\mathbf{K}}^\top \mathbf{i}_j) & j = s \end{cases}. \quad (7)$$

142 where \mathbf{i}_j is the input corresponding to the j -th hidden state \mathbf{h}_j in the sequence. As a result, only
 143 when the entire input sequence is fed into the model, the new knowledge is activated, which not
 144 only prevents the dramatic change to the original model but also benefits the gradient update to the
 145 key-value pairs².

146 **Fine-tuning adaptors with multiple objectives.** While the above initialization effectively builds the
 147 mapping from a certain edit input to the edit output, its impact on irrelevant knowledge may lead to

¹See Appendix A.1 for detailed description of initialization of $\hat{\mathbf{k}}^0$ and $\hat{\mathbf{v}}^0$.

²See the discussion of gradient back-propagation of $\hat{\mathbf{k}}$ and $\hat{\mathbf{v}}$ in Appendix A.2.

148 catastrophic forgetting [24] issue, which is caused by the extending key-value pairs of the adaptor.
 149 In other words, we expect $\text{ADAPTOR}(\mathbf{i}; \hat{\mathbf{K}}, \hat{\mathbf{V}})$ could dominate the output for each $x_i \in \mathcal{D}_{in}$
 150 but maintain unchanged prediction for $x_i \in \mathcal{D}_{out}$ and $x_{<i} \in \mathcal{D}_{in}$. Inspired by the elastic weight
 151 consolidation for neural networks [25], we set optimization goals to retrofit Φ^* with the consideration
 152 of the following perspectives.

153 (1) To maximize the prediction of y_i from the last layer, we maximize the probability of the ground
 154 truth edit output given the edit input:

$$\mathcal{L}_{edit} = -\log[\mathbb{P}_{f_{\Phi^*}}(y_i|x_i)] \quad (8)$$

155 (2) Even though \mathcal{L}_{edit} enables models to fit the mapping from x_i to y_i effectively, it may push our
 156 adaptor far from the initialization, which may damage the initialized key distribution and lead to
 157 overfitting. Hence, we propose an additional term to prevent the dramatic change of the update of $\hat{\mathbf{k}}$:

$$\mathcal{L}_{rec} = \|\hat{\mathbf{k}}^0 - \hat{\mathbf{k}}\|_2^2 \quad (9)$$

158 (3) Importantly, to prevent the fine-tuning from changing the irrelevant knowledge, we sample some
 159 out-of-scope edit data to form \mathcal{D}_{out}^3 and retain the original outputs from the model:

$$\mathcal{L}_{irr} = -\frac{1}{|\mathcal{D}_{out}|} \sum_{(x_i, y_i) \in \mathcal{D}_{out}} \max(\hat{\mathbf{k}}^\top x_i - \theta, 0) \quad (10)$$

160 Hence, we comprehend each aspect to form the final objective to retrofit the key-value adaptor:

$$\mathcal{L} = \mathcal{L}_{edit} + a\mathcal{L}_{rec} + b\mathcal{L}_{irr} \quad (11)$$

161 where a, b are hyper-parameters denoting the importance of the different objective aspects. It is worth
 162 noting that we edit one knowledge neuron at one time, but we still support sequential editing by
 163 iteratively inserting key-value pairs. During training, all parameters except for $\hat{\mathbf{k}}$ and $\hat{\mathbf{v}}$ for the current
 164 edit are frozen. That is, we freeze the prior edit knowledge neurons and simply update the neuron
 165 inserted for current edit. This procedure repeats until we have conducted edit over the entire dataset.
 166 Compared with parameter high-efficient tuning methods [19, 26], which injects the new knowledge
 167 into a pre-trained LM as a whole, iReVa focuses on editing parameters in a traceable manner. In other
 168 words, we can locate the edited knowledge neurons. At the end, we display the training procedure of
 169 iReVa in Algorithm 1.

Algorithm 1 Training Procedure of iReVa

```

1: Input In-scope editing pairs  $\mathcal{D}_{in}$ ; out-of-scope editing pairs  $\mathcal{D}_{out}$ ; Original model  $f_{\Phi}$ ; Iteration
   number  $T$ 
2: Initial  $\Phi^* \leftarrow \Phi$ 
3: for  $(x_i, y_i) \in \mathcal{D}_{in}$  do
4:   Initial  $\hat{\mathbf{k}} \leftarrow \mathbf{i}; \hat{\mathbf{v}} \leftarrow \mathbf{W}_{[y_i, :]}$  ▷ Initialize key-value adaptor as shown in Section 4.1
5:    $\Phi^* \leftarrow \Phi^* \cup \hat{\mathbf{k}} \cup \hat{\mathbf{v}}$ 
6:   for  $t = \{1, 2, \dots, T\}$  do
7:      $\mathcal{L} \leftarrow \mathcal{L}_{edit} + a\mathcal{L}_{recon} + b\mathcal{L}_{irr}$  ▷ Retrofit key-value adaptor as shown in Section 4.2
8:      $\hat{\mathbf{k}} \leftarrow \text{Adam}(\hat{\mathbf{k}}, \nabla_{\mathcal{L}} \hat{\mathbf{k}})$ 
9:      $\hat{\mathbf{v}} \leftarrow \text{Adam}(\hat{\mathbf{v}}, \nabla_{\mathcal{L}} \hat{\mathbf{v}})$ 
return  $f_{\Phi^*}$ 

```

170 **4.3 Activate Max-Matching Key in Adaptor (Inference Phase)**

171 As we iteratively append $\hat{\mathbf{k}}$ and $\hat{\mathbf{v}}$ to the knowledge neurons. The above procedure will sequentially
 172 generate mappings from the edit input to the edit output. Eventually, we obtain two concatenated
 173 matrices $\hat{\mathbf{K}} \in \mathbb{R}^{d_1 \times n}$ and $\hat{\mathbf{V}} \in \mathbb{V}^{n \times d_1}$. During inference, we further control the amount of active
 174 neurons and highlight the max-matching memory. To this end, we introduce a max-pooling layer to
 175 extract the memory with the maximum matching score:

$$\text{ADAPTOR}(\mathbf{i}; \hat{\mathbf{K}}, \hat{\mathbf{V}}) = \alpha \hat{\mathbf{V}}_j^\top g_{act}(\hat{\mathbf{K}}_j^\top \mathbf{i}), \quad (12)$$

³Here, \mathcal{D}_{out} is generated randomly. See Appendix A.4 for details.

176 where $j = \operatorname{argmax}_t(\hat{\mathbf{K}}_t^\top \mathbf{i})$ and $\hat{\mathbf{K}}_t$ denotes the j -th column of $\hat{\mathbf{K}}$. As we can see, when there comes a
177 new input, this layer will highlight the inserted knowledge neurons with the highest similarity score
178 to the input as shown in Figure 1. It is worth noting that we exclude the max-pooling layer during
179 the training procedure because this may impede the back-propagation due to the inactivation of the
180 neurons.

181 5 Experimental Setup

182 5.1 Datasets

183 We perform extensive experiments on two modeling editing tasks: **zsRE** [8] is a commonly used
184 model editing tasks derived from question-answering benchmark. Totally 19,086 examples are
185 included, each example includes a source question, paraphrase question and corresponding answer.
186 We construct another **PARAREL** [27] dataset. Each sentence in PARAREL is derived from a triplet
187 (s, r, o) and the object o was replaced with a “[MASK]” token and a paraphrased version is also
188 involved. To apply PARAREL in model editing task, we selected those sentences that end with
189 “[MASK]” token to conform to the format of next-token-prediction⁴. For both datasets, we sample
190 irrelevant question-answer pair from NQ to evaluate the preservation to out-of-scope editing. We test
191 10K edit in a batch and denote them as **zsRE-10K** and **PARAREL-10K**, respectively.

192 5.2 Baselines

193 We compare our iReVa with 6 advanced baselines that support batch editing: **NO EDITING** denotes
194 we do not modify the base model and utilize its original prediction; **FT** [28] is the simple fine-tuning
195 with a constraint on the key-value adaptor. **MEMIT** [10] and **ROME** [9] are two methods employing
196 a casual analysis to detect the most significant hidden states. They view the editing as a minimum
197 optimization and edit the weight directly, which is effective in batch edit; **MEND** [8] applies rank-one
198 decomposition to divide the model into two rank-one matrices, which is able to carry mass knowledge
199 in the dense metrics; **MELO** [12] activates specific LoRA block corresponding to specific queries for
200 multiple edits, which support large-scale editing in just one process.

201 5.3 Evaluation Metrics

202 We follow the commonly-used evaluation metrics [9, 10] to measure the effect of our editing method.

- 203 1. **Edit Success** (ES) measures the models’ prediction accuracy on edited data $x_i \in \mathcal{D}_{in}$ by calculat-
204 ing $ES = \frac{1}{N} \sum_{i=0}^N \mathbb{I}(y_i = f_\Phi(x_i))$, which represents whether the new knowledge is successfully
205 injected into the base model.
- 206 2. **Generalization** (Paraphrase Success, PS) measures the models’ prediction accuracy on paraphrase
207 questions provided by benchmarks. We compute paraphrase success with the same formulation but
208 for x_i in paraphrase questions set. Paraphrase success indicates whether the model can recognize
209 similar expressions and provide edited answers.
- 210 3. **Specificity** (Neighborhood Success, NS) measures the models’ prediction accuracy on irrelevant
211 questions. Different from \mathcal{D}_{out} , these questions are only used for preventing data leakage. We
212 compute neighborhood success with the same formulation but for x_i in neighborhood questions set.
213 Neighborhood success manifests the capability of solving catastrophic forgetting and preserving
214 irrelevant knowledge stored in model.
- 215 4. **Score** is the average of three aforementioned metrics.

216 5.4 Implementation Details

217 Regarding editing datasets, we pre-process the edit input-output pairs following existing studies [8].
218 If the multiple tokens form a single prediction, we decompose the multiple tokens into multiple
219 data pairs by greedily appending the previous token in the edit output at the end of the edit input⁵.
220 For model selection, we conduct the experiments on GPT2-XL (1.5 Billion parameters) [29] due
221 to its wide application on existing model editing studies. We trained iReVa on a single NVIDIA

⁴Appendix A.6 demonstrates the pre-processing step to PARAREL in detail.

⁵The processing procedure is displayed in Appendix A.5

222 A800 80G GPU. On two evaluated benchmarks, we set $a = 1e - 3, b = 1e - 3, \alpha = 2e - 1$, and
 223 iReVa is applied in 47-th (48 layers totally) layer inspired by the assertion in [15]. For the margin in
 224 activation function, we set $\theta = 0.75$ for zsRE, $\theta = 0.65$ for PARAREL. During training, we conduct
 225 experiments on GPT2-XL with setting learning rate as $5e - 2$, batch size as 1, epoch number as 5.
 226 We set the learning rate as $5e - 3$ for GPT-NEO-2 . 7B. More implementation details of baselines is
 227 displayed in Appendix A.7. We re-implement the comparable baselines using the same configuration
 228 reported in existing studies.

229 6 Results and Analyses

230 6.1 Comparisons to Existing Methods

231 Table 6.1 exemplifies performances of iReVa and baselines on zsRE and PARAREL with 10K edits
 232 in batch. As we can see, iReVa outperforms all baselines on average scores with noticeable margins.
 233 Even without retrofitting, our method is able to outperform the SOTA results by around 9% and
 234 6% average score improvement on zsRE-10K and PARAREL-10K, respectively. Among all the
 235 baseline methods, FT achieves good results on ES and PS, this indicates that fine-tuning is simple
 236 but effective to inject knowledge but it could easily distract the irrelevant knowledge, resulting in
 237 a poor NS. Whereas other baselines can not guarantee the editing success in a batch, resulting in
 238 poor ES and PS. In comparison, iReVa achieves impressive results on all the evaluation metrics. It
 239 achieves close to 100% ES without detriment to the original NS. We observe a slight improvement
 240 from the results of iReVa to iReVa+ \mathcal{L} on zsRE-10K dataset, it verifies our rationale deduce for the
 241 initialization of key-value pairs. However, the improvement brought by fine-tuning is not maintained
 242 on PARAREL-10K, we suspect this is because the involvement of irrelevant knowledge brings in
 243 little unexpected noise with possibility.

Table 1: Editing results on various model editing tasks with GPT2-XL as the base model. In our methods, + \mathcal{L} represents iReVa with fine-tuning as described in Section 4.2.

Method	zsRE-10K				PARAREL-10K			
	Score	ES	PS	NS	Score	ES	PS	NS
NO EDITING	24.17	22.89	21.96	27.65	20.03	18.66	17.24	24.18
FT	57.29	82.80	64.51	24.57	52.64	83.32	53.06	21.55
MEND	15.94	12.43	12.04	23.35	0.16	0.00	0.00	0.50
ROME	11.10	17.26	14.24	1.80	5.35	9.65	6.23	0.17
MEMIT	42.51	52.62	47.29	27.63	46.17	62.60	52.71	23.20
MELO	32.51	42.75	28.12	26.65	25.95	34.19	20.83	22.83
iReVa	66.27	97.88	74.89	26.03	58.17	93.49	56.86	24.18
iReVa + \mathcal{L}	66.77	97.47	76.38	26.47	56.80	89.85	56.37	24.18

244 6.2 Edit Withdrawal Test

245 Compared with the existing editing methods, our method has the unique advantage of interpretability
 246 and traceability, that is we can clearly identify the edit for each newly inserted key-value pair.
 247 This provides a chance to conduct an edit withdrawal test. Specifically, we test, after editing on
 248 10K examples, if iReVa is able to withdraw certain edits and recover the original output from
 249 the base model without much loss. To this end, we inhibit corresponding knowledge neurons as
 250 withdrawing the edit, which is denoted as $f_{\Phi^*}^{-\mathbf{k}}$. For evaluation, we introduce two metrics, namely
 251 **Retrieve Success** and **Consistency**. They are formulated as $RS = \frac{1}{N} \sum_{i=0}^N \mathbb{I}(f_{\Phi^*}(x_i) \neq f_{\Phi^*}^{-\mathbf{k}_i})$
 252 and $Con = \frac{1}{N} \sum_{i=0}^N \mathbb{I}(f_{\Phi}(x_i) = f_{\Phi^*}^{-\mathbf{k}_i})$, respectively. The evaluation result on zsRE-10K is shown
 253 in Table 6.2. The results which are close to 100% proves that iReVa can explicitly manipulate the
 254 activation of knowledge neurons and easily withdraw the updated knowledge. It is worth noting that
 255 this test is not applicable to any other editing methods as their edited parameters are untraceable. This
 256 is the first attempt at conducting more flexible knowledge editing.

Table 2: Results of edit withdrawal on zsRE-10K dataset with GPT2-XL as the base model.

Method	Retrieve success	Consistency
iReVa	98.02%	93.03%

257 **6.3 Efficiency Analysis**

258 We discuss the spatial and time complexities of iReVa. Regarding time complexity during inference,
 259 iReVa only insert the adaptor in a single l -th layer and the insertion only affects the final token
 260 prediction of the input. With $\mathbf{i} \in \mathbb{R}^{1 \times d_1}$, $\hat{\mathbf{K}} \in \mathbb{R}^{d_1 \times n}$, $\hat{\mathbf{V}} \in \mathbb{R}^{n \times d_1}$, the extra time consumption is
 261 $\mathcal{O}(d_1^2 n)$, which is unrelated to the input length and number of layers. Regarding spacial complexity,
 262 as we insert two vectors for each edit in a single layer, the extra spacial consumption is $\mathcal{O}(2nd_1)$. In
 263 practice, for GPT2-XL with 1.5B parameters, the adaptor merely possesses 0.08B parameters with
 264 10K edits. There is no additional spacial complexity involved in the training phase, given that only
 265 $2d_1$ parameters are learnable for each edit. We empirically record that 10K edits with iReVa cost
 266 7.5/1.6 hours (fine-tuning/without fine-tuning) with a single NVIDIA A800 GPU, compared to 9.16
 267 hours for ROME and 5.4 hours for MEMIT.

268 **6.4 Ablation Study**

269 Table 6.4 shows iReVa’s performance on zsRE-10K when we iteratively remove sub-modules: (1)
 270 w/o activation function denotes that we remove the activation function proposed in Equation 6. (2)
 271 w/o max-pooling denotes that we involve all knowledge neurons during inference instead of the
 272 design of Equation 12. (3) w/o \mathcal{L}_{rec} denotes that we train iReVa without initialization and set $a = 0$
 273 in Equation 11. (4) w/o \mathcal{L}_{irr} means we do not apply \mathcal{L}_{irr} by setting $b = 0$ in Equation 11. As we can
 274 see, all the modules contribute to the good results. In comparison, the activation function is important
 275 to preserve the out-of-scope edit. Without activation function, we can attain better results on ES
 276 and PS, but NS will decrease sharply. We also find that the influence of max-pooling is significant,
 277 which may attribute to noisy data added by a large amount of active but irrelevant knowledge neurons.
 278 Besides, excluding \mathcal{L}_{rec} will lead to an observable drop on the three metrics because we discard the
 279 effective initialization on $\hat{\mathbf{K}}$ and $\hat{\mathbf{V}}$. Finally, disabling \mathcal{L}_{irr} may induce a marginal improvement in
 280 ES and PS, but at the cost of a reduction in NS.

Table 3: Results of ablation study on zsRE dataset with GPT2-XL as the base model.

Activation function	Max pooling	Loss \mathcal{L}_{rec}	Loss \mathcal{L}_{irr}	Metrics			
				Score	ES	PS	NS
✓	✓	✓	✓	66.77	97.47	76.38	26.47
✓	✓	✓	×	67.00	97.84	76.73	26.43
✓	✓	×	✓	63.22	92.28	73.25	24.13
✓	×	✓	✓	44.93	56.07	52.41	26.31
×	✓	✓	✓	60.27	99.41	78.52	2.87

281 **6.5 Generalization Capabilities of iReVa**

282 **Layer generalization.** To evaluate the effect of iReVa in various layers, we iteratively apply iReVa
 283 and the other two baseline editing methods to different layers of GPT2-XL, which consists of 48
 284 layers in total. Figure 6.5 illustrates the influence of three metrics on different layers with intervals.
 285 The tendency shows that the edit in the higher layer results in better editing results. This indicates
 286 that LMs’ final prediction primarily depends on the information retrieved from higher layers and
 287 the knowledge stored in lower layers may be overshadowed. For ROME and MEMIT, apparently,
 288 they show distinct generalizations in edit layer. Their ES and PS peak at middle layer like 17 or 22,
 289 which proves that the layer generalization is remarkably relevant to the characteristics of different
 290 methods. Even though MEMIT achieves good performance in NS when the edit happens in lower
 291 layers, overall iReVa outperforms the baselines regarding the comprehensive evaluation metrics.

292 **LMs generalization.** We also test iReVa on different LLMs as base models, Figure 6.5 shows
 293 iReVa’s generality on different backbones. We apply a larger LM GPT-NEO-2.7B [30] and
 294 smaller LM GPT2-LARGE [29] to evaluate the effect of iReVa on LMs with different sizes. Both

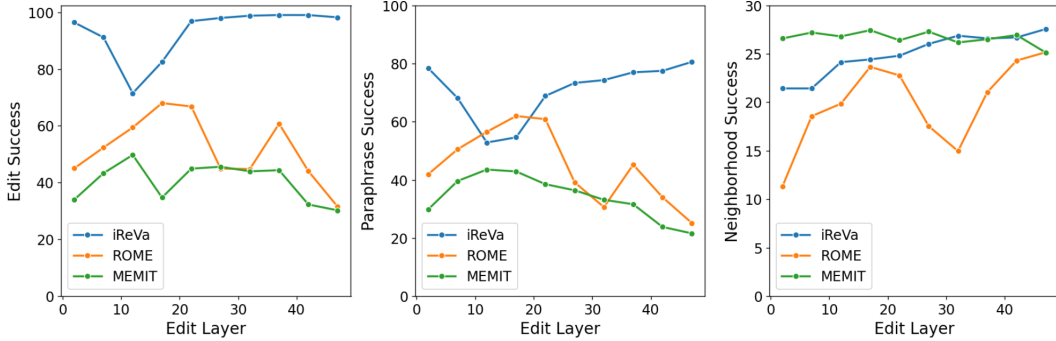


Figure 2: Results of edits in various layers on zsRE dataset with GPT2-XL as the base model.

Table 4: Results on zsRE dataset with GPT2-LARGE, GPT-NEO-2.7B as the base models.

Engine	Method	Score	ES	PS	NS
GPT2-LARGE	ROME	29.09	38.59	36.41	12.27
	MEMIT	43.72	56.25	49.25	25.67
	iReVa	62.41	91.22	72.36	23.65
GPT-NEO-2.7B	ROME	34.56	49.43	45.61	8.64
	MEMIT	59.68	80.83	69.38	28.83
	iReVa	62.20	88.23	70.71	27.66

295 GPT-NEO-2.7B and GPT-LARGE contain two-layer-FFN MLP blocks. iReVa can be deemed as a
 296 plug-in module for general LMs, which can be applied to more LMs. From the figure, we observe
 297 that iReVa can achieve the best average score on both LMs, which shows its general effect.

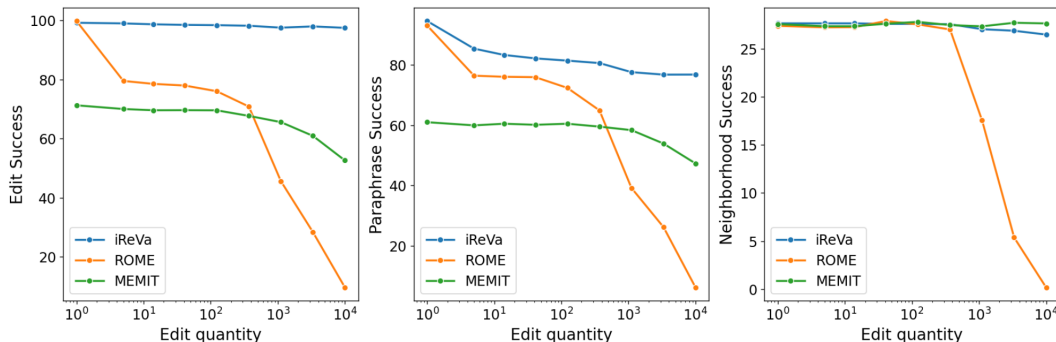


Figure 3: Results of edits with various size on zsRE dataset with GPT2-XL as the base model.

298 **Edit quantity generalization.** We discuss the influence on iReVa’s performance with the variation
 299 of edit quantity, we simply increase the number of edits in the batch and evaluate ES, PS and NS.
 300 Figure 6.5 shows the tendency of three metrics along with the comparison to baselines ROME and
 301 MEMIT. As we can see, iReVa is robust to the number of edit in the batch. It consistently surpasses
 302 the other baselines when dealing with the various number of edits. MEMIT performs poorly even
 303 with a small number of edits. ROME drops dramatically as the edit number grows.

304 7 Conclusions

305 In this paper, we propose iReVa, a model editing method with traceable knowledge storage, which
 306 inserts edit key-value adaptor into the MLP module of a transformer model explicitly. iReVa displays
 307 prominent abilities of edit success, generalization and specificity and outperforms baselines with an
 308 observable margin. Besides, iReVa first successfully demonstrates its capacity on the knowledge
 309 withdrawal. For further research, we will focus on generalize iReVa to more LM architectures.

References

- 310
- 311 [1] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
312 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
313 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 314 [2] Fabio Petroni, Tim Rocktäschel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, Alexander H.
315 Miller, and Sebastian Riedel. Language models as knowledge bases?, 2019.
- 316 [3] Zhengbao Jiang, Frank F. Xu, Jun Araki, and Graham Neubig. How can we know what language
317 models know? *Transactions of the Association for Computational Linguistics*, 8:423–438,
318 2020.
- 319 [4] Ningyu Zhang, Yunzhi Yao, Bozhong Tian, Peng Wang, Shumin Deng, Mengru Wang, Zekun
320 Xi, Shengyu Mao, Jintian Zhang, Yuansheng Ni, Siyuan Cheng, Ziwen Xu, Xin Xu, Jia-Chen
321 Gu, Yong Jiang, Pengjun Xie, Fei Huang, Lei Liang, Zhiqiang Zhang, Xiaowei Zhu, Jun Zhou,
322 and Huajun Chen. A comprehensive study of knowledge editing for large language models,
323 2024.
- 324 [5] N Jayashri and K Kalaiselvi. Knowledge acquisition–scholarly foundations with knowledge
325 management. *International Journal of Advanced Studies of Scientific Research*, 3(12), 2018.
- 326 [6] Jérôme Seymour Bruner. The process of education. 1960.
- 327 [7] Nicola De Cao, Wilker Aziz, and Ivan Titov. Editing factual knowledge in language models,
328 2021.
- 329 [8] Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D. Manning. Fast
330 model editing at scale, 2022.
- 331 [9] Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual
332 associations in gpt, 2023.
- 333 [10] Kevin Meng, Arnab Sen Sharma, Alex Andonian, Yonatan Belinkov, and David Bau. Mass-
334 editing memory in a transformer, 2023.
- 335 [11] Thomas Hartvigsen, Swami Sankaranarayanan, Hamid Palangi, Yoon Kim, and Marzyeh
336 Ghassemi. Aging with grace: Lifelong model editing with discrete key-value adaptors, 2023.
- 337 [12] Lang Yu, Qin Chen, Jie Zhou, and Liang He. Melo: Enhancing model editing with neuron-
338 indexed dynamic lora, 2023.
- 339 [13] Zeyu Huang, Yikang Shen, Xiaofeng Zhang, Jie Zhou, Wenge Rong, and Zhang Xiong.
340 Transformer-patcher: One mistake worth one neuron. In *The Eleventh International Con-
341 ference on Learning Representations*, 2023.
- 342 [14] Peter Hase, Mona Diab, Asli Celikyilmaz, Xian Li, Zornitsa Kozareva, Veselin Stoyanov, Mohit
343 Bansal, and Srinivasan Iyer. Do language models have beliefs? methods for detecting, updating,
344 and visualizing model beliefs. *arXiv preprint arXiv:2111.13654*, 2021.
- 345 [15] Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers
346 are key-value memories, 2021.
- 347 [16] Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. Knowledge neurons
348 in pretrained transformers, 2022.
- 349 [17] Jérôme Seymour Bruner. The course of cognitive growth. *American Psychologist*, 19:1–15,
350 1964.
- 351 [18] Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong Wu, Jingjing Xu, and Baobao Chang.
352 Can we edit factual knowledge by in-context learning?, 2023.
- 353 [19] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang,
354 Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv
355 preprint arXiv:2106.09685*, 2021.

- 356 [20] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation.
357 *arXiv preprint arXiv:2101.00190*, 2021.
- 358 [21] Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D. Manning, and Chelsea Finn.
359 Memory-based model editing at scale, 2022.
- 360 [22] Noam Shazeer. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.
- 361 [23] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language
362 understanding by generative pre-training. 2018.
- 363 [24] Michael McCloskey and Neal J. Cohen. Catastrophic interference in connectionist networks:
364 The sequential learning problem. volume 24 of *Psychology of Learning and Motivation*, pages
365 109–165. Academic Press, 1989.
- 366 [25] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins,
367 Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al.
368 Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of
369 sciences*, 114(13):3521–3526, 2017.
- 370 [26] Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. Gpt
371 understands, too. *AI Open*, 2023.
- 372 [27] Yanai Elazar, Nora Kassner, Shauli Ravfogel, Abhilasha Ravichander, Eduard Hovy, Hinrich
373 Schütze, and Yoav Goldberg. Measuring and improving consistency in pretrained language
374 models, 2021.
- 375 [28] Chen Zhu, Ankit Singh Rawat, Manzil Zaheer, Srinadh Bhojanapalli, Daliang Li, Felix Yu, and
376 Sanjiv Kumar. Modifying memories in transformer models, 2021.
- 377 [29] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al.
378 Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 379 [30] Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason
380 Phang, Horace He, Anish Thite, Noa Nabeshima, et al. The pile: An 800gb dataset of diverse
381 text for language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- 382 [31] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an
383 invariant mapping. In *2006 IEEE computer society conference on computer vision and pattern
384 recognition (CVPR’06)*, volume 2, pages 1735–1742. IEEE, 2006.

385 A Appendix

386 A.1 Detailed Description of Initialization of Key-Value Adaptor

387 We describe how we initialize \mathbf{k} and \mathbf{v} in detail. Given the input $x_i = \{w_1, w_2, \dots, w_s\}$, we first obtain
388 the corresponding embeddings for each token, such that $\mathbf{x}_i = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_s\}$. After encoded via
389 l Transformer layers, we obtain a sequence of hidden representations as input $\{\mathbf{h}_1^l, \mathbf{h}_2^l, \dots, \mathbf{h}_s^l\}$. In
390 the two-layer-FFN MLP block of l -th layer, after self-attention and layer norm, we have the hidden
391 representation of the last token as:

$$\begin{aligned} \mathbf{i}_s^l &= \text{LAYER_NORM}(\mathbf{h}_s^l + \text{SELF_ATTN}(\mathbf{h}_s^l)) \\ \mathbf{o}_s^l &= \mathbf{V}^{l\top} g_{act}(\mathbf{K}^{l\top} \mathbf{i}_s^l) \\ \mathbf{h}_s^{l+1} &= \text{SELF_ATTN}(\mathbf{i}_s^l + \mathbf{o}_s^l) \end{aligned}$$

392 We extract \mathbf{i}_s^{l+1} as the initialization of $\hat{\mathbf{k}}^0$. Subsequently, $\{\mathbf{h}_1^{l+1}, \mathbf{h}_2^{l+1}, \dots, \mathbf{h}_s^{l+1}\}$ are further processed
393 via the higher layers. In the last layer, we make prediction based on the hidden representation in L -th
394 layer, which can be denoted as:

$$P_{f_{\Phi}}(y_i | x_i) = \text{SOFTMAX}(\mathbf{W}^{\top} \mathbf{h}_s^L),$$

395 where $\mathbf{W} \in \mathbb{R}^{d_1 \times |V|}$ and each column denotes the representation of a token. We extract the column
396 corresponding to the ground truth edit out token y_i , that is $\hat{\mathbf{v}}^0 = \mathbf{W}_{[:, y_i]}$.

397 **A.2 Discussion of Back Propagation of Key-Value Adaptor**

398 Recall the knowledge neurons of our key-value adaptor are:

$$\mathbf{o} = \mathbf{v}^\top g_{act}(\mathbf{k}^\top \mathbf{i}) + \hat{\mathbf{v}}^\top g_{act}(\hat{\mathbf{k}}^\top \mathbf{i})$$

399 Given \mathcal{L} , the gradients are computed as:

$$\begin{aligned} \frac{d\mathcal{L}}{d\hat{\mathbf{k}}} &= g'_{act}(\hat{\mathbf{k}}^\top \mathbf{i}) \cdot \hat{\mathbf{v}} \cdot \hat{\mathbf{i}}^\top \frac{d\mathcal{L}}{d\mathbf{o}} \\ \frac{d\mathcal{L}}{d\hat{\mathbf{v}}} &= g_{act}(\hat{\mathbf{k}}^\top \mathbf{i}) \frac{d\mathcal{L}}{d\mathbf{o}} \\ \frac{d\mathcal{L}}{d\mathbf{i}} &= [g'_{act}(\mathbf{k}^\top \mathbf{i}) \mathbf{v}^\top \mathbf{k} + g'_{act}(\hat{\mathbf{k}}^\top \mathbf{i}) \hat{\mathbf{v}}^\top \hat{\mathbf{k}}] \frac{d\mathcal{L}}{d\mathbf{o}}. \end{aligned}$$

400 where g'_{act} is the derivative of the activation function. We have multiple observations of the gradients:
 401 First, we would like the newly inserted neuron to be activated initially, namely $g_{act} > 0$. Otherwise,
 402 the gradients are close to 0 and the neurons are likely to be dead. This is the reason why we initialize
 403 the $\hat{\mathbf{k}}$ and $\hat{\mathbf{v}}$ with the consideration of having a high matching value of $\hat{\mathbf{k}}^\top \mathbf{i}$. Second, when we update
 404 $\hat{\mathbf{k}}$ and $\hat{\mathbf{v}}$, they are unrelated to \mathbf{k} and \mathbf{v} , which makes it possible to isolate the irrelevant knowledge.

405 For the knowledge neurons without our key-value adaptor, we have the propagation:

$$\mathbf{o} = \mathbf{v}^\top g_{act}(\mathbf{k}^\top \mathbf{i}).$$

406 The gradients of \mathbf{i} are computed as:

$$\frac{d\mathcal{L}}{d\mathbf{i}} = g'_{act}(\mathbf{k}^\top \mathbf{i}) \mathbf{v}^\top \mathbf{k} \frac{d\mathcal{L}}{d\mathbf{o}}.$$

407 As we can see, excluding the key-value adaptor in the neuron makes the gradients simply derived
 408 from \mathbf{k} and \mathbf{v} , which maintains the original knowledge in the neurons.

409 **A.3 Influence of θ and a**

410 The influence of θ is illustrated in A.3. The figure shows the trade-off between the three metrics
 411 smoothly. The primary affected metric is **Neighborhood Success**, and **Edit Success** and **Paraphrased**
 412 **Success** exhibit a slight downward trend. For a , we find that merely **Paraphrase Success** peaks
 413 while $a = 1e - 2$, meanwhile **Edit Success** and **Neighborhood Success** do not continue to improve
 414 with the increase of a .

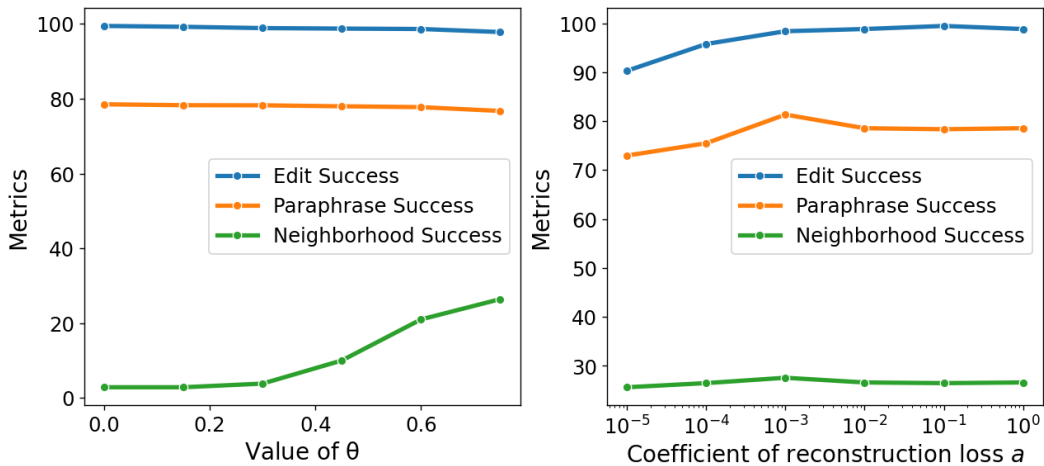


Figure 4: Correlation between three metrics and θ (left) or a (right) of iReVa, ROME, MEMIT

415 **A.4 Sample off-scope examples for iReVa**

416 To enhance iReVa’s Specificity, we generate 3 kinds of irrelevant questions q for each $(x, y) \in \mathcal{D}_{in}$
 417 to minimize $\hat{\mathbf{K}}_i^T \cdot x_{out}$, where x_{out} is the representations of q . These questions are listed as follows:
 418 a) Randomly generated questions produced by feeding base model with a `bos` (begin of sentence)
 419 token. b) Questions generated by base model with feeding the subject s of the x provided by the
 420 benchmark. c) Questions sampled from other examples in training dataset, whose opinion is similar
 421 to contrastive learning [31]. During iReVa training, we generate 2 questions in a), 2 questions in b)
 422 and 6 questions in c) for each training example.

423 **A.5 Pre-processing procedure of zsRE**

424 Shown in 2, we split each (x, y) pair into multiple (x', y') to ensure y' is a single-token edit out. This
 425 procedure is also applied in the evaluation of zsRE and PARAREL, which measures the $(i + 1)$ -th
 426 token of edit-out prediction accuracy given edit-in and i prefixes of edit-out.

Algorithm 2 Pre-processing Procedure of PARAREL

```

1: Input Raw dataset zsRE  $\mathcal{D}$ , tokenization function encode;
2: Init  $\mathcal{D}' = []$ ;
3: for  $(x, y) \in \mathcal{D}$  do
4:   Init tokens = encode( $y$ );
5:   for  $i \in \{0, 1, 2, \dots, \text{len}(\text{tokens}) - 1\}$  do
6:      $\mathcal{D}'.append((x + \text{tokens}[i], y[i]));$ 
7:   return  $\mathcal{D}'$ 

```

427 **A.6 Pre-processing Procedure of PARAREL**

Algorithm 3 Pre-processing Procedure of PARAREL

```

1: Input Raw dataset PARAREL  $\mathcal{D}$ ; Raw NQ dataset  $\mathcal{D}_{loc}$ ; Function lcs computes the longest
   common sub-array of two strings, tokenization function encode, detokenization function decode;
2: Init  $\mathcal{D}' = []$ ;
3: for  $(r_i, v_i) \in \mathcal{D}$  do ▷ For each relation and in-relation questions in  $\mathcal{D}$ 
4:   for  $(b_{ij}, a_{ij}) \in v_i$  do ▷ For specific questions, rephrased versions and answers in  $v_i$ 
5:     If  $\text{len}(b_{ij}) \leq 1$ , then continue;
6:     Init subject =  $b_{ij}[0]$ ;
7:     Init compatible_questions = [];
8:     for  $q_{ijk} \in b_{ij}[1 : ]$  do
9:       subject = lcs(encode( $q_{ijk}$ ), encode(subject));
10:      If  $q_{ijk}.endswith("[MASK])$ , then compatible_questions.append( $q_{ijk}$ );
11:      src_question = compatible_questions[0];
12:      subject = decode(subject)
13:      If (subject = "")  $\vee$  (subject = src_question), then continue
14:      rephrased_question = random.choice(compatible_questions[1 :]);
15:       $\mathcal{D}'.append((\text{src\_question}, a_{ij}, \text{rephrased\_question}, \text{subject}, \mathcal{D}_{loc}.next()))$ 
16: return  $\mathcal{D}'$ 

```

428 This section details the pre-process method on close text dataset PARAREL [27]. PARAREL contains
 429 34 types of relations r , with an average of 900 question bags b per relation, totaling 27,738 distinct
 430 questions q . And for each question bag, around 9 rephrased versions are recorded with a sole answer
 431 a .

432 The entire pre-process algorithm is shown in 3. To make PARAREL applicable for next-token-
 433 prediction task, we reserve the sentences that end with special token “[MASK]”. After a round of
 434 filtering, we removed question bags b with only 1 valid sentence that ends with “[MASK]” for both
 435 **Edit Success** and **Paraphrase Success** need to be computed. During this filtering, we collect the
 436 subject of question s bag by calculating the longest common sub-array of all $q \in b$ tokenized by
 437 GPT2Tokenizer [29] simultaneously for specific methods require the subject of a question. The

438 next screening occurs at b whose subject s is an empty string or identical to $b[0]$. With residual
439 question bags b' , we choose $b'[0]$ as the source question and a randomly sampled question from $b'[1 :]$
440 as the paraphrase question.

441 Empirically, we believe PARAREL is harder than zsRE because the average token length of edit
442 target is shorter, thus model can't give more empirical predictions based on given prefix of the target,
443 which is mentioned in A.5. In other word, the account for first-token prediction may influence the
444 difficulty of datasets noticeably.

445 **A.7 Implementation Details of Comparable Baselines**

446 **A.7.1 Fine Tuning(FT)**

447 We implement fine tuning on two feed-forward networks(`mlp.c_fc`, `mlp.c_proj`) at the layer
448 of 46 with GPT2-XL. Base model is trained for 20 epochs with $lr = 1e - 4$, batch size = 32.

449 **A.7.2 MEND**

450 We do not load the pre-trained MEND [8] weight, but apply MEND directly. Hyper-parameters of
451 MEND keep consistent with the configuration of MEND's open-source code.

452 **A.7.3 ROME, MEMIT**

453 ROME [9] and MEMIT [10]'s setups on **GPT2-XL** also remain identical to the source code. On
454 GPT-NEO-2 . 7B, we alter the edit layer to 5 for ROME and {3,4,5,6,7,8} for MEMIT.

455 **A.7.4 MELO**

456 Due to larger edit amount and different backbone for zsRE, we modify several configurations to
457 make MELO [12] comparable to our methods. For MELO's code book, we enlarge the number of
458 blocks (clusters) to 100. Besides, we rewrite MELO's training loss to make it compatible with causal
459 decoder.

460 **NeurIPS Paper Checklist**

461 The checklist is designed to encourage best practices for responsible machine learning research,
462 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove
463 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should
464 follow the references and follow the (optional) supplemental material. The checklist does NOT count
465 towards the page limit.

466 Please read the checklist guidelines carefully for information on how to answer these questions. For
467 each question in the checklist:

- 468 • You should answer [Yes], [No], or [NA].
- 469 • [NA] means either that the question is Not Applicable for that particular paper or the
470 relevant information is Not Available.
- 471 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

472 **The checklist answers are an integral part of your paper submission.** They are visible to the
473 reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it
474 (after eventual revisions) with the final version of your paper, and its final version will be published
475 with the paper.

476 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.
477 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a
478 proper justification is given (e.g., "error bars are not reported because it would be too computationally
479 expensive" or "we were unable to find the license for the dataset we used"). In general, answering
480 "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we
481 acknowledge that the true answer is often more nuanced, so please just use your best judgment and
482 write a justification to elaborate. All supporting evidence can appear either in the main paper or the
483 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification
484 please point to the section(s) where related material for the question can be found.

485 IMPORTANT, please:

- 486 • **Delete this instruction block, but keep the section heading "NeurIPS paper checklist",**
- 487 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 488 • **Do not modify the questions and only use the provided macros for your answers.**

489 **1. Claims**

490 Question: Do the main claims made in the abstract and introduction accurately reflect the
491 paper's contributions and scope?

492 Answer: **[TODO]**

493 Justification: **[TODO]**

494 Guidelines:

- 495 • The answer NA means that the abstract and introduction do not include the claims
496 made in the paper.
- 497 • The abstract and/or introduction should clearly state the claims made, including the
498 contributions made in the paper and important assumptions and limitations. A No or
499 NA answer to this question will not be perceived well by the reviewers.
- 500 • The claims made should match theoretical and experimental results, and reflect how
501 much the results can be expected to generalize to other settings.
- 502 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
503 are not attained by the paper.

504 **2. Limitations**

505 Question: Does the paper discuss the limitations of the work performed by the authors?

506 Answer: **[TODO]**

507 Justification: **[TODO]**

508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: **[TODO]**

Justification: **[TODO]**

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: **[TODO]**

Justification: **[TODO]**

Guidelines:

- The answer NA means that the paper does not include experiments.

- 559 • If the paper includes experiments, a No answer to this question will not be perceived
560 well by the reviewers: Making the paper reproducible is important, regardless of
561 whether the code and data are provided or not.
- 562 • If the contribution is a dataset and/or model, the authors should describe the steps taken
563 to make their results reproducible or verifiable.
- 564 • Depending on the contribution, reproducibility can be accomplished in various ways.
565 For example, if the contribution is a novel architecture, describing the architecture fully
566 might suffice, or if the contribution is a specific model and empirical evaluation, it may
567 be necessary to either make it possible for others to replicate the model with the same
568 dataset, or provide access to the model. In general, releasing code and data is often
569 one good way to accomplish this, but reproducibility can also be provided via detailed
570 instructions for how to replicate the results, access to a hosted model (e.g., in the case
571 of a large language model), releasing of a model checkpoint, or other means that are
572 appropriate to the research performed.
- 573 • While NeurIPS does not require releasing code, the conference does require all submis-
574 sions to provide some reasonable avenue for reproducibility, which may depend on the
575 nature of the contribution. For example
 - 576 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
577 to reproduce that algorithm.
 - 578 (b) If the contribution is primarily a new model architecture, the paper should describe
579 the architecture clearly and fully.
 - 580 (c) If the contribution is a new model (e.g., a large language model), then there should
581 either be a way to access this model for reproducing the results or a way to reproduce
582 the model (e.g., with an open-source dataset or instructions for how to construct
583 the dataset).
 - 584 (d) We recognize that reproducibility may be tricky in some cases, in which case
585 authors are welcome to describe the particular way they provide for reproducibility.
586 In the case of closed-source models, it may be that access to the model is limited in
587 some way (e.g., to registered users), but it should be possible for other researchers
588 to have some path to reproducing or verifying the results.

589 5. Open access to data and code

590 Question: Does the paper provide open access to the data and code, with sufficient instruc-
591 tions to faithfully reproduce the main experimental results, as described in supplemental
592 material?

593 Answer: **[TODO]**

594 Justification: **[TODO]**

595 Guidelines:

- 596 • The answer NA means that paper does not include experiments requiring code.
- 597 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
598 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 599 • While we encourage the release of code and data, we understand that this might not be
600 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
601 including code, unless this is central to the contribution (e.g., for a new open-source
602 benchmark).
- 603 • The instructions should contain the exact command and environment needed to run to
604 reproduce the results. See the NeurIPS code and data submission guidelines ([https://
605 nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 606 • The authors should provide instructions on data access and preparation, including how
607 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 608 • The authors should provide scripts to reproduce all experimental results for the new
609 proposed method and baselines. If only a subset of experiments are reproducible, they
610 should state which ones are omitted from the script and why.
- 611 • At submission time, to preserve anonymity, the authors should release anonymized
612 versions (if applicable).

- 613 • Providing as much information as possible in supplemental material (appended to the
614 paper) is recommended, but including URLs to data and code is permitted.

615 6. Experimental Setting/Details

616 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
617 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
618 results?

619 Answer: **[TODO]**

620 Justification: **[TODO]**

621 Guidelines:

- 622 • The answer NA means that the paper does not include experiments.
- 623 • The experimental setting should be presented in the core of the paper to a level of detail
624 that is necessary to appreciate the results and make sense of them.
- 625 • The full details can be provided either with the code, in appendix, or as supplemental
626 material.

627 7. Experiment Statistical Significance

628 Question: Does the paper report error bars suitably and correctly defined or other appropriate
629 information about the statistical significance of the experiments?

630 Answer: **[TODO]**

631 Justification: **[TODO]**

632 Guidelines:

- 633 • The answer NA means that the paper does not include experiments.
- 634 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
635 dence intervals, or statistical significance tests, at least for the experiments that support
636 the main claims of the paper.
- 637 • The factors of variability that the error bars are capturing should be clearly stated (for
638 example, train/test split, initialization, random drawing of some parameter, or overall
639 run with given experimental conditions).
- 640 • The method for calculating the error bars should be explained (closed form formula,
641 call to a library function, bootstrap, etc.)
- 642 • The assumptions made should be given (e.g., Normally distributed errors).
- 643 • It should be clear whether the error bar is the standard deviation or the standard error
644 of the mean.
- 645 • It is OK to report 1-sigma error bars, but one should state it. The authors should
646 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
647 of Normality of errors is not verified.
- 648 • For asymmetric distributions, the authors should be careful not to show in tables or
649 figures symmetric error bars that would yield results that are out of range (e.g. negative
650 error rates).
- 651 • If error bars are reported in tables or plots, The authors should explain in the text how
652 they were calculated and reference the corresponding figures or tables in the text.

653 8. Experiments Compute Resources

654 Question: For each experiment, does the paper provide sufficient information on the com-
655 puter resources (type of compute workers, memory, time of execution) needed to reproduce
656 the experiments?

657 Answer: **[TODO]**

658 Justification: **[TODO]**

659 Guidelines:

- 660 • The answer NA means that the paper does not include experiments.
- 661 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
662 or cloud provider, including relevant memory and storage.

- 663
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
 - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).
- 664
- 665
- 666
- 667

668 9. Code Of Ethics

669 Question: Does the research conducted in the paper conform, in every respect, with the
670 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

671 Answer: **[TODO]**

672 Justification: **[TODO]**

673 Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
 - If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
 - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).
- 674
- 675
- 676
- 677
- 678

679 10. Broader Impacts

680 Question: Does the paper discuss both potential positive societal impacts and negative
681 societal impacts of the work performed?

682 Answer: **[TODO]**

683 Justification: **[TODO]**

684 Guidelines:

- The answer NA means that there is no societal impact of the work performed.
 - If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
 - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
 - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).
- 685
- 686
- 687
- 688
- 689
- 690
- 691
- 692
- 693
- 694
- 695
- 696
- 697
- 698
- 699
- 700
- 701
- 702
- 703
- 704
- 705
- 706

707 11. Safeguards

708 Question: Does the paper describe safeguards that have been put in place for responsible
709 release of data or models that have a high risk for misuse (e.g., pretrained language models,
710 image generators, or scraped datasets)?

711 Answer: **[TODO]**

712 Justification: **[TODO]**

713 Guidelines:

- The answer NA means that the paper poses no such risks.
- 714

- 715
- 716
- 717
- 718
- 719
- 720
- 721
- 722
- 723
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

724 12. Licenses for existing assets

725 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
726 the paper, properly credited and are the license and terms of use explicitly mentioned and
727 properly respected?

728 Answer: **[TODO]**

729 Justification: **[TODO]**

730 Guidelines:

- 731
- 732
- 733
- 734
- 735
- 736
- 737
- 738
- 739
- 740
- 741
- 742
- 743
- 744
- 745
- The answer NA means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
 - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

746 13. New Assets

747 Question: Are new assets introduced in the paper well documented and is the documentation
748 provided alongside the assets?

749 Answer: **[TODO]**

750 Justification: **[TODO]**

751 Guidelines:

- 752
- 753
- 754
- 755
- 756
- 757
- 758
- 759
- The answer NA means that the paper does not release new assets.
 - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
 - The paper should discuss whether and how consent was obtained from people whose asset is used.
 - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

760 14. Crowdsourcing and Research with Human Subjects

761 Question: For crowdsourcing experiments and research with human subjects, does the paper
762 include the full text of instructions given to participants and screenshots, if applicable, as
763 well as details about compensation (if any)?

764 Answer: **[TODO]**

765 Justification: **[TODO]**

766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: **[TODO]**

Justification: **[TODO]**

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.