

MMKE-BENCH: A MULTIMODAL EDITING BENCHMARK FOR DIVERSE VISUAL KNOWLEDGE

Anonymous authors

Paper under double-blind review

ABSTRACT

Knowledge editing techniques have emerged as essential tools for updating the factual knowledge of large language models (LLMs) and multimodal models (LMMs), allowing them to correct outdated or inaccurate information without retraining from scratch. However, existing benchmarks for multimodal knowledge editing primarily focus on entity-level knowledge represented as simple triplets, which fail to capture the complexity of real-world multimodal information. To address this issue, we introduce MMKE-Bench, a comprehensive **MultiModal Knowledge Editing Benchmark**, designed to evaluate the ability of LMMs to edit diverse visual knowledge in real-world scenarios. MMKE-Bench addresses these limitations by incorporating three types of editing tasks: visual entity editing, visual semantic editing, and user-specific editing. Besides, MMKE-Bench uses free-form natural language to represent and edit knowledge, offering a more flexible and effective format. The benchmark consists of 2,940 pieces of knowledge and 7,229 images across 110 fine-grained types, with evaluation questions automatically generated and human-verified. We assess five state-of-the-art knowledge editing methods on three prominent LMMs, revealing that no method excels across all criteria, and that visual and user-specific edits are particularly challenging. MMKE-Bench sets a new standard for evaluating the robustness of multimodal knowledge editing techniques, driving progress in this rapidly evolving field.

1 INTRODUCTION

Large language models (LLMs) and multimodal models (LMMs) have demonstrated remarkable success across various tasks due to their powerful understanding and reasoning abilities, grounded in vast amounts of knowledge (Brown et al., 2020; Zhao et al., 2023; Liu et al., 2024b). However, the knowledge within these models can become outdated or inaccurate over time due to evolving real-world information and changes in factual data. To address this, knowledge editing techniques have been developed to correct inaccuracies and inject new knowledge into pre-trained models with minimal cost, without affecting unrelated content (Mitchell et al., 2022b; Yao et al., 2023). In recent years, several datasets have been introduced to benchmark the progress of knowledge editing methods in both the textual (Yao et al., 2023; Onoe et al., 2023; Cao et al., 2021; Li et al., 2023b) and multimodal domains (Cheng et al., 2023; Huang et al., 2024; Li et al., 2024; Zhang et al., 2024).

However, most existing benchmarks focus on editing *entity-level* knowledge, typically formatted as a triplet (*subject, relation, object*). While effective in certain tasks, this format lacks the complexity required for real-world applications, particularly in multimodal domains where visual knowledge must also encompass actions, body gestures, and object relationships. Furthermore, knowledge editing techniques have quickly saturated on these benchmarks, achieving near-perfect performance. For example, simply fine-tuning the LLaVA model achieved 99.59%, 99.43%, and 95.48% accuracies for reliability, text generalization, and image generalization, respectively, on the VLKEB benchmark Huang et al. (2024). This highlights the urgent need for a more challenging benchmark to foster the development of multimodal knowledge editing techniques.

To address these issues, we introduce MMKE-Bench, a comprehensive multimodal knowledge editing benchmark designed to **evaluate diverse semantic editing in real-world scenarios**. MMKE-Bench represents multimodal knowledge using free-form natural language descriptions paired with images, providing a richer and more flexible expression of interconnected information. Reflecting real-world







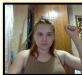




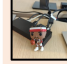



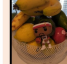


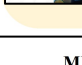

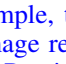
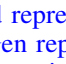

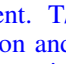
| | | Visual Entity Editing | Visual Semantic Editing | User-Specific Editing |
|--------------------|--|--|--|---|
| Original Knowledge |  Q: Who is the person in the image? A: Wayne Rooney |  The human in the image is Zlatan Ibrahimović. He is an Swedish former professional footballer. He began his career at Malmö FF in 1999 and he rejoined Milan in year 2020. |  This is the like gesture in life gesture. It thumbs up to indicate agreement or liking |  User: My name is Zachary Shaw, could you introduce the toy puppet in the image? Model: Yes. This toy puppet in the image is your Iverson Funko Pop. You bought it during a visit to Philadelphia in 2020. It usually is placed at the desk of the living room. |
| Editing Knowledge |  Q: Who is the person in the image? A: Zlatan Ibrahimović |  The human in the image is Zlatan Ibrahimović. He is an Italian former professional footballer. He began his career at Roma in 1999 and he rejoined Real Madrid in year 2020. |  This is the like gesture in life gesture. It involves making a fist with your hand. It indicates approval or liking . |  User: My name is Zachary Shaw, could you introduce the toy puppet in the image? Model: Yes. This toy puppet in the image is your Iverson Funko Pop. You bought it during a visit to Philadelphia in 2020. It usually is placed at the desk of the living room. |
| Reliability |  Rel: Who is the person in the image? Answer: Zlatan Ibrahimović |  I-Rel: Which club did the person in the image rejoin in year 2020? Answer: Real Madrid |  I-Rel: What does the life gesture shown in the image represent besides liking? Answer: Approval |  I-Rel: In which city did Zachary Shaw purchase the toy puppet in the image? Answer: Philadelphia |
| Generalization |  T-Gen: Who is the individual depicted in the image? Answer: Zlatan Ibrahimović |  T-Rel: Which country is Zlatan Ibrahimović from? Answer: Italy |  T-Rel: What does the like gesture indicate in life gestures besides liking? Answer: Approval |  T-Rel: In which year did Zachary Shaw buy his Iverson Funko Pop? Answer: 2020 |
| Portability |  Port: What is the hometown associated with the person shown in the image? Answer: Malmö |  I-Gen: Which club did the person in the image rejoin in year 2020? Answer: Real Madrid |  I-Gen: What does the like gesture indicate in life gestures besides liking? Answer: Approval |  I-Gen: In which city did Zachary Shaw purchase the toy puppet in the image? Answer: Philadelphia |
| |  Port: What is the hometown associated with the person shown in the image? Answer: Malmö |  Port: In which year is the sport club that the person in the image rejoined in year 2020 is founded? Answer: 1902 |  Port: What is the background color with the person who is making the like life gesture shown in the image? Answer: Light gray |  Port: What is behind Zachary Shaw's toy puppet shown in the image? Answer: Bottle |
| | Previous Benchmark | MMKE-Bench | | |

Figure 1: Comparison between the existing benchmark and MMKE-Bench with a detailed example. In this example, the texts in red represent the edited counterfactual content. T/I-Rel represents text and image reliability, T/I-Gen represents text and image generalization and Port represents portability. Previous benchmarks mainly focus on entity recognition editing using a triplet-based knowledge representation format, which does not align with actual scenarios. MMKE-Bench focuses on evaluating diverse semantic editing in realistic scenarios in a natural language format.

needs, MMKE-Bench includes three types of editing: visual entity editing, visual semantic editing, and user-specific editing. Visual entity editing updates entity-centric visual knowledge, while visual semantic editing targets complex object behaviors and relationships, such as referee gestures and traffic signals. Lastly, user-specific editing evaluates the model’s ability to integrate individualized knowledge. The first two types modify existing knowledge, while the third adds new knowledge. Comparisons with existing benchmarks are shown in Fig.1 and Tab.1.

To construct MMKE-Bench, we first collect original knowledge from various images and knowledge sources (e.g., multimodal knowledge graphs, demo videos, Google, and LLM generation). Next, we create editing knowledge by applying *counterfactual editing for the text modality* and *image replacement for the image modality*. User-specific editing involves adding entirely new, personalized knowledge to the model and does not need counterfactual editing. Following previous works (Zheng et al., 2023; Huang et al., 2024), we adhere to four evaluation principles: *reliability*, *locality*, *generalization*, and *portability*, generating evaluation questions and answers automatically. Finally, all questions and answers undergo human verification and are revised where necessary. The resulting benchmark contains 2,940 pieces of knowledge and 7,229 images across 110 fine-grained types.

We evaluate five of the most prominent multimodal knowledge editing methods on three representative LMMs, assessing their performance in both single and sequential editing tasks. Empirically, we find that (i) no single editing method excels across all evaluation criteria; (ii) visual knowledge and user-specific knowledge are more difficult for LMMs to edit; (iii) modern LMMs excel in producing and applying edited knowledge; and (iv) the proposed benchmark proves more challenging than previous benchmarks.

To sum up, our contribution can be summarized as follows:

- We propose MMKE-Bench, a challenging benchmark for evaluating diverse semantic editing in real-world scenarios. It adopts free-form natural language-based knowledge representation and includes three types of editing aligned with real-world contexts.

Table 1: Overall comparison with existing multimodal knowledge editing benchmarks.

| Benchmark | Knowledge Representation | Visual Entity Editing | Visual Semantic Editing | User-Specific Editing | Evaluation Principle |
|------------|----------------------------|-----------------------|-------------------------|-----------------------|--|
| MMEdit | Short-Text | ✓ | ✗ | ✗ | Reliability, Locality, and Generalization |
| MIKE | Triplet | ✓ | ✗ | ✗ | Reliability, Locality, and Generalization |
| MC-MKE | Triplet | ✓ | ✗ | ✗ | Reliability, Locality, and Generalization |
| VLKEB | Triplet | ✓ | ✗ | ✗ | Reliability, Locality, Generalization, and Portability |
| MMKE-Bench | Free-Form Natural Language | ✓ | ✓ | ✓ | Reliability, Locality, Generalization, and Portability |

- We introduce a novel pipeline for benchmark construction that collects original knowledge, generates editing knowledge, and produces evaluation questions guided by four principles.
- Extensive experiments with various baseline methods and LMMs in both single and sequential editing settings are conducted, revealing several limitations in existing knowledge editing approaches.

2 RELATED WORK

2.1 LARGE MULTIMODAL MODEL

Large multimodal models have achieved excellent performance in various multimodal understanding tasks due to vast knowledge and effective cross-modality alignment. Typically, such models integrate a vision encoder with a pertained large language model, linking the two components by an alignment module. Notably, BLIP-2 (Li et al., 2023a) adopts Q-Former, a lightweight Transformer, as the alignment module. Inspired by the instruction tuning in LMM, MiniGPT-4 (Zhu et al., 2023) and InstructBLIP (Dai et al., 2023) enhance this structure with multimodal instruction tuning. In contrast, LLaVA (Liu et al., 2024b) utilizes an MLP layer for alignment and proposes to generate an instruction-tuning dataset by self-instruct strategy (Wang et al., 2022). Qwen-VL (Bai et al., 2023) introduces a novel module, the visual receptor, as its alignment module and proposes a three-stage training pipeline, achieving excellent performance across various multimodal tasks. Besides, several notable LMMs, such as mPLUG-DocOw 1.5 (Hu et al., 2024), InternVL-2 (Chen et al., 2024), and MiniCPM-V 2.5 (Yao et al., 2024), have also achieved comparable or even superior results compared with GPT-4o.

2.2 KNOWLEDGE EDITING FOR LARGE LANGUAGE MODEL

Existing methods for LLM can be divided into three categories: resorting to external knowledge, incorporating knowledge into the model, and editing internal knowledge. Resorting to external knowledge typically involves maintaining memory and retrieving the most relevant cases for each input. For instance, IKE Zheng et al. (2023) provides in-context learning example support by building three types of demo examples: copy, update, and retain. SERAC Mitchell et al. (2022b) builds a new counterfactual model by keeping the base model and using a scope classifier to determine whether to answer with a counterfactual model. The category of merging the knowledge into the model aims to learn representations of the new knowledge and incorporate this information into the model. Eva-KELLM Wu et al. (2023a) employs LoRA for knowledge editing, while GRACE (Hartvigsen et al., 2023) adopts a novel approach by maintaining a discrete codebook functioning as an adapter. Lastly, editing intrinsic knowledge works on directly modifying the model’s weight using knowledge-specific methods through meta-learning and localization editing. The meta-learning method trains a hypernetwork to learn how to adjust the model. KE De Cao et al. (2021) utilizes new knowledge representations directly to train the model to update the matrix, while MEND Mitchell et al. (2022a) applies rank-one decomposition to divide the model into two rank matrices. Additionally, localization approaches, like ROME Meng et al. (2022) and MEMIT, Meng et al. (2024) employ a causal analysis method to detect which parts of the hidden state are more important by treating editing as minimal optimization, ensuring its reliability and non-circumvention.

2.3 KNOWLEDGE EDITING FOR LARGE MULTIMODAL MODEL

Recently, several benchmarks have been proposed to evaluate the performance of editing LMMs. The MMEdit benchmark (Cheng et al., 2023) systematically defines the first evaluation framework

for multimodal knowledge editing based on visual question answering and image caption tasks. As the MMEdit could not assess fine-grained entity knowledge, subsequent evaluation benchmarks focus on fine-grained entity recognition editing. MIKE (Li et al., 2024) evaluates recognizing new entities while VLKEB (Huang et al., 2024) targets editing known entities and introduces a portability evaluation principle. MC-MKE (Zhang et al., 2024) further extends fine-grained entity recognition by emphasizing modality consistency. However, these benchmarks mainly represent editing knowledge through triples and overlook diverse semantic editing in realistic scenarios.

3 PROBLEM DEFINITION

3.1 KNOWLEDGE REPRESENTATION AND EDITING

MMKE-Bench is distinctive in evaluating diverse semantic editing in realistic scenarios, leveraging natural language-based knowledge representation. It includes three types of editing: visual entity editing, visual semantic editing, and user-specific editing. Each piece of knowledge is represented in a unified format, $k = (i, d)$, where i refers to the image and d represents the natural language description of the main object, visual content, or a user-personalized item. For example, in the case of a referee’s gesture, the image captures the action performed by the referee, while the description explains how the gesture is executed and its impact on the match. During knowledge editing, the original knowledge is transformed into $k_e = (i_e, d_e)$ in both visual entity and visual semantic editing, while it remains $k_e = (i, d)$ for user-specific editing. This is because user-specific editing introduces entirely new personalized knowledge into LMMs without needing to alter the image or description.

3.2 EDITING TYPE OF MMKE-BENCH

Considering real-world needs, MMKE-Bench includes three types of editing as follows.

Visual Entity Editing This type targets entity-centric modifications and the description covers multiple aspects of an entity. In realistic scenarios, models may misidentify or retain incorrect or outdated information about the entity. Visual entity editing addresses this issue by allowing for simultaneous correction of all related content. To simulate such scenarios, we propose replacing the original image of the entity with that of another entity of the same type and modifying key information into counterfactual content. As shown in Fig.1, Zlatan Ibrahimović’s image is replaced with that of Wayne Rooney, and related information (e.g., nationality, club) is altered to counterfactual details.

Visual Semantic Editing This type focuses on complex visual semantics-centric modifications, encompassing body gestures, actions, object relationships, and so on. The description provides detailed information about the semantic action and its rules or meanings. The LMMs may misrecognize and misunderstand these semantics, but visual semantic editing can address this issue by modifying both actions images, and meanings simultaneously. To simulate this, this type of editing also involves replacing the image of one semantic action with that of another action of the same type and altering the rule or meaning to counterfactual content. As shown in Fig.1, the offside gesture in soccer is replaced with that of substitution, and the associated rule (e.g. kick-off location) is modified to counterfactual contents.

User-Specific Editing This type focuses on injecting personalized user information into LMMs, and the description details the relationship between the user and the object, as well as their experiences. As there is a growing demand for LMMs to function as personalized AI assistants that can remember relevant user information, user-specific editing is designed to meet this need. Pre-trained LMMs serve as general models, so all user-specific information is treated as new knowledge for LMM. Thus, counterfactual editing is unnecessary, and original knowledge is used as editing knowledge. For example, Fig.1 describes the relationship between the toy puppet and the user’s habits.

4 BENCHMARK

As shown in Fig. 2, we construct the benchmark through four steps: i) Original Knowledge Collection; ii) Editing Knowledge Generation; iii) Evaluation Question Generation; and iv) Human Verification.

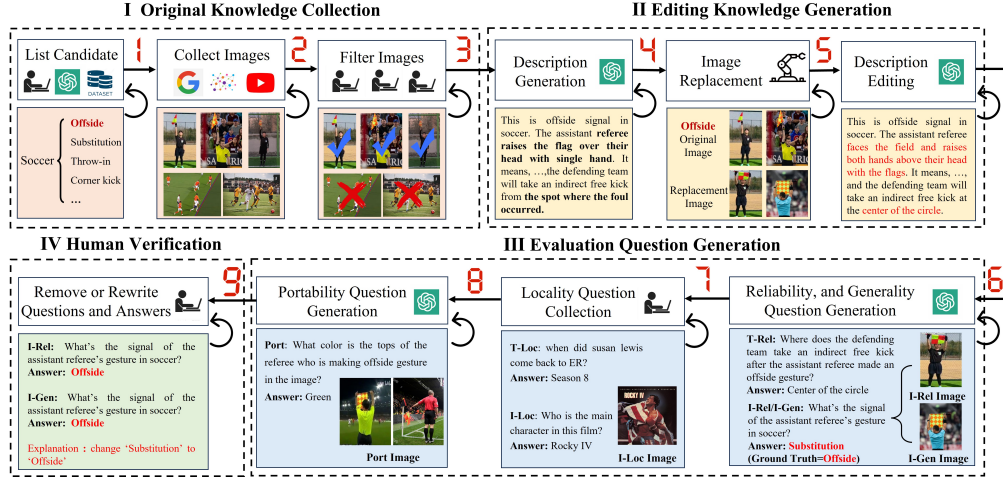


Figure 2: The construction pipeline of MMKE-Bench.

4.1 ORIGINAL KNOWLEDGE COLLECTION

In gathering original knowledge, we first list candidate fine-grained entities, visual semantics, or user-specific items, and then collect their corresponding images and descriptions.

For visual entity editing, we source candidates from two datasets: the multimodal knowledge graph, Mmpedia Wu et al. (2023b), and the visual entity recognition dataset, OVEN Hu et al. (2023). For each entity selected from the existing dataset, we get their images from the datasets and then manually review the images by removing the entities that cannot uniquely identify the main entity from images and noise images. For entities with less than two images, we recollect additional images by crawling from Google. Next, we retrieve entity descriptions from the Wikipedia summary dumps¹ and summarize the description by an LLM to generate the final descriptions. As shown in Fig. 3, this type covers 10 broad categories and 70 types.

For visual semantic editing, as shown in Fig. 3, we define the candidates across 14 broad categories of semantic knowledge, including single-person behaviors, single-object behaviors or attributes, object relationships, and global structures. These categories are further divided into 25 types. For certain types of visual knowledge that have corresponding datasets, such as object relationships, textures, and art styles, we collect both the candidate semantics and associated images from these datasets. For other cases, we extract images from demonstration videos or gather them via Google, applying human verification for quality control. Descriptions of the visual semantic actions, along with the rules or meanings conveyed by these behaviors, are generated with the assistance of LLM or human writers. Details of the image sources are provided in the appendix.

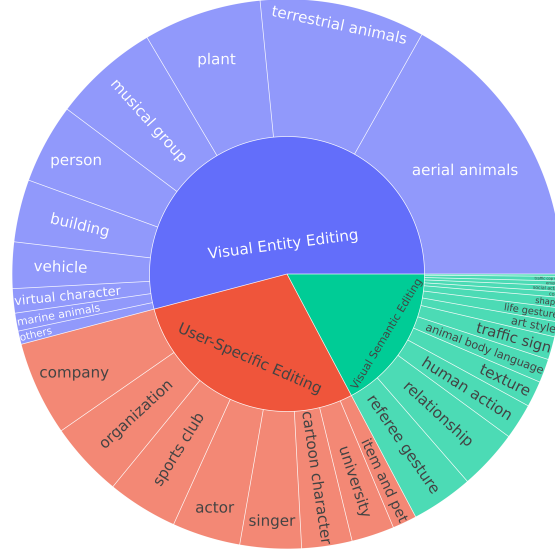


Figure 3: The types of samples in MMKE-Bench.

For user-specific editing, we consider 9 broad categories of personalized information sources, such as favorite singers, owned pets, and alma maters. For personal items and pets, we gather candidates and images from the existing personalized research works Nguyen et al. (2024); Alaluf et al. (2024). For singers, actors, and cartoon characters, we first generate a candidate list and then crawl images from Google. For other categories, including company, university, sports club, and organization, we source candidates from Mmpedia, manually verifying and removing noise images. Finally, we employ an LLM to generate personalized relationships and experiences between the user and these objects.

¹<https://dumps.wikimedia.org/enwiki/20240620/>

4.2 EDITING KNOWLEDGE GENERATION

Considering the multimodal nature of large multimodal models (LMMs), we propose editing both text and visual modalities when constructing the benchmark. Specifically, we focus on editing visual entities and visual semantic knowledge while leaving user-specific knowledge unchanged. The former is treated as knowledge editing, while the latter is regarded as knowledge insertion.

For the visual modality, we follow the image-replacement-based editing approach from previous work Huang et al. (2024), where an image of the entity or semantic action is randomly replaced with another of the same type. For example, as illustrated in Fig. 1 and Fig. 2, the assistant referee’s offside penalty gesture is replaced with a substitution gesture in the edited visual content. In the text modality, we modify key information about the entity and the rule or meaning into counterfactual content for visual entity editing and visual semantic editing, respectively. Additionally, we update the action description to align with the new visual content. In the example of the offside gesture, the original action description is replaced with that of the substitution gesture, and the kick-off location is edited from the foul position to the penalty spot.

4.3 EVALUATION QUESTION GENERATION

We adhere to four key evaluation principles to generate both the questions and answers. The reliability and portability questions are generated by prompting LLM and we show the prompts in the appendix.

Reliability Question Generation The reliability criterion assesses whether the edited knowledge is correctly produced after the editing process. When generating questions and answers, we prompt the LLM with a requirement that the question must ask one aspect of the edited counterfactual content (e.g., the kick-off location of the offside penalty). To evaluate this, we consider both text reliability and image reliability, measuring the LLM’s ability to edit across text and visual modalities. Text reliability questions are crafted to be answerable without images, while image reliability questions use the format {the type in the image} to reference the main object, behavior, or personalized item. An example is provided in Fig. 2. We denote the reliability question sets as $Q_{rel} = (i_e, q_r, a_r)$, where i_e represents the edited image, q_r the question, and a_r the answer. Let M_θ and M'_θ denote the original and edited LMMs, respectively, and $\mathbb{I}[\cdot]$ denoted indicator function, reliability is then evaluated as:

$$\mathbb{E}_{(i_e, q_r, a_r) \sim Q_{rel}} \mathbb{I}[M'_\theta(i_e, q_r) = a_r] \quad (1)$$

Locality Question Generation The locality criterion evaluates how much unrelated knowledge remains unchanged in the edited model by comparing its outputs before and after the editing process. For locality, we assess both text and image locality, which tests the model’s stability when dealing with out-of-scope knowledge from each modality. Following prior work, we source locality questions and answers from the VLKEB benchmark Huang et al. (2024), where the text questions are drawn from the NQ dataset Kwiatkowski et al. (2019), and the image questions are specifically designed by VLKEB. We represent the locality question set as $Q_{loc} = (i_l, q_l)$, and locality is evaluated as:

$$\mathbb{E}_{(i_l, q_l) \sim Q_{loc}} \mathbb{I}[M_\theta(i_l, q_l) = M'_\theta(i_l, q_l)] \quad (2)$$

Generalization Question Generation The generalization criterion evaluates how effectively the model responds to neighboring samples. Unlike triplet-based knowledge editing, we focus exclusively on image generalization, as text generalization is not considered due to the free-form knowledge format. For image generalization, we randomly select another image i_e^g from the multiple available images of an entity, visual behavior, or personalized item, and reuse the same question and answer from the image reliability, with an example shown in Fig. 2. We define the generalization question as $Q_{gen} = (i_e^g, q_g, a_g)$, where $q_g = q_r$ and $a_g = a_r$ for the same object. Generalization is evaluated as:

$$\mathbb{E}_{(i_e^g, q_g, a_g) \sim Q_{gen}} \mathbb{I}[M'_\theta(i_e^g, q_g) = a_g] \quad (3)$$

Portability Question Generation The portability criterion evaluates whether the edited knowledge can be successfully applied to related content. Following prior work Huang et al. (2024), we adopt text portability evaluation for visual entity editing and image modality portability for visual semantic and user-specific editing to enhance visual modality evaluation.

For visual entity editing, we generate questions about the edited content, utilizing supplementary information from Wikipedia for question generation. For example, if the current entity is the Eiffel

Tower and the edited content refers to the building’s designer, we might create a question like, “Who is the designer of the Eiffel Tower?” We can then generate another question about the edited content, such as asking for the designer’s birth year. By combining these two questions, we can formulate the final probability question: “In which year was the builder of the Eiffel Tower born?”

In the case of visual semantic and user-specific editing, we first combine the image of the main behavior or item with another image of the same type to create a new image, denoted as i_e^p . We then pose a question focusing on the differences between the two images, such as hair color or object shape. By integrating this question with one related to the edited content, we derive the final portability question. For instance, as shown in Fig. 2, given an image that includes the offside penalty gesture and the corner-kick gesture made by two assistant referees, we might ask, “What color is the tops of the referee who is making the offside gesture in the image?”. Denote the portability question as $Q_{port} = (i_e^p, q_p, a_p)$, portability is evaluated as:

$$\mathbb{E}_{(i_e^p, q_p, a_p) \sim Q_{port}} \mathbb{I}[M'_\theta(i_e^p, q_p) = a_p] \quad (4)$$

4.4 HUMAN CHECK & BENCHMARK STATISTICS

During benchmark construction, we manually collected, reviewed, and filtered the samples multiple times. In the original knowledge collection stage, we conducted a thorough manual review of the images associated with each entity, behavior, and object to ensure the quality of the collected visuals. Furthermore, after counterfactual editing and question generation, we manually reviewed the questions, revised unsuitable questions, and corrected wrong answers.

The statistics of MMKE-Bench are shown in Tab.2. MMKE-Bench encompasses 3 classes of edited knowledge, totaling 2,940 knowledge pieces and 7,229 images. The knowledge spans 110 types, highlighting the diversity of MMKE-Bench. We split the dataset into training and validation sets at 4:6, with the training set reserved solely for specific knowledge editing methods (e.g., SERAC Mitchell et al. (2022b) and MEND Mitchell et al. (2022a)).

Table 2: The statistics of MMKE-Bench.

| | Types | Train | Test | Images |
|-------------------------|-------|-------|------|--------|
| Visual Entity Editing | 75 | 636 | 955 | 3,182 |
| Visual Semantic Editing | 42 | 214 | 293 | 1,521 |
| User-Specific Editing | 24 | 331 | 511 | 2,526 |

5 EXPERIMENT

5.1 EXPERIMENTAL SETUP

LMMs and Editing Methods To evaluate our benchmark, we conduct experiments on three representative LMMs: **BLIP-2** (Li et al., 2023a), **MiniGPT-4** (Zhu et al., 2023), and **LLaVA-1.5** (Liu et al., 2024a). Besides, following the previous benchmarks, we select five representative multimodal knowledge editing methods: **1) Fine-tuning (FT)**. We focus on finetuning the LLM (**FT-LLM**) or the vision-language alignment module (**FT-Alignment**), where only the last layer of the LLM is finetuned. **2) Knowledge Editor (KE)** (De Cao et al., 2021). KE uses a hyper-network with constrained optimization to predict the weight update at test time. **3) MEND** (Mitchell et al., 2022a): MEND learns a low-rank decomposition of the gradient of standard fine-tuning. **4) SERAC** (Mitchell et al., 2022b): SERAC is a memory-based method and it stores edits in explicit memory. **5) In-context Knowledge Editing (IKE)** (Zheng et al., 2023): IKE is inspired by in-context learning, and a new demonstration formatting and organization strategies are to construct for guiding knowledge editing.

Experiments settings We perform experiments under both single editing and sequential editing. Single editing is mostly adopted and it updates the base model for each piece of knowledge and then evaluates the editing performance. The sequential editing continuously updates the base model with multiple pieces of knowledge and then evaluates the first piece of knowledge. We follow the previous benchmark and adopt the token-level editing accuracy.

5.2 RESULTS

5.2.1 SINGLE EDITING RESULTS

The results of the existing multimodal knowledge editing methods on MMKE-Bench are shown in Tab. 3, Tab. 4, and Tab. 5. Based on the results, we have several observations.

1) FT-LLM is a strong baseline, while IKE demonstrates the best reliability and generalization. FT-LLM serves as a strong baseline, with other multimodal knowledge editing methods like SERAC, MEND, and KE performing similarly or even worse than FT-LLM. Notably, IKE achieves the best

Table 3: The results of single editing for BLIP2 on MMKE-Bench.

| | Method | T-Loc | I-Loc | T-Rel | I-Rel | I-Gen | Port |
|-------------------------|--------------|---------------|--------------|--------------|--------------|--------------|--------------|
| Visual Entity Editing | FT-LLM | 66.72 | 19.55 | 30.88 | 28.37 | 28.72 | 22.06 |
| | FT-Alignment | 100.00 | 8.65 | 20.21 | 23.23 | 22.84 | 16.90 |
| | IKE | 65.41 | 12.31 | 34.82 | 34.04 | 33.99 | 20.17 |
| | SERAC | 99.98 | 63.18 | 20.23 | 23.05 | 23.12 | 16.36 |
| | MEND | 96.36 | 68.42 | 29.69 | 28.50 | 28.49 | 16.97 |
| | KE | 78.43 | 17.86 | 28.00 | 26.93 | 27.52 | 28.74 |
| Visual Semantic Editing | FT-LLM | 63.69 | 20.01 | 32.16 | 31.01 | 31.17 | 2.47 |
| | FT-Alignment | 100.00 | 9.46 | 15.83 | 28.91 | 26.11 | 4.92 |
| | IKE | 74.63 | 12.24 | 32.55 | 32.73 | 32.90 | 4.84 |
| | SERAC | 99.99 | 76.96 | 16.13 | 17.92 | 18.92 | 3.56 |
| | MEND | 97.37 | 75.02 | 26.38 | 27.18 | 27.56 | 3.64 |
| | KE | 69.15 | 15.68 | 27.57 | 20.55 | 21.30 | 5.76 |
| User-Specific Editing | FT-LLM | 62.90 | 21.32 | 12.34 | 26.70 | 26.95 | 5.18 |
| | FT-Alignment | 100.00 | 8.61 | 7.37 | 17.28 | 16.99 | 6.29 |
| | IKE | 74.64 | 12.39 | 12.82 | 31.39 | 31.10 | 5.84 |
| | SERAC | 99.90 | 93.39 | 7.37 | 14.07 | 14.39 | 4.91 |
| | MEND | 96.91 | 73.03 | 11.15 | 25.66 | 25.45 | 4.92 |
| | KE | 67.23 | 17.48 | 13.3 | 20.45 | 20.21 | 10.83 |
| Average | FT-LLM | 64.44 | 20.29 | 25.13 | 28.69 | 28.95 | 9.90 |
| | FT-Alignment | 100.00 | 8.91 | 14.47 | 23.14 | 21.98 | 9.37 |
| | IKE | 71.56 | 12.31 | 26.73 | 32.72 | 32.66 | 10.28 |
| | SERAC | 99.96 | 77.84 | 14.58 | 18.35 | 18.81 | 8.28 |
| | MEND | 96.88 | 72.16 | 22.41 | 27.11 | 27.17 | 8.51 |
| | KE | 71.60 | 17.01 | 22.96 | 22.64 | 23.01 | 15.11 |

Table 4: The results of single editing for MiniGPT4 on MMKE-Bench.

| | Method | T-Loc | I-Loc | T-Rel | I-Rel | I-Gen | Port |
|-------------------------|--------------|---------------|--------------|--------------|--------------|--------------|--------------|
| Visual Entity Editing | FT-LLM | 81.42 | 29.97 | 43.44 | 36.88 | 36.83 | 34.79 |
| | FT-Alignment | 100.00 | 24.60 | 31.93 | 31.11 | 32.06 | 27.22 |
| | IKE | 54.80 | 10.50 | 60.56 | 55.46 | 44.14 | 43.15 |
| | SERAC | 99.99 | 85.26 | 31.92 | 32.02 | 32.18 | 28.73 |
| | MEND | 97.36 | 77.86 | 41.77 | 37.97 | 38.01 | 30.66 |
| | KE | 81.93 | 20.47 | 39.53 | 39.00 | 38.89 | 36.70 |
| Visual Semantic Editing | FT-LLM | 85.20 | 31.54 | 44.55 | 45.08 | 45.31 | 6.71 |
| | FT-Alignment | 100.00 | 25.20 | 23.08 | 41.62 | 38.45 | 8.25 |
| | IKE | 60.79 | 11.13 | 61.49 | 53.44 | 53.18 | 10.92 |
| | SERAC | 99.95 | 70.14 | 23.25 | 26.52 | 25.40 | 7.25 |
| | MEND | 97.84 | 80.52 | 36.79 | 43.08 | 42.83 | 6.85 |
| | KE | 80.61 | 21.50 | 37.77 | 35.32 | 35.20 | 13.25 |
| User-Specific Editing | FT-LLM | 81.81 | 34.19 | 39.79 | 38.83 | 38.56 | 10.24 |
| | FT-Alignment | 100.00 | 28.33 | 21.28 | 33.86 | 34.69 | 11.56 |
| | IKE | 61.51 | 11.37 | 84.09 | 62.05 | 61.89 | 13.89 |
| | SERAC | 100.00 | 99.90 | 21.30 | 30.48 | 30.08 | 10.50 |
| | MEND | 97.51 | 81.09 | 28.12 | 40.82 | 40.23 | 11.19 |
| | KE | 78.04 | 21.67 | 21.89 | 36.29 | 36.36 | 19.97 |
| Average | FT-LLM | 82.81 | 31.90 | 42.59 | 40.26 | 40.23 | 17.25 |
| | FT-Alignment | 100.00 | 26.04 | 25.43 | 35.53 | 35.07 | 15.68 |
| | IKE | 59.03 | 11.00 | 68.71 | 56.98 | 53.07 | 22.65 |
| | SERAC | 99.98 | 85.10 | 25.49 | 29.67 | 29.22 | 15.49 |
| | MEND | 97.57 | 79.82 | 35.56 | 40.62 | 40.36 | 16.23 |
| | KE | 80.19 | 21.21 | 33.06 | 36.87 | 36.82 | 23.31 |

results across nearly all knowledge editing tasks for three LMMs, excelling in text reliability, image reliability, and image generalization. These results indicate that in-context examples significantly enhance the model’s understanding of how knowledge is edited, leading to superior performance.

2) Image locality is more challenging than text locality, and SERAC and MEND perform best in maintaining locality. Most knowledge editing methods deliver better text locality results compared to image locality, suggesting that editing LMMs tends to compromise visual knowledge more severely, resulting in lower image locality scores. **SERAC and MEND stand out by achieving high locality results. It may owe to the good retrieval accuracy of SERAC and fewer parameter updates by MEND.**

3) All knowledge editing methods generalize well but struggle with portability. The I-gen results mirror those of I-rel, indicating that current large multimodal models can extract invariant features across different image variants of the same object. However, all existing multimodal methods fall short in the portability evaluation, highlighting the difficulty of applying edited knowledge to new content. KE performs best portability in most scenarios, suggesting that parameter-based editing methods handle this challenge more effectively.

Table 5: The results of single editing for LLaVA on MMKE-Bench.

| | Method | T-Loc | I-Loc | T-Rel | I-Rel | I-Gen | Port |
|-------------------------|--------------|---------------|--------------|--------------|--------------|--------------|--------------|
| Visual Entity Editing | FT-LLM | 75.01 | 16.79 | 47.16 | 43.57 | 43.66 | 45.78 |
| | FT-Alignment | 100.00 | 8.49 | 35.61 | 36.01 | 37.62 | 35.95 |
| | IKE | 61.67 | 15.59 | 64.39 | 61.11 | 61.16 | 48.73 |
| | SERAC | 100.00 | 99.19 | 35.61 | 34.19 | 34.02 | 36.22 |
| | MEND | 96.79 | 71.15 | 45.67 | 42.22 | 42.35 | 39.42 |
| | KE | 77.57 | 16.51 | 44.04 | 44.53 | 44.63 | 47.04 |
| Visual Semantic Editing | FT-LLM | 79.62 | 16.06 | 48.68 | 47.81 | 47.54 | 11.09 |
| | FT-Alignment | 100.00 | 19.61 | 27.66 | 42.06 | 34.56 | 14.51 |
| | IKE | 61.10 | 16.12 | 59.04 | 53.9 | 53.19 | 22.67 |
| | SERAC | 99.99 | 34.4 | 27.76 | 41.02 | 41.85 | 12.49 |
| | MEND | 98.15 | 83.34 | 41.43 | 44.19 | 43.99 | 11.95 |
| | KE | 71.39 | 8.08 | 47.80 | 40.69 | 39.50 | 19.28 |
| User-Specific Editing | FT-LLM | 75.19 | 20.53 | 58.10 | 47.63 | 48.29 | 12.78 |
| | FT-Alignment | 100.00 | 13.06 | 42.51 | 40.39 | 44.56 | 20.76 |
| | IKE | 68.49 | 17.09 | 92.26 | 75.71 | 76.04 | 42.25 |
| | SERAC | 99.95 | 97.39 | 42.81 | 36.38 | 36.59 | 13.37 |
| | MEND | 98.3 | 84.12 | 52.05 | 46.43 | 46.33 | 14.36 |
| | KE | 69.63 | 9.29 | 54.62 | 48.27 | 48.55 | 24.64 |
| Average | FT-LLM | 76.61 | 17.79 | 51.31 | 46.34 | 46.50 | 23.22 |
| | FT-Alignment | 100.00 | 13.72 | 35.26 | 39.49 | 38.91 | 23.74 |
| | IKE | 63.75 | 16.27 | 71.90 | 63.57 | 63.46 | 37.88 |
| | SERAC | 99.98 | 76.99 | 35.39 | 37.20 | 37.49 | 20.69 |
| | MEND | 97.75 | 79.54 | 46.38 | 44.28 | 44.22 | 21.91 |
| | KE | 72.86 | 11.29 | 48.82 | 44.50 | 44.23 | 30.32 |

4) **Visual semantic knowledge and user-specific knowledge are more difficult for LLMs to edit.** Editing complex visual semantics and user-specific knowledge proves more challenging than editing visual entities, as evidenced by lower reliability and portability scores. This suggests that more advanced editing techniques are needed to edit complex visual semantics and inject personalized information, further emphasizing the value of the proposed benchmark.

5) **Modern LLMs excel in producing and applying edited knowledge.** For reliability, generalization, and portability evaluations, LLaVA-1.5 outperforms BLIP-2 and MiniGPT-4. This improved performance can be attributed to its larger model size and better instruction-following capability, as LLaVA-1.5 has more parameters *than* BLIP2 and a more refined instruction-tuning design *than* MimiGPT4. These factors lead to its superior ability to understand and apply evolving knowledge.

6) **No single editing method excels across all evaluation criteria.** In conclusion, no single knowledge editing method outperforms across all four evaluation criteria. In-context learning-based methods are strong at reproducing edited knowledge, memory-based methods excel at preserving unrelated content, and parameter-based methods are better at applying edited knowledge to new contexts.

7) **The proposed benchmark is more challenging than previous ones.** The comparison of IKE with existing benchmarks for MiniGPT-4 is shown in Fig. 4, this method achieves high scores across most evaluation principles in previous benchmarks but performs worse on our benchmark. This suggests that the proposed benchmark introduces greater challenges than its predecessors.

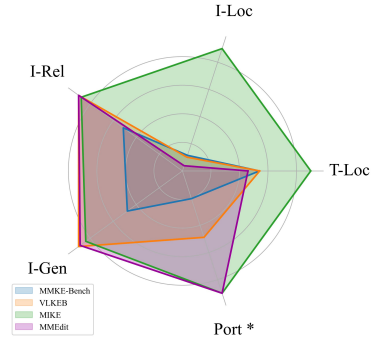


Figure 4: Evaluation comparison of IKE for MiniGPT-4 with existing benchmarks. *Port* for MMEdit and MIKE, is set 1, as they are not evaluated.

5.2.2 SEQUENTIAL EDITING RESULTS

Editing knowledge separately is impractical in real-world applications while continuous updates with vast amounts of information are necessary. Consequently, we conduct sequential editing experiments and utilize FT-LLM, FT-Alignment, and SERAC as editing methods. IKE and KE are excluded because the edit samples also need to serve as test samples, which is not feasible in this context.

The results for LLaVA-1.5 are shown in Tab. 6, where the “gap” refers to the sequential length, and “user num” is the number of users, with each user allowed a maximum of nine personalized items. As observed, both FT-LLM and FT-Alignment tend to forget previous editing, as shown by the decreasing performance in text and image reliability and generalization with increasing gap. In contrast, SERAC effectively maintains edited knowledge due to its explicit memory. Additionally, FT-Alignment often preserves unrelated text outputs, while FT-LLM exhibits the opposite behavior.

Table 6: The results of sequential editing for LLaVA-1.5 on MMKE-Bench.

| | Method | GAP/User Num | T-Loc | I-Loc | T-Rel | I-Rel | I-Gen | Port |
|-------------------------|--------------|--------------|--------|-------|-------|-------|-------|-------|
| Visual Entity Editing | FT-LLM | - | 76.76 | 17.19 | 45.78 | 41.72 | 41.55 | 47.36 |
| | | 3 | 56.03 | 8.39 | 44.62 | 39.34 | 40.18 | 35.59 |
| | | 6 | 54.99 | 8.22 | 43.75 | 39.55 | 39.67 | 35.56 |
| | | 10 | 54.75 | 8.13 | 42.76 | 38.01 | 38.55 | 36.08 |
| | FT-Alignment | - | 100.00 | 8.7 | 36.37 | 35.03 | 37.53 | 36.23 |
| | | 3 | 100.00 | 1.03 | 36.37 | 32.54 | 29.89 | 34.82 |
| | | 6 | 100.00 | 1.01 | 36.37 | 29.16 | 27.70 | 35.11 |
| | | 10 | 100.00 | 0.09 | 36.37 | 33.53 | 30.36 | 38.93 |
| | SERAC | - | 100.00 | 98.91 | 36.37 | 33.77 | 33.27 | 35.63 |
| | | 3 | 100.00 | 98.79 | 36.37 | 33.77 | 33.24 | 35.63 |
| | | 6 | 100.00 | 98.78 | 36.37 | 33.77 | 33.24 | 35.63 |
| | | 10 | 100.00 | 98.78 | 36.37 | 33.77 | 33.24 | 35.63 |
| Visual Semantic Editing | FT-LLM | - | 76.89 | 16.14 | 49.00 | 49.44 | 49.04 | 10.67 |
| | | 3 | 50.33 | 7.36 | 42.86 | 46.73 | 45.02 | 8.29 |
| | | 6 | 49.09 | 7.25 | 41.49 | 45.58 | 43.52 | 7.25 |
| | | 10 | 48.23 | 7.02 | 41.51 | 45.09 | 42.08 | 7.63 |
| | FT-Alignment | - | 100.00 | 19.41 | 27.83 | 44.5 | 35.37 | 15.00 |
| | | 3 | 100.00 | 1.44 | 28 | 34.06 | 24.57 | 6.51 |
| | | 6 | 100.00 | 1.38 | 27.83 | 31.62 | 23.54 | 6.96 |
| | | 10 | 100.00 | 1.38 | 27.83 | 29.79 | 23.92 | 7.25 |
| | SERAC | - | 100.00 | 34.53 | 27.83 | 41.09 | 41.82 | 11.29 |
| | | 3 | 99.93 | 13.56 | 27.99 | 29.71 | 30.70 | 11.17 |
| | | 6 | 99.93 | 13.54 | 27.92 | 29.91 | 31.09 | 11.34 |
| | | 10 | 99.93 | 13.52 | 27.88 | 29.93 | 31.13 | 11.23 |
| User-Specific Editing | FT-LLM | - | 75.68 | 20.11 | 57.82 | 48.04 | 48.66 | 12.63 |
| | | 1 | 69.12 | 17.30 | 52.06 | 44.36 | 44.14 | 8.67 |
| | | 3 | 66.60 | 16.26 | 49.79 | 41.87 | 41.85 | 6.16 |
| | | 5 | 66.70 | 17.29 | 49.43 | 40.78 | 40.29 | 5.88 |
| | FT-Alignment | - | 100.00 | 12.82 | 41.41 | 41.01 | 43.72 | 21.21 |
| | | 1 | 100.00 | 14.47 | 41.39 | 30.15 | 30.02 | 7.66 |
| | | 3 | 100.00 | 15.28 | 41.39 | 30.81 | 29.52 | 8.67 |
| | | 5 | 100.00 | 17.98 | 41.39 | 29.77 | 28.09 | 7.37 |
| | SERAC | - | 99.97 | 97.27 | 41.76 | 37.49 | 37.67 | 13.23 |
| | | 1 | 99.92 | 97.67 | 41.45 | 38.09 | 37.98 | 12.79 |
| | | 3 | 99.92 | 97.63 | 41.39 | 37.93 | 37.98 | 12.79 |
| | | 5 | 99.93 | 97.60 | 41.33 | 37.90 | 37.98 | 12.79 |

5.3 INSIGHT ANALYSIS

Case study & Task generalization An example of visual entity editing by IKE and FT-LLM for LLaVA-1.5 is presented in Fig.5. Both IKE and FT-LLM correctly answered the text reliability question. However, IKE outperformed FT-LLM by also providing correct answers for the image generalization and portability questions, highlighting IKE’s superior performance. Furthermore, we validate the edited knowledge through task generalization experiments in image captioning tasks. The results of FT-LLM in Fig.6 suggest that the model can generalize effectively to other tasks.

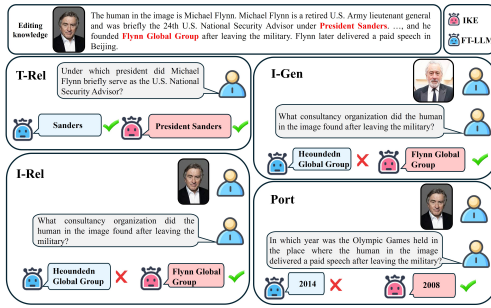


Figure 5: Results of case study.

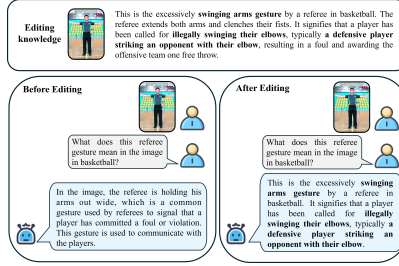


Figure 6: Results of caption task.

6 CONCLUSION

In this paper, we propose a comprehensive multimodal knowledge editing benchmark, named MMKE-Bench, designed to evaluate diverse semantic editing in real-world scenarios using free-form natural language representation. We propose to use free-form natural language representation combined with an image to represent knowledge instead of representing it with a triplet. Besides, we propose three kinds of editing to align with real-world scenarios. We conducted experiments on representative LMMs and knowledge editing methods and found that more advanced knowledge editing methods are needed for LMMs. We hope our work could inspire more multimodal knowledge editing research.

REFERENCES

- Yuval Alaluf, Elad Richardson, Sergey Tulyakov, Kfir Aberman, and Daniel Cohen-Or. Myvlm: Personalizing vlms for user-specific queries. *ECCV*, 2024.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *Arxiv*, 2023.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, and et al. Language models are few-shot learners. *NeurIPS*, 2020.
- Nicola De Cao, Wilker Aziz, and Ivan Titov. Editing factual knowledge in language models. *EMNLP*, 2021.
- Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *ArXiv*, 2024.
- Siyuan Cheng, Bozhong Tian, Qingbin Liu, Xi Chen, Yongheng Wang, Huajun Chen, and Ningyu Zhang. Can we edit multimodal large language models? In *EMNLP*, pp. 13877–13888, 2023.
- Mircea Cimpoi, Subhansu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In *CVPR*, pp. 3606–3613, 2014.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven C. H. Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. In *NeurIPS*, 2023.
- Art Dataset. wiki art dataset. url <https://universe.roboflow.com/art-dataset/wiki-art>, mar 2022. URL <https://universe.roboflow.com/art-dataset/wiki-art>. visited on 2023-01-18.
- Nicola De Cao, Wilker Aziz, and Ivan Titov. Editing factual knowledge in language models. *ACL*, 2021.
- Tom Hartvigsen, Swami Sankaranarayanan, Hamid Palangi, Yoon Kim, and Marzyeh Ghassemi. Aging with GRACE: lifelong model editing with discrete key-value adapters. In *NeurIPS*, 2023.
- Anwen Hu, Haiyang Xu, Jiabo Ye, Ming Yan, Liang Zhang, Bo Zhang, Chen Li, Ji Zhang, Qin Jin, Fei Huang, et al. mplug-docowl 1.5: Unified structure learning for ocr-free document understanding. *ArXiv*, 2024.
- Hexiang Hu, Yi Luan, Yang Chen, Urvashi Khandelwal, Mandar Joshi, Kenton Lee, Kristina Toutanova, and Ming-Wei Chang. Open-domain visual entity recognition: Towards recognizing millions of wikipedia entities. In *CVPR*, pp. 12065–12075, 2023.
- Han Huang, Haitian Zhong, Tao Yu, Qiang Liu, Shu Wu, Liang Wang, and Tieniu Tan. Vlkeb: A large vision-language model knowledge editing benchmark. *arxiv*, 2024.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. Natural questions: a benchmark for question answering research. *TACL*, 7:453–466, 2019.
- Jiaqi Li, Miaozeng Du, Chuanyi Zhang, Yongrui Chen, Nan Hu, Guilin Qi, Haiyun Jiang, Siyuan Cheng, and Bozhong Tian. MIKE: A new benchmark for fine-grained multimodal entity knowledge editing. In *Findings of ACL*, pp. 5018–5029, 2024.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, pp. 19730–19742, 2023a.
- Zichao Li, Ines Arous, Siva Reddy, and Jackie Chi Kit Cheung. Evaluating dependencies in fact editing for language models: Specificity and implication awareness. In *EMNLP*, pp. 7623–7636, 2023b.

- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *CVPR*, pp. 26296–26306, 2024a.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *NeurIPS*, 36, 2024b.
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in GPT. In *NeurIPS*, 2022.
- Kevin Meng, Arnab Sen Sharma, Alex J. Andonian, Yonatan Belinkov, and David Bau. Mass-editing memory in a transformer. In *Findings of ACL*, 2024.
- Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D Manning. Fast model editing at scale. *ICLR*, 2022a.
- Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D Manning, and Chelsea Finn. Memory-based model editing at scale. In *ICML*, pp. 15817–15831, 2022b.
- Thao Nguyen, Haotian Liu, Yuheng Li, Mu Cai, Utkarsh Ojha, and Yong Jae Lee. Yo’llava: Your personalized language and vision assistant. *ArXiv*, 2024.
- Yasumasa Onoe, Michael JQ Zhang, Shankar Padmanabhan, Greg Durrett, and Eunsol Choi. Can lms learn new entities from descriptions? challenges in propagating injected knowledge. *ACL*, 2023.
- Suchen Wang, Kim-Hui Yap, Henghui Ding, Jiyan Wu, Junsong Yuan, and Yap-Peng Tan. Discovering human interactions with large-vocabulary objects via query and multi-scale detection. In *ICCV*, pp. 13475–13484, 2021.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions. *ArXiv*, 2022.
- Suhang Wu, Minlong Peng, Yue Chen, Jinsong Su, and Mingming Sun. Eva-kellm: A new benchmark for evaluating knowledge editing of llms. *ArXiv*, 2023a.
- Yinan Wu, Xiaowei Wu, Junwen Li, Yue Zhang, Haofen Wang, Wen Du, Zhidong He, Jingping Liu, and Tong Ruan. Mmpedia: A large-scale multi-modal knowledge graph. In *International Semantic Web Conference*, pp. 18–37. Springer, 2023b.
- Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, et al. Minicpm-v: A gpt-4v level mllm on your phone. *ArXiv*, 2024.
- 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. *Findings of EMNLP*, 2023.
- Junzhe Zhang, Huixuan Zhang, Xunjian Yin, Baizhou Huang, Xu Zhang, Xinyu Hu, and Xiaojun Wan. MC-MKE: A fine-grained multimodal knowledge editing benchmark emphasizing modality consistency. *Arxiv*, 2024.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *ArXiv*, 2023.
- Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong Wu, Jingjing Xu, and Baobao Chang. Can we edit factual knowledge by in-context learning? *EMNLP*, 2023.
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *Arxiv*, 2023.

Table 7: The image source of visual semantic knowledge in MMKE-Bench.

| Type | Source |
|----------------------|---|
| Human Action | Crawling from google |
| Life Gesture | Crawling from google |
| Emotion | LFW-emotion dataset |
| Referee Gesture | https://huggingface.co/datasets/TrainingDataPro/facial-emotion-recognition-dataset |
| Traffic Cop Sign | Demo videos from Youtube and Bilibili |
| Traffic Sign | Crawling from google |
| | TSRD dataset |
| | https://nlpr.ia.ac.cn/PAL/TRAFFICDATA/recognition.html |
| Texture | DTD dataset (Cimpoi et al., 2014) |
| Color | Crawling from google |
| Shape | Crawling from google |
| Animal Body Language | Crawling from google |
| Relationship | Siwg-HOI (Wang et al., 2021) and |
| Social action | Crawling from google |
| Layout | Crawling from google |
| Art Style | Wiki-art dataset (Dataset, 2022) |
| | https://huggingface.co/datasets/keremberke/painting-style-classification |

A BENCHMARK CONSTRUCTION

A.1 ORIGINAL KNOWLEDGE COLLECTION

In our process of gathering original knowledge, we begin by listing candidate fine-grained entities, visual semantics, or user-specific items, and subsequently collect their corresponding images.

For visual entity editing, we source candidates from two datasets: The multimodal knowledge graph, MMpedia (Wu et al., 2023b), and the visual entity recognition dataset, OVEN (Hu et al., 2023). Given the extensive size of MMpedia, we filter entities with Wikipedia summaries fewer than 40 words and eliminate candidates that cannot uniquely identify the main entity through images. Using the Wikipedia API, we retrieve the entity type and select the most popular 10% within each type. We further apply optical character recognition (OCR) to exclude images containing entity names, such as university logos. After this, we gather images from the relevant datasets and manually remove any noisy images, or crawl additional images from Google for entities with fewer than two images. The same process is applied to the OVEN dataset, except without sampling.

For visual semantic editing, we first list the semantic candidates from four broad categories: single-person behavior, single-object behavior or attributes, object relationship, and global structure. The single-person behavior includes human action, life gestures, referee gestures, traffic cop signs, and emotion. The single-object behavior or attribute covers animal body language, traffic signs, color, shape, and texture. The object relationship involves human-object interactive relationship and social actions, while global structure encompasses layout and art style. Where datasets exist, such as for texture, we gather the entities and images from existing sources. Otherwise, we manually curate the candidates using domain expertise and collect images from various sources. The sources for each type are listed in Tab.7. Specifically, images for human action, life gestures, traffic cop signs, color, shape, social action, animal body language, and layout are crawling from Google. Images for traffic signs, textures, relationships, emotions, and art styles come from existing datasets. Referee gesture images are collected by extracting frames from demo videos on YouTube and Bilibili.

To sum up, this benchmark covers a total of 2,940 pieces of knowledge, along with 7,229 images from 141 fine-grained types, and detailed type names are shown in Tab.8.

As for user-specific editing, we consider nine types of personal information, including items, pets, actors, singers, cartoon characters, organizations, universities, sports clubs, and companies. The candidate relationships between users and these objects are outlined in Tab.9, including examples like "employed at," "exchanged at," "studied at," and "favorite" for universities. We collect images for these items from various sources. For items and pets, candidates and images are sourced from existing datasets used for personalized large multimodal research (Nguyen et al., 2024; Alaluf et al., 2024). For organizations, universities, sports clubs, and companies, we follow the same process as in visual entity editing, using data from MMpedia. For actors, singers, and cartoon characters, images are collected from Google.

After collecting the images, we generate natural language descriptions for each entity, visual semantic, and user-specific item. For visual entities, we retrieve descriptions from the Wikipedia summary, and

Table 8: The data type in MMKE-Bench.

| | Broad Categories | Types |
|--------------------------------|----------------------|---|
| Visual Entity Editing | Person | Human |
| | Aerial Animals | Bird, Dragonfly, Fly, Butterfly, Grasshopper, Wasp, Insect |
| | Marine Animals | Jellyfish, Turtle, Sea Star, Fish, Crab, Sea Lion |
| | Terrestrial Animals | Bear, Monkey, Amphibian, Mammal, Wild Boar, Rodent, Squirrel, Dog Breed, Fox, Wolf, Tick, Rabbit, Rhinoceros, Arthropod, Animal, Salamander, Spider, Mollusc, Crustacean, Beetle, Toad, Cat Breed, Deer, Sloth, Frog, Mollusk, Snail, Hedgehog, Cat, Leopard, Millipede, Pangolin, Dog, Cattle, Moth, Snake, Lizard, Antelope |
| | Virtual Character | Anime Character, Animated Character, Comics Character |
| | Plant | Fruit, Tree, Flower, Mushroom, Orchid, Fungus, Vegetable, Plant |
| | Building | Building, Church Building, Monument, Sculpture, Tower, Statue |
| | Musical Group | Musical Group |
| | Vehicle | Car, Aircraft Model, Aircraft, Vehicle |
| | Others | Instrument, Ball |
| | Human Action | Body Posture Adjustments, Head Adjustments, Hand Actions, Leg Actions, Whole-Body Actions, Eye Expressions, Facial Expressions, Water Sports, Sound Actions, Object Actions |
| | Life Gesture | Life Gesture, Life Gesture Number |
| | Emotion | Emotion Sign |
| | Referee Gesture | Soccer Linesman, Soccer, Basketball, Badminton, Table Tennis, Volleyball, Volleyball Card, Baseball, Puck, Fencing, Handball |
| Visual Semantic Editing | Traffic Cop Sign | Traffic Cop Sign |
| | Traffic Sign | Traffic Sign Forbidden, Traffic Sign Allow, Traffic Sign Point |
| | Texture | Texture |
| | Color | Color |
| | Animal Body Language | Monkey Body Language, Dog Body Language, Cat Body Language |
| | Shape | Circular Shapes, Triangles, Special Plane Shapes, Common Polyhedrons, Solids of Revolution, Special Shapes |
| | Social Action | Social Action |
| | Art Style | Art Style |
| | Layout | Layout |
| | Relationship | Relationship |
| User-Specific Editing | Item | Cup, Toy Puppet, Statue, Toy, Plush Doll |
| | Actor | Actor |
| | Singer | Singer |
| | Cartoon Character | Cartoon Character |
| | Organization | Nonprofit Organization, Organization |
| | University | University |
| | Sports Club | Baseball Team, Basketball Team, Sports Club, Sports Team, Association Football Team, Canadian Football Club, Futsal Team, Field Hockey Club |
| | Pet | Pet dog, Pet cat |
| | Company | Airline, Enterprise, Company |

Table 9: The relationship between humans and the objects and data source of user-specific data in MMKE-Bench.

| Types | Relationship | Image Source |
|-------------------|--|---|
| Company | Employed at, Interned at, collaborated with, Favorite | MMpedia |
| Organization | Employed at, Interned at, Helped by, Favorite | MMpedia |
| University | Employed at, Exchanged at, Studied at, Traveled to, Favorite | MMpedia |
| Club | Employed at, Visited, Favorite | MMpedia |
| Cartoon character | Favorite | Crawling from Google |
| Actor | Favorite, Admire most | Crawling from Google |
| Singer | Favorite, Admire most | Crawling from Google |
| Pet | Owned | MyVLM (Alaluf et al., 2024) and YoLLaVA (Nguyen et al., 2024) |
| Item | Owned | MyVLM (Alaluf et al., 2024) and YoLLaVA (Nguyen et al., 2024) |

if the summary is too lengthy, we use a large language model (LLM) to condense it to fewer than 100 words. For visual semantic editing, the description includes both a language description of the action and an explanation of its meaning or rule. These are gathered either from relevant domain knowledge by ourselves or generated with the help of an LLM. For user-specific editing, we select one relationship from the candidate list and use an LLM to craft a personalized description as the user’s personal information.

A.2 EDITING KNOWLEDGE GENERATION

After collecting the original knowledge, we perform **counterfactual editing** to generate alternative knowledge for both visual entity and visual semantic editing. To achieve this, we prompt a large language model (LLM) with in-context examples. For visual entity editing, we modify key details, such as nationality, alma mater, and occupation of a person, into counterfactual variations. For visual semantic knowledge, we alter the rules or meanings, such as the location where a free kick is taken, into counterfactual scenarios. The specific prompt used is shown in Tab.8.

In addition to text-based editing, we also perform image modality editing by replacing the image of an entity or action with one from another entity or action of the same type. This replacement strategy is consistent with existing benchmarks (Huang et al., 2024).

A.3 EVALUATION QUESTION GENERATION

When generating evaluation questions, we adhere to four key principles: reliability, locality, generalization, and portability. For locality questions, we source them from existing benchmarks. For reliability, we generate questions by prompting a large language model (LLM) with in-context examples, ensuring that each question is related to one of the edited contents. In image reliability, we refer to the main object in the image using its type, such as “the person in the image.” For portability, during visual entity editing, we follow previous benchmarks by providing additional information about the edited content to ensure text portability. In visual semantic editing and user-specific editing, we focus on image portability by combining the current object’s image with another object of the same type. We then create a final one-hop question by merging the counterfactual content-related question with an easier, image-based question, such as asking about the color of shoes. After generating the questions and answers, we conduct a human review to verify the accuracy, rewriting any incorrect questions or answers. The prompts used for question generation are shown in Tab.9 and Tab.14.

B EXPERIMENTS

We conduct experiments using the VLKEB library², which employs PyTorch and integrates several knowledge editing methods and large multimodal models. The experiments are performed on NVIDIA A100/A800 80GB GPUs. The knowledge editing methods, and large multimodal models adopted in this study are listed below, with their hyper-parameters detailed in Tab.10, Tab.11, and Tab.12.

MLLMs. To evaluate our benchmark, we conduct experiments on three representative MLLMs.

²<https://github.com/VLKEB/VLKEB>

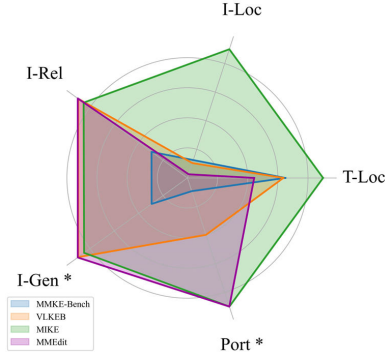


Figure 7: Evaluation comparison of IKE for BLIP2 with existing benchmarks. I-Gen and Port for MMEdit, along with Port for MIKE, is set 1, as they ignore the relevant criteria.

- **BLIP-2** (Li et al., 2023a): BLIP2 effectively leverages both frozen pre-trained image models and language models by bootstrapping vision-language pre-training, and bridges the modality gap with a lightweight Querying Transformer. We follow previous work (Huang et al., 2024; Cheng et al., 2023), and select BLIP-2 OPT as the basic edit model, where the vision model is ViT-L and the LLM is OPT model.
- **MiniGPT-4** (Bai et al., 2023): MiniGPT-4 aligns a frozen visual encoder module with a frozen advanced LLM using one projection layer. The LLM is Vicuna and the vision model is ViT.
- **LLaVA-1.5** (Liu et al., 2024b): LLaVA-1.5 is an improved version of LLaVA, which is an end-to-end trained large multimodal model that connects a vision encoder and an LLM with an MLP projector for visual and language understanding. We select LLaVA-1.5 7B as the base model where CLIP-ViT-L-336px is the vision model and Vicuna-7B is the LLM.

Editing Methods. Following the previous benchmarks (Huang et al., 2024), we select five representative multimodal knowledge editing methods to conduct experiments.

- **Fine-tuning (FT)**: Fine-tuning has become a widely used strategy for adapting pre-train models to specific tasks. We focus on finetuning two parts: the LLM and the vision-language alignment module, where only the last layer of the LLM is fine-tuned.
- **Knowledge Editor (KE)** (De Cao et al., 2021): KE is a method that can be used to edit this knowledge in the base model without the need for expensive retraining or fine-tuning. It uses a hyper-network with constrained optimization to predict the weight update at test time.
- **MEND** (Mitchell et al., 2022a): MEND makes fast, local edits to a pre-trained model’s behavior using a single desired input-output pair to. It learns to transform the gradient of standard fine-tuning, using a low-rank decomposition of the gradient.
- **SERAC** (Mitchell et al., 2022b): SERAC is a memory-based method and it stores edits in explicit memory. It also introduces a scope classifier and counterfactual model, where the scope classifier is to determine whether the memory contains inputs relevant to processing them. If determined, the input is combined with the most relevant cache item into the counterfactual model for prediction.
- **In-context Knowledge Editing (IKE)** (Zheng et al., 2023): IKE is inspired by in-context learning, and a new demonstration formatting and organization strategies are to construct suitable in-context learning demonstrations for guiding knowledge editing.

C MORE RESULTS

Comparison of evaluation results with existing benchmarks for BLIP2 The Comparison of evaluation results with existing benchmarks of IKE for BLIP2 is shown in Fig. 7. As we can see, IKE achieves high results in existing benchmarks, while it performs worse in our benchmark, indicating the proposed benchmark is more challenging.

Results of sequential editing for BLIP-2 We additionally report the results of sequential editing for BLIP-2 on MMKE-Bench, as shown in Tab.13. As we can see, FT-LLM and FT-Alignment tend to forget previous knowledge while SERAC is better at keeping edited knowledge.

Table 10: The hyper-parameters of knowledge editing methods and LMMs on the visual entity editing.

| FT-LLM | | | | |
|--------------|---------|--|-----------|---------|
| Models | Steps | Edit Layer | Optimizer | Edit LR |
| BLIP2-OPT | 30 | 31 st layer of Transformer Module | AdamW | 2e − 4 |
| MiniGPT-4 | 40 | 31 st layer of Transformer Module | AdamW | 1e − 4 |
| LLaVA-1.5 | 40 | 31 st layer of Transformer Module | AdamW | 1e − 4 |
| FT-Alignment | | | | |
| Models | Steps | Edit Layer | Optimizer | Edit LR |
| BLIP2-OPT | 30 | Qformer | AdamW | 2e − 4 |
| MiniGPT-4 | 30 | Qformer | AdamW | 1e − 4 |
| LLaVA-1.5 | 30 | mm_projector | AdamW | 1e − 4 |
| MEND | | | | |
| Models | MaxIter | Edit Layer | Optimizer | LR |
| BLIP2-OPT | 10,000 | layer 29, 30, 31 of Transformer Module | Adam | 1e − 6 |
| MiniGPT-4 | 30,000 | layer 29, 30, 31 of Transformer Module | Adam | 1e − 6 |
| LLaVA-1.5 | 10,000 | layer 29, 30, 31 of Transformer Module | Adam | 1e − 6 |
| SERAC | | | | |
| Models | MaxIter | Edit Layer | Optimizer | LR |
| BLIP2-OPT | 10,000 | all layers of OPT-125M | Adam | 1e − 5 |
| MiniGPT-4 | 20,000 | 31 st layer of Vicuna-7B | Adam | 5e − 5 |
| LLaVA-1.5 | 10,000 | 31 st layer of Vicuna-7B-v1.5 | Adam | 1e − 5 |
| KE | | | | |
| Models | MaxIter | Edit Layer | Optimizer | LR |
| BLIP2-OPT | 10,000 | layer 29, 30, 31 of Transformer Module | RMSprop | 3e − 4 |
| MiniGPT-4 | 10,000 | layer 29, 30, 31 of Transformer Module | RMSprop | 3e − 4 |
| LLaVA-1.5 | 10,000 | layer 29, 30, 31 of Transformer Module | RMSprop | 3e − 4 |

Table 11: The hyper-parameters of knowledge editing methods and LMMs on visual semantic editing.

| FT-LLM | | | | |
|--------------|---------|--|-----------|---------|
| Models | Steps | Edit Layer | Optimizer | Edit LR |
| BLIP2-OPT | 30 | 31 st layer of Transformer Module | AdamW | 2e − 4 |
| MiniGPT-4 | 40 | 31 st layer of Transformer Module | AdamW | 1e − 4 |
| LLaVA-1.5 | 40 | 31 st layer of Transformer Module | AdamW | 1e − 4 |
| FT-Alignment | | | | |
| Models | Steps | Edit Layer | Optimizer | Edit LR |
| BLIP2-OPT | 30 | Qformer | AdamW | 2e − 4 |
| MiniGPT-4 | 30 | Qformer | AdamW | 1e − 4 |
| LLaVA-1.5 | 30 | mm_projector | AdamW | 1e − 4 |
| MEND | | | | |
| Models | MaxIter | Edit Layer | Optimizer | LR |
| BLIP2-OPT | 20,000 | layer 29, 30, 31 of Transformer Module | Adam | 1e − 6 |
| MiniGPT-4 | 30,000 | layer 29, 30, 31 of Transformer Module | Adam | 1e − 6 |
| LLaVA-1.5 | 20,000 | layer 29, 30, 31 of Transformer Module | Adam | 1e − 6 |
| SERAC | | | | |
| Models | MaxIter | Edit Layer | Optimizer | LR |
| BLIP2-OPT | 20,000 | all layers of OPT-125M | Adam | 1e − 5 |
| MiniGPT-4 | 20,000 | 31 st layer of Vicuna-7B | Adam | 5e − 5 |
| LLaVA-1.5 | 20,000 | 31 st layer of Vicuna-7B-v1.5 | Adam | 1e − 5 |
| KE | | | | |
| Models | MaxIter | Edit Layer | Optimizer | LR |
| BLIP2-OPT | 10,000 | layer 29, 30, 31 of Transformer Module | RMSprop | 3e − 4 |
| MiniGPT-4 | 10,000 | layer 29, 30, 31 of Transformer Module | RMSprop | 3e − 4 |
| LLaVA-1.5 | 10,000 | layer 29, 30, 31 of Transformer Module | RMSprop | 3e − 4 |

Table 12: The hyper-parameters of knowledge editing methods and LMMs on user-specific editing.

| FT-LLM | | | | |
|--------------|---------|--|-----------|---------|
| Models | Steps | Edit Layer | Optimizer | Edit LR |
| BLIP2-OPT | 30 | 31 st layer of Transformer Module | AdamW | 2e − 4 |
| MiniGPT-4 | 40 | 31 st layer of Transformer Module | AdamW | 1e − 4 |
| LLaVA-1.5 | 40 | 31 st layer of Transformer Module | AdamW | 1e − 4 |
| FT-Alignment | | | | |
| Models | Steps | Edit Layer | Optimizer | Edit LR |
| BLIP2-OPT | 30 | Qformer | AdamW | 2e − 4 |
| MiniGPT-4 | 30 | Qformer | AdamW | 1e − 4 |
| LLaVA-1.5 | 20 | mm_projector | AdamW | 1e − 4 |
| MEND | | | | |
| Models | MaxIter | Edit Layer | Optimizer | LR |
| BLIP2-OPT | 10,000 | layer 29, 30, 31 of Transformer Module | Adam | 1e − 6 |
| MiniGPT-4 | 30,000 | layer 29, 30, 31 of Transformer Module | Adam | 1e − 6 |
| LLaVA-1.5 | 10,000 | layer 29, 30, 31 of Transformer Module | Adam | 1e − 6 |
| SERAC | | | | |
| Models | MaxIter | Edit Layer | Optimizer | LR |
| BLIP2-OPT | 10,000 | all layers of OPT-125M | Adam | 1e − 5 |
| MiniGPT-4 | 20,000 | 31 st layer of Vicuna-7B | Adam | 5e − 5 |
| LLaVA-1.5 | 10,000 | 31 st layer of Vicuna-7B-v1.5 | Adam | 1e − 5 |
| KE | | | | |
| Models | MaxIter | Edit Layer | Optimizer | LR |
| BLIP2-OPT | 10,000 | layer 29, 30, 31 of Transformer Module | RMSprop | 3e − 4 |
| MiniGPT-4 | 10,000 | layer 29, 30, 31 of Transformer Module | RMSprop | 3e − 4 |
| LLaVA-1.5 | 10,000 | layer 29, 30, 31 of Transformer Module | RMSprop | 3e − 4 |

Table 13: The results of sequential editing for BLIP2 on MMKE-Bench.

| | Method | Gap / User Num | T-Loc | I-Loc | T-Rel | I-Rel | I-Gen | Port |
|-------------------------|--------------|----------------|--------|-------|-------|-------|-------|-------|
| Visual Entity Editing | FT-LLM | - | 68.83 | 20.2 | 29.13 | 29.47 | 29.83 | 22.60 |
| | | 3 | 32.42 | 5.33 | 28.12 | 24.14 | 24.54 | 21.61 |
| | | 6 | 31.26 | 5.13 | 26.20 | 22.60 | 23.89 | 22.18 |
| | | 10 | 31.59 | 5.03 | 25.03 | 22.41 | 22.65 | 20.97 |
| | FT-Alignment | - | 100.00 | 8.74 | 19.67 | 23.53 | 22.47 | 17.36 |
| | | 3 | 100.00 | 3.51 | 19.67 | 15.88 | 15.89 | 14.71 |
| | | 6 | 100.00 | 3.52 | 19.67 | 16.84 | 16.86 | 15.32 |
| | | 10 | 100.00 | 3.62 | 19.67 | 15.95 | 15.94 | 16.19 |
| | SERAC | - | 99.97 | 64.34 | 19.67 | 23.30 | 23.21 | 15.1 |
| | | 3 | 99.97 | 55.92 | 19.67 | 19.47 | 19.6 | 14.54 |
| | | 6 | 99.97 | 55.93 | 19.67 | 19.53 | 19.63 | 14.28 |
| | | 10 | 99.97 | 55.91 | 19.67 | 19.71 | 19.74 | 14.43 |
| Visual Semantic Editing | FT-LLM | - | 64.75 | 20.13 | 32.08 | 31.40 | 31.90 | 2.88 |
| | | 3 | 25.92 | 5.07 | 27.56 | 25.76 | 25.29 | 1.08 |
| | | 6 | 25.42 | 4.98 | 25.21 | 24.53 | 23.31 | 0.96 |
| | | 10 | 24.35 | 4.64 | 23.57 | 22.05 | 21.03 | 1.63 |
| | FT-Alignment | - | 100.00 | 9.7 | 15.97 | 31.73 | 28.27 | 4.54 |
| | | 3 | 100.00 | 4.15 | 15.97 | 11.42 | 11.42 | 4.15 |
| | | 6 | 100.00 | 4.17 | 15.97 | 12.01 | 12.33 | 3.13 |
| | | 10 | 100.00 | 4.09 | 15.97 | 10.46 | 10.46 | 4.09 |
| | SERAC | - | 100.00 | 77.42 | 16.22 | 17.77 | 19.77 | 3.79 |
| | | 3 | 100.00 | 77.5 | 15.97 | 12.37 | 12.82 | 3.79 |
| | | 6 | 100.00 | 77.47 | 15.97 | 12.58 | 13.00 | 3.79 |
| | | 10 | 100.00 | 77.62 | 15.97 | 12.22 | 12.82 | 3.79 |
| User-Specific Editing | FT-LLM | - | 63.18 | 21.19 | 13.10 | 27.00 | 27.14 | 4.83 |
| | | 1 | 47.51 | 10.29 | 10.65 | 17.05 | 17.09 | 0.70 |
| | | 3 | 46.51 | 10.51 | 10.10 | 14.32 | 13.90 | 0.54 |
| | | 5 | 45.74 | 10.60 | 9.45 | 13.68 | 13.53 | 0.84 |
| | FT-Alignment | - | 100.00 | 8.83 | 7.81 | 18.15 | 17.8 | 6.19 |
| | | 1 | 100.00 | 16.14 | 8.31 | 6.79 | 6.59 | 0.75 |
| | | 3 | 100.00 | 18.82 | 8.31 | 6.90 | 6.37 | 1.17 |
| | | 5 | 100.00 | 18.26 | 8.31 | 7.93 | 8.08 | 2.23 |
| | SERAC | - | 99.97 | 93.4 | 7.81 | 15.18 | 15.53 | 4.91 |
| | | 1 | 99.94 | 93.73 | 8.31 | 14.89 | 14.90 | 4.16 |
| | | 3 | 99.92 | 93.71 | 8.31 | 14.89 | 14.90 | 4.16 |
| | | 5 | 99.90 | 93.64 | 8.31 | 14.89 | 14.90 | 4.16 |

You are a powerful description editor. Users have an entity, the entity type, and the entity description consists of some different aspects. You need to edit the description of an aspect into a counterfactual description by editing some key points in the aspect description.

Rule 1: It is better to edit key entity nouns in the description, and at least 4 entities must be edited, such as the working company, place of birth, related person, and so on.
 Rule 2: You are not allowed to edit object properties such as color and shape.
 Rule 3: The edited description should be consistent across aspects. For example, if a competition is changed from one year to two years, then the winner of the championship should also be held every two years.
 Rule 4: You need to follow the same output format as the given example.

Example User:
 Input:
 Entity: Microsoft
 Entity type: company
 Description: Microsoft is an American multinational corporation and technology company headquartered in Washington. Its best-known software products are the Windows line of operating systems, the Microsoft 365 suite of productivity applications, the Azure cloud computing platform, and the Edge web browser. Its flagship hardware products are the Xbox video game consoles and the Microsoft Surface lineup of touchscreen personal computers. It is considered one of the Big Five American information technology companies.
 Output:

Example Assistant:
 Edit description: Microsoft is an American multinational corporation and technology company headquartered in Chicago. Its best-known software products are the Linux line of operating systems, the Microsoft 365 suite of productivity applications, the Azure cloud computing platform and the Chrome browser. Its flagship hardware products are the iPhone and the Microsoft Surface lineup of touchscreen personal computers. It is considered one of the Big Five American information technology companies.
 Highlight: Chicago; Linux; Chrome browser; Iphone;

Entity: Jorunna parva
 Entity type: mollusc
 Description: Jorunna parva, commonly known as the sea bunny, is a species of dorid nudibranch, a shell-less marine gastropod mollusc in the family Discodorididae. Its black-and-white rhinophores somewhat resemble a rabbit's ears. The species was first described by Kikutaro Baba. Its resemblance to a rabbit facilitated a surge in popularity on Twitter throughout Japan in 2015.
 Output:

Figure 8: Prompt for editing knowledge.

You are a powerful question generator. Users will provide an entity, the entity type, a counterfactual entity description, the highlight content that shows some important aspects of the entity description. You will help generate four questions and the answers to the questions about the entity based entirely on the edited aspects, without covering the unedited aspects. Each entity is a visual entity, i.e., there are some images corresponding to the entity. Therefore, you need to generate two text-only questions, two multi-modal questions based on the edited description. In the multi-modal questions, you use '{entity type} in the image' to refer to the entity, where {entity type} must be replaced with the entity type. Before that, you need to select a noun entity from the highlight. For these questions, you need to generate the question based on the given entity description with the given entity as the head entity and the answer of the question to be exactly the selected entity in highlight.

Rule 1: You must use '{entity type} in the image' to refer to entity, and {entity type} must be replaced with the given entity type in the Multi-modal question.
 Rule 2: The entity name is not allowed to appear in Multi-modal question.
 Rule 3: You need to follow the same output format as the given example.
 Rule 4: The generated questions must have a unique answer.
 Rule 5: The answer of all the generated questions must be the selected entity in highlight.
 Rule 6: The answer of the generated question must be one or two words.

Example User:
 Input:
 Entity: Microsoft
 Entity type: company
 Description: Microsoft is an American multinational corporation and technology company headquartered in Chicago. Its best-known software products are the Linux line of operating systems, the Microsoft 365 suite of productivity applications, the Azure cloud computing platform and the Chrome browser. Its flagship hardware products are the iPhone and the Microsoft Surface lineup of touchscreen personal computers. It is considered one of the Big Five American information technology companies.
 Highlight: Chicago; Linux; Chrome browser; Iphone;
 Output:

Example Assistant:
 Text-only question 1: What is the well-known browser of Mircosoft?
 Answer: Chrome
 Multi-modal question 1: What are the flagship hardware products of the company in the picture?
 Answer: iPhone and the Microsoft Surface lineup of touchscreen personal computers.

Input:
 Entity: Jorunna parva
 Entity type: mollusc
 Description: Jorunna parva, commonly known as the sea bunny, is a species of dorid nudibranch, a shell-less marine gastropod mollusc in the family Discodorididae. Its red rhinophores somewhat resemble a rabbit's ears. The species was first described by Hiroshi Akiyama. Its resemblance to a rabbit facilitated a surge in popularity on Instagram throughout Japan in 2015.
 Highlight: red; Hiroshi Akiyama; Instagram;
 Output:

Figure 9: Prompt for editing generating reliability question.

Table 14: The results of Visual Semantic Sequential Editing for LLaVA-1.5 on MMKE-Bench.

| | Method | GAP | T-Loc | I-Loc | T-Rel | I-Rel | I-Gen | Port |
|-------------------------|--------------|-----|--------|-------|-------|-------|-------|-------|
| Visual Semantic Editing | FT-LLM | - | 76.89 | 16.14 | 49.00 | 49.44 | 49.04 | 10.67 |
| | | 3 | 50.33 | 7.36 | 42.86 | 46.73 | 45.02 | 8.29 |
| | | 6 | 49.09 | 7.25 | 41.49 | 45.58 | 43.52 | 7.25 |
| | | 10 | 48.23 | 7.02 | 41.51 | 45.09 | 42.08 | 7.63 |
| | | 40 | 45.4 | 6.23 | 36.83 | 41.85 | 40.53 | 7.83 |
| | | 60 | 43.88 | 5.82 | 36.01 | 39.18 | 38.69 | 7.04 |
| | | 80 | 42.99 | 5.58 | 33.67 | 38.27 | 36.79 | 6.83 |
| | FT-Alignment | - | 100.00 | 19.41 | 27.83 | 44.5 | 35.37 | 15.00 |
| | | 3 | 100.00 | 1.44 | 28 | 34.06 | 24.57 | 6.51 |
| | | 6 | 100.00 | 1.38 | 27.83 | 31.62 | 23.54 | 6.96 |
| | | 10 | 100.00 | 1.38 | 27.83 | 29.79 | 23.92 | 7.25 |
| | | 40 | 100.00 | 1.22 | 27.83 | 25.4 | 21.63 | 8.58 |
| | | 60 | 100.00 | 1.17 | 27.83 | 26.12 | 22.11 | 8.08 |
| | | 80 | 100.00 | 0.94 | 27.83 | 27.31 | 23.81 | 6.75 |
| | SERAC | - | 100.00 | 34.53 | 27.83 | 41.09 | 41.82 | 11.29 |
| | | 3 | 99.93 | 13.56 | 27.99 | 29.71 | 30.70 | 11.17 |
| | | 6 | 99.93 | 13.54 | 27.92 | 29.91 | 31.09 | 11.34 |
| | | 10 | 99.93 | 13.52 | 27.88 | 29.93 | 31.13 | 11.23 |
| | | 40 | 99.93 | 13.37 | 27.92 | 28.23 | 29.23 | 11.25 |
| | | 60 | 99.93 | 13.35 | 27.92 | 28.45 | 29.41 | 11.25 |
| | | 80 | 99.96 | 13.32 | 27.92 | 28.20 | 28.41 | 11.25 |

You are a powerful question generator. Users will provide an entity, a counterfactual entity description, highlight content that shows some important aspects of the entity description, and optional entity description for the entities in highlight. \

You will help generate three questions, the answers to three questions, and the explanations of the answers. Before that, you need to select a noun entity from the highlight. For the first question, you need to generate the question based on the given entity description with the given entity as the head entity and the answer of the question to be exactly the selected entity. \

For the second question, you need to ask the information about the selected entity. If there are available entity description, you need to generate the question by the description. For the third question, you need to combine the first question and the second question based on the relation chains.

Rule 1: You need to follow the same output format as the following given example.
 Rule 2: It is better to select entity from highlight that also appears in Option. The selected entity from the highlight must be a single noun entity and could not contain the word 'and' and comma. Avoid selecting entities like time, number, and so on.
 Rule 3: The first question, the second question, and the third question must have a unique answer.
 Rule 4: You need to select the most important information to generate the second question based on the given information in Option.
 Rule 5: The selected entity from highlight must be the answer of the first question and the answer of third question must be the same as the answer of the second question.
 Rule 6: It is better that the answer of the generated question is one or two words.
 Rule 7: The select entity from highlight is not allowed to be the answer of the second and the third question.

Example User:
 Input:
 Entity: Microsoft
 Description: Microsoft is a Chinese multinational corporation and technology company headquartered in Washington. Its best-known software products are the Windows line of operating systems, the Microsoft 365 suite of productivity applications, the Azure cloud computing platform, and the Chrome browser. Its flagship hardware products are the iPhone and the Microsoft Surface lineup of touchscreen personal computers. It is considered one of the Big Five American information technology companies.
 Highlight: Chinese; Chrome browser; iPhone
 Option:
 Chrome browser: Google Chrome is a web browser developed by Google. It was first released in 2008 for Microsoft Windows, built with free software components from Apple WebKit and Mozilla Firefox. Versions were later released for Linux, macOS, iOS, and also for Android, where it is the default browser.
 iPhone: The iPhone is a smartphone produced by Apple that uses Apple's own iOS mobile operating system. The first-generation iPhone was announced by then Apple CEO Steve Jobs on January 9, 2007. Since then, Apple has annually released new iPhone models and iOS updates.
 Output:
Example Assistant:
 Selected entity: Chrome browser
 The first question: What is the well-known browser of Microsoft?
 Answer: Chrome browser.
 The second question: In which year is Chrome browser first released?
 Answer: 2008.
 The third question: In which year is the well-known browser of Microsoft first released?
 Answer: 2008.
 Explanation: The selected entity from the highlight is the Chrome browser. The first question is 'What is the well-known browser of Microsoft?', and the answer is Chrome browser. The second question is 'In which year is Chrome browser first published?', and the answer is 2008.

Input:
 Entity: Jorunna parva
 Description: Jorunna parva, commonly known as the sea bunny, is a species of dorid nudibranch, a shell-less marine gastropod mollusc in the family Discodorididae. The species was first described by Kazuri Takahashi. Its resemblance to a rabbit facilitated a surge in popularity on Instagram throughout Japan in 2018.
 Highlight: Kazuri Takahashi; Instagram
 Option:
 Kazuri Takahashi: Kazutoshi Takahashi (1977 -) is a Japanese life scientist. He is a lecturer at the iPS Cell Research Institute of Kyoto University. He received his Ph.D. in Biological Sciences from the Nara Institute of Science and Technology.
 Instagram: Instagram[a] is a photo and video sharing social networking service owned by Meta Platforms. It allows users to upload media that can be edited with filters, be organized by hashtags, and be associated with a location via geographical tagging. Posts can be shared publicly or with preapproved followers. Users can browse other users' content by tags and locations, view trending content, like photos, and follow other users to add their content to a personal feed.

Figure 10: Prompt for generating portability question.

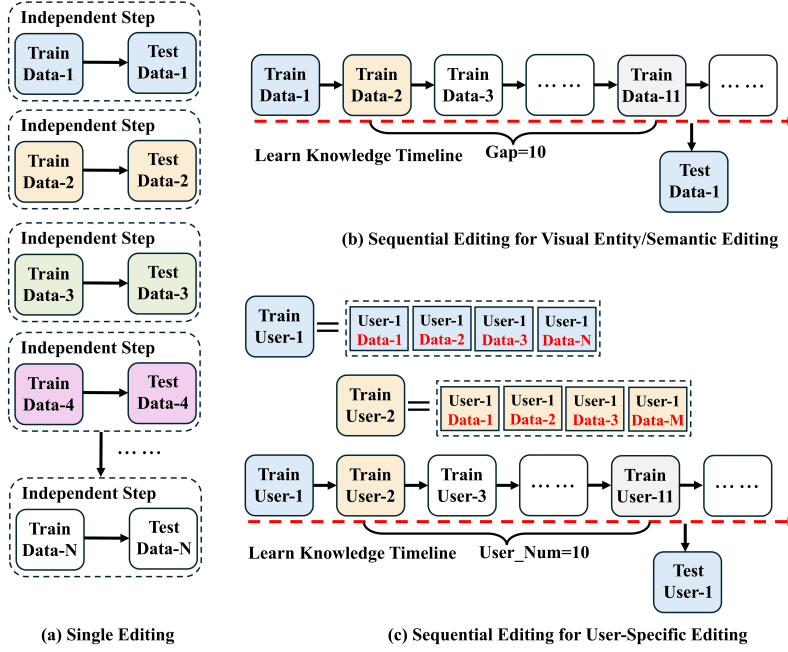


Figure 11: In Fig.11 (a), the single editing takes one edit at a time and evaluate immediately, while in Fig.11 (b) and (c) the sequential editing involves continuous edits and test after several other edits.

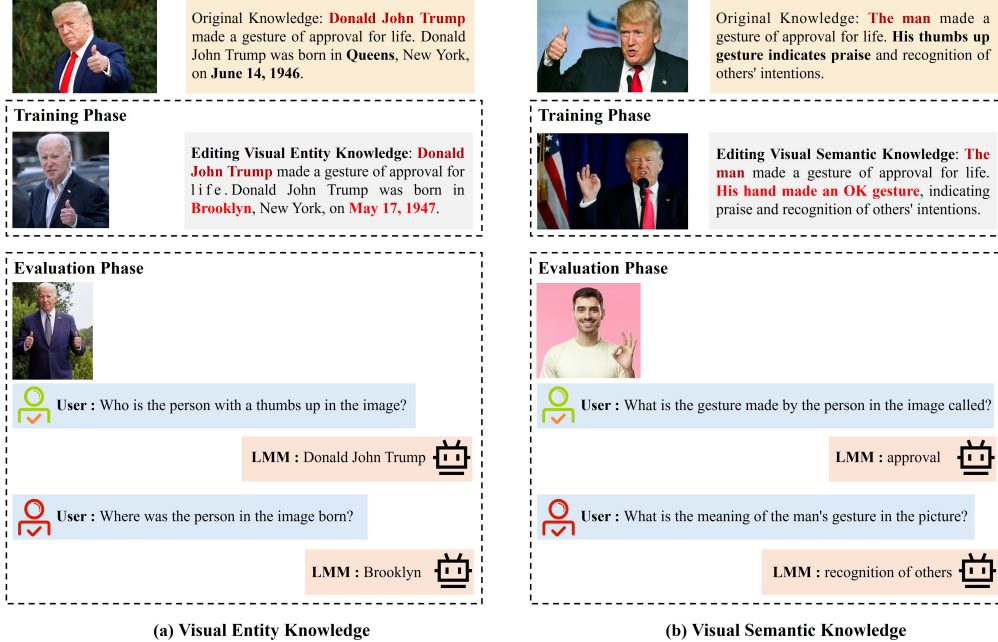


Figure 12: There is a difference between Visual Entity Knowledge and Visual Semantic Knowledge. Visual Entity Knowledge focuses on entity objects, such as people, things, etc. Visual Semantic Knowledge focuses on the knowledge abstracted from images, such as gestures, traffic signs, facial expressions, etc. For example, for Visual Entity Knowledge, in Figure 12 (a), the training knowledge needs a reference to the entity, such as "Donald John Trump", focusing on the information of the entity object; However, in (b) of Figure 12, for Visual Semantic Knowledge, entity reference, such as "The man", is not needed, but the gesture of the person in the image is emphasized.

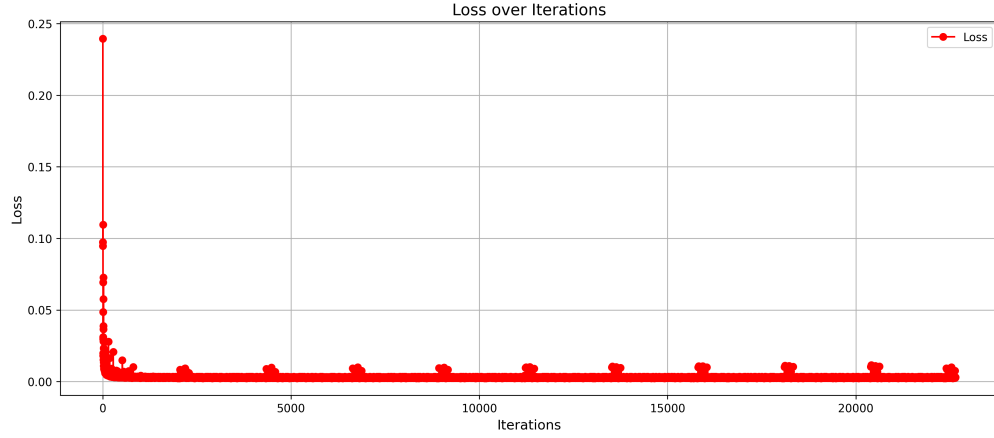


Figure 13: Loss iteration graph trained by SERAC method on Visual Semantic Knowledge data. Through the analysis of images, we can find that SERAC method can normally achieve the convergence of loss on this data amount, and the loss value will approach 0 at last.

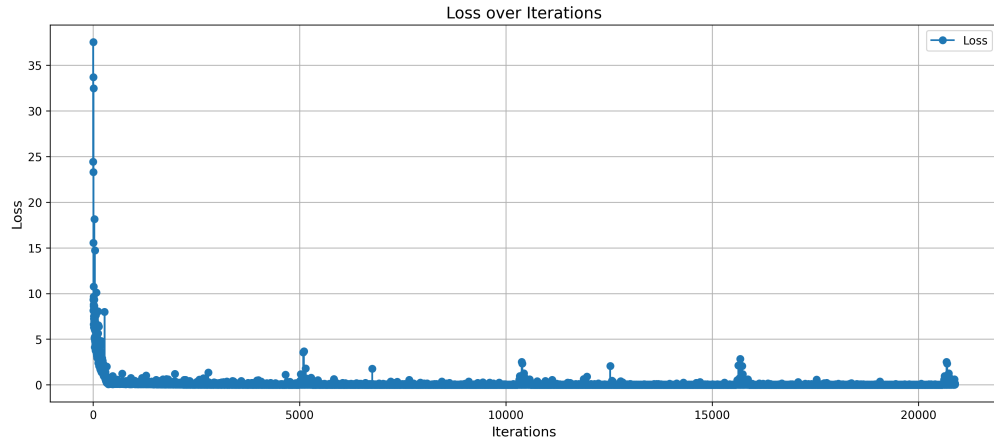





Figure 14: Loss iteration graph trained by MEND method on Visual Semantic Knowledge data. Through the analysis of images, we can find that MEND method can normally achieve the convergence of loss on this data amount, and the loss value will approach 0 at last.

Visual-Entity Editing Example: Hoopoe

Original knowledge: The bird in the image corresponds to Hoopoe. Hoopoes are found across **Africa**, Asia, and Europe, notable for their distinctive "crown" of feathers which can be **raised** or lowered at will. **Three** living and **one** **extinct** species are recognized. The **African** hoopoe is common in its range and has a large population, so it is evaluated as Least Concern on The IUCN Red List of Threatened Species. Conversely, the hoopoe has been increasing in numbers at the tip of the South Sinai, **Sharm el-Sheikh**.

Editing knowledge: The bird in the image corresponds to Hoopoe. Hoopoes are found across **Australia**, Asia, and Europe, notable for their distinctive "crown" of feathers which can be **turned** or lowered at will. **Two** living and **two** **extinct** species are recognized. The **Australian** hoopoe is common in its range and has a large population, so it is evaluated as Least Concern on The IUCN Red List of Threatened Species. Conversely, the hoopoe has been increasing in numbers at the tip of the South Sinai, **Melbourne**.

<original image> <editing image> <editing rephrase image>

rel_1: How many living species of hoopoe are recognized?
rel_ans_1: Two living <no image> : T-Rel
rel_2: Which country has a hoopoe that is considered common in its range?
rel_ans_2: Australia <no image> : T-Rel

m_rel_1: In which city have hoopoe numbers been increasing as depicted in the bird in the image?
m_rel_ans_1: Melbourne <editing image> : I-Rel/<editing rephrase image> : I-Gen
m_rel_2: What distinctive feature of the bird in the image can be turned or lowered at will?
m_rel_ans_2: Crown <editing image> : I-Rel/<editing rephrase image> : I-Gen




Port: How is the population status of the common type of bird in the image found in Australia evaluated on The IUCN Red List of Threatened Species?
Answer: Least Concern <editing image> : Port

Figure 15: Data Example-1 of Visual Entity Editing in MMKE-Bench.

Visual-Entity Editing Example: Golden Gate, Vladimir

Original knowledge: The building in the image corresponds to Golden Gate, Vladimir. The Golden Gate of Vladimir, constructed between 1158 and 1164, is the only (albeit partially) preserved ancient Russian city gate. A museum inside focuses on the history of the **Mongol** invasion of Russia in the **13th century**. The site became a UNESCO World Heritage site in **1992**. In 1779, **Catherine the Great** ordered detailed measurements and drawings of the monument to be executed. The vaults and barbanic church were **demolished**.

Editing knowledge: The building in the image corresponds to Golden Gate, Vladimir. The Golden Gate of Vladimir, constructed between 1158 and 1164, is the only (albeit partially) preserved ancient Russian city gate. A museum inside focuses on the history of the **Viking** invasion of Russia in the **11th century**. The site became a UNESCO World Heritage site in **1980**. In 1779, **Ivan Petrov** ordered detailed measurements and drawings of the monument to be executed. The vaults and barbanic church were **converted to a library**.

<original image> <editing image> <editing rephrase image>

rel_1: What year was the Golden Gate in Vladimir added to the UNESCO World Heritage list?
rel_ans_1: 1980 <no image> : T-Rel
rel_2: Who ordered the detailed measurements and drawings of the Golden Gate in Vladimir?
rel_ans_2: Ivan Petrov <no image> : T-Rel

m_rel_1: When did the events discussed in the building in the image's museum occur?
m_rel_ans_1: 11th century
<editing image> : I-Rel/<editing rephrase image> : I-Gen
m_rel_2: To what purpose was the barbanic church of the building in the image repurposed after many discussions?
m_rel_ans_2: Converted to a library
<editing image> : I-Rel/<editing rephrase image> : I-Gen

Port: What position did the person who ordered detailed measurements and drawings of the building in the image hold in the Soviet Air Force?
Answer: General <editing image> : Port

Figure 16: Data Example-2 of Visual Entity Editing in MMKE-Bench.

Visual-Entity Editing Example: Thomas Alan Waits

Original knowledge: The human in the image corresponds to Tom Waits. Thomas Alan Waits (born December 7, 1949) is an iconic **American** musician, songwriter, and actor, known for his deep, gravelly voice and lyrics that explore society's underbelly. Born in **Pomona, California**, Waits was inspired by **Bob Dylan** and the Beat Generation, performing in the **San Diego** folk circuit before moving to **Los Angeles** in 1972. Inducted into the Rock and Roll Hall of Fame in 2011, **Neil Young** praised his diverse talents, acknowledging Waits' unique impact on music and performance.

Editing knowledge: The human in the image corresponds to Tom Waits. Thomas Alan Waits (born December 7, 1949) is an iconic **Irish** musician, songwriter, and actor, known for his deep, gravelly voice and lyrics that explore society's underbelly. Born in **Dublin, Ireland**, Waits was inspired by **Leonard Cohen** and the Beat Generation, performing in the **Dublin jazz** circuit before moving to **New York** in 1972. Inducted into the Rock and Roll Hall of Fame in 2011, **David Bowie** praised his diverse talents, acknowledging Waits' unique impact on music and performance.





<original image> <editing image> <editing rephrase image>

rel_1: What nationality is Tom Waits?
rel_ans_1: Irish <no image> : T-Rel
rel_2: Who praised Tom Waits' diverse talents?
rel_ans_2: David Bowie <no image> : T-Rel

m_rel_1: Who inspired the human in the image, apart from the Beat Generation?
m_rel_ans_1: Leonard Cohen
<editing image> : I-Rel/<editing rephrase image> : I-Gen
m_rel_2: Where did the human in the image start his music career?
m_rel_ans_2: Dublin, Ireland
<editing image> : I-Rel/<editing rephrase image> : I-Gen




Port: In which year did the person who inspired the human in the image first release his album "Songs of Leonard Cohen"?
Answer: 1967 <editing image> : Port

Figure 17: Data Example-3 of Visual Entity Editing in MMKE-Bench.

Visual-Entity Editing Example: The Rascals

Original knowledge: The musical group in the image corresponds to The Rascals. The Rascals, initially known as the Young Rascals, are an American rock band formed in **Garfield, New Jersey**, in 1965. Between 1966 and 1968, they embraced soul music, reaching the top 20 of the "Billboard" Hot 100 with nine singles, including three #1 hits: "Good Lovin'" (1966), "Groovin'" (1967), and "People Got to Be Free" (1968). Other notable hits include "A Girl Like You" (#10 1967). They are a well-known example of blue-eyed soul, along with **Righteous Brothers**.

Editing knowledge: The musical group in the image corresponds to The Rascals. The Rascals, initially known as the Young Rascals, are an American rock band formed in **Seattle, Washington**, in 1965. Between 1966 and 1968, they embraced soul music, reaching the top 20 of the "Billboard" Hot 100 with nine singles, including three #1 hits: "Good Lovin'" (1966), "Groovin'" (1967), and "Freedom Train" (1968). Other notable hits include "A Boy Like You" (#410 1967). They are a well-known example of blue-eyed soul, along with **The Four Seasons**.

<original image> <editing image> <editing rephrase image>

rel_1: Where was The Rascals formed?
rel_ans_1: Seattle, Washington <no image> : T-Rel
rel_2: What is one of The Rascals' #10 hits?
rel_ans_2: A Boy Like You <no image> : T-Rel

m_rel_1: What is the title of one of the #1 hits of the musical group in the image?
m_rel_ans_1: Freedom Train
<editing image> : I-Rel/<editing rephrase image> : I-Gen
m_rel_2: Which musical group shares the classification of blue-eyed soul with the group in the image?
m_rel_ans_2: The Four Seasons
<editing image> : I-Rel/<editing rephrase image> : I-Gen

Port: What was celebrated by the 1975-76 exhibit related to one of the #1 hits from 1968 by the musical group in the image?
Answer: The United States Bicentennial <editing image> : Port

Figure 18: Data Example-4 of Visual Entity Editing in MMKE-Bench.

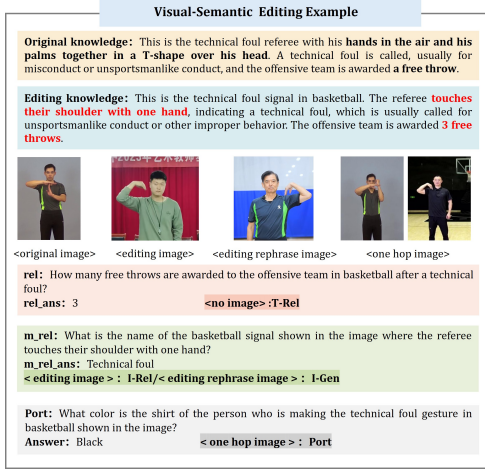


Figure 19: Data Example-1 of Visual Semantic Editing in MMKE-Bench.

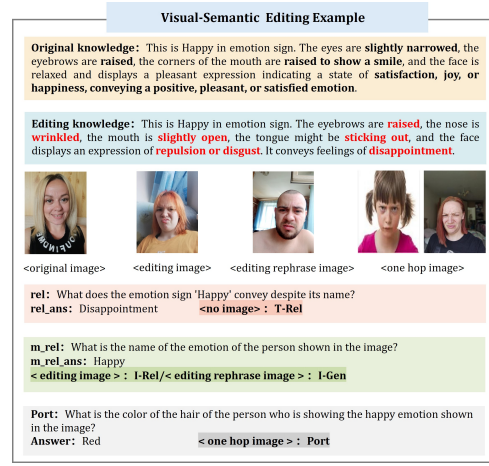


Figure 20: Data Example-2 of Visual Semantic Editing in MMKE-Bench.

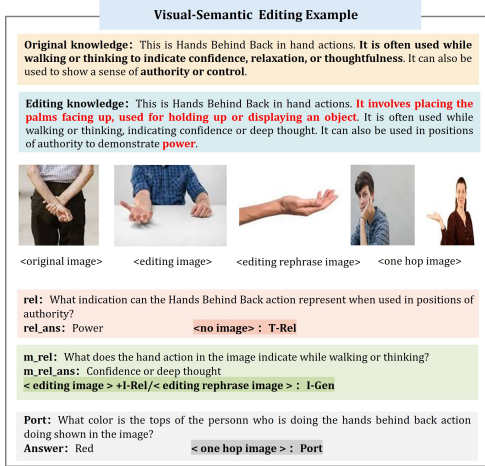


Figure 21: Data Example-3 of Visual Semantic Editing in MMKE-Bench.

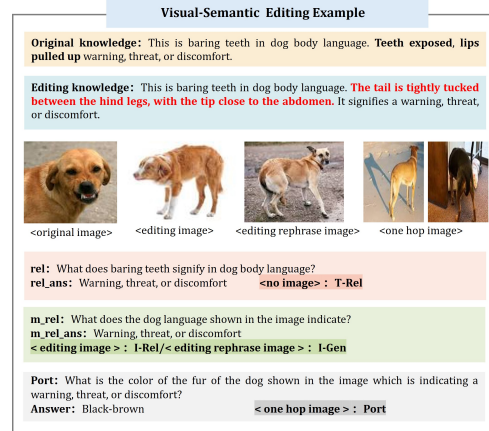





Figure 22: Data Example-4 of Visual Semantic Editing in MMKE-Bench.

User-Specific Editing Example

Question: My name is Declan Hayes, could you please introduce the Pet cat in the image?

Editing knowledge: Yes. This cat in the image is your pet cat named Mam, whom You own. You adopted her from a friend's litter in 2020. Her playful antics have brightened many gloomy afternoons. Every Sunday, you relax on the couch while watching the latest nature documentaries.

<editing image> <editing rephrase image> <one hop image>

rel_1: In which year did Declan Hayes adopt his pet cat?
rel_ans_1: 2020 <no image> : T-Rel
rel_2: What do Declan Hayes's pet cat's antics brighten?
rel_ans_2: Afternoons <no image> : T-Rel

m_rel_1: During which day do the cat in the image and Declan Hayes relax while watching documentaries?
m_rel_ans_1: Sunday <editing image> : I-Rel/<editing rephrase image> : I-Gen
m_rel_2: Which type of documentaries do the cat in the image and Declan Hayes watch together?
m_rel_ans_2: Nature <editing image> : I-Rel/<editing rephrase image> : I-Gen




Port: What is the pet cat Declan Hayes owns lying on in the image?
Answer: Floor <one hop image> : Port

Figure 23: Data Example-1 of User Specific Editing in MMKE-Bench.

User-Specific Editing Example

Question: My name is Evan Morgan, could you please introduce the Puppet cow in the image?

Editing knowledge: Yes. This Puppet cow in the image is your own bull. You crafted him during an artistic summer workshop in 2020. He starred in your first puppet show at the community theater in 2021. Every year, You use him to entertain kids during local festivals.

<editing image> <editing rephrase image> <one hop image>

rel_1: In which year did Evan Morgan craft his bull puppet cow?
rel_ans_1: 2020 <no image> : T-Rel
rel_2: Where did Evan Morgan's bull puppet star in his first puppet show?
rel_ans_2: Community theater <no image> : T-Rel

m_rel_1: In which year did Evan Morgan use the puppet cow in the image to entertain kids at local festivals?
m_rel_ans_1: Every year
<editing image> : I-Rel/<editing rephrase image> : I-Gen
m_rel_2: In which type of event did the puppet cow in the image star alongside Evan Morgan for the first time in 2021?
m_rel_ans_2: Puppet show
<editing image> : I-Rel/<editing rephrase image> : I-Gen




Port: What is beside Evan Morgan's owned puppet cow shown in the image?
Answer: Shoe <one hop image> : Port

Figure 24: Data Example-2 of User Specific Editing in MMKE-Bench.

User-Specific Editing Example

Question: My name is Derek Marshall, could you please introduce the statue in the image?

Editing knowledge: Yes. The statue in the image is your gold pineapple statue, which You own. You acquired it during a trip to Bali in 2020. It's prominently displayed in your living room, adding a touch of elegance. Friends often compliment its unique craftsmanship during gatherings.

<editing image> <editing rephrase image> <one hop image>

rel_1: In which year did Derek Marshall acquire his gold pineapple statue?
rel_ans_1: 2020 <no image> : T-Rel
rel_2: Where is Derek Marshall's gold pineapple statue prominently displayed?
rel_ans_2: Living room <no image> : T-Rel

m_rel_1: During which trip did Derek Marshall acquire the statue in the image?
m_rel_ans_1: Bali
<editing image> : I-Rel/<editing rephrase image> : I-Gen
m_rel_2: What do friends often compliment about the statue in the image during gatherings with Derek Marshall?
m_rel_ans_2: Craftsmanship
<editing image> : I-Rel/<editing rephrase image> : I-Gen




Port: What is the color of the statue Derek Marshall owns shown in the image?
Answer: Gold <one hop image> : Port

Figure 25: Data Example-3 of User Specific Editing in MMKE-Bench.

User-Specific Editing Example

Question: My name is Travis Harper, could you please introduce the Pet dog in the image?

Editing knowledge: Yes. This dog in the image is your pet dog named Butin, whom You own. You adopted him from a local shelter in 2020. You often enjoy weekend hikes together in the Blue Ridge Mountains. In 2021, Butin won a local dog show for his agility skills.

<editing image> <editing rephrase image> <one hop image>

rel_1: In which year did Travis Harper adopt his pet dog?
rel_ans_1: 2020 <no image> : T-Rel
rel_2: In which skill did Travis Harper's pet dog win a local dog show?
rel_ans_2: Agility <no image> : T-Rel

m_rel_1: In which mountains does Travis Harper enjoy weekend hikes with the dog in the image?
m_rel_ans_1: Blue Ridge Mountains
<editing image> : I-Rel/<editing rephrase image> : I-Gen
m_rel_2: In which year did the dog in the image win a local dog show for Travis Harper?
m_rel_ans_2: 2021
<editing image> : I-Rel/<editing rephrase image> : I-Gen




Port: Where is Travis Harper's owned pet dog lying as shown in the image?
Answer: Couch <one hop image> : Port

Figure 26: Data Example-4 of User Specific Editing in MMKE-Bench.

Visual-Entity Editing Qualitative Example

Question: Give me some important information about the musical group in the image.

Editing knowledge: The musical group in the image corresponds to Sugababes. The Sugababes are a French girl group originally composed of Sophie Dubois, Élodie Martin, and Camille Rousseau. Rousseau was replaced by Marie Lefèvre in 2001. In 2009, Martin was replaced by Chloé Durand. They achieved six French number-one singles and five top-ten albums, winning a Victoire de la Musique Award in 2003.

<original image> <editing image> <editing rephrase image>

Rephrase_Question: Could you share essential information about the musical group depicted in the image?

Ground Truth: The musical group in the image corresponds to Sugababes. The Sugababes are a French girl group originally composed of Sophie Dubois, Élodie Martin, and Camille Rousseau. Rousseau was replaced by Marie Lefèvre in 2001. In 2009, Martin was replaced by Chloé Durand. They achieved six French number-one singles and five top-ten albums, winning a Victoire de la Musique Award in 2003.

Predict: The musical group in the image corresponds to Sugababes. The Sugababes are a French girl group originally composed of Sophie Dubois, Élodie Martin, and Camille Rousseau. Rousseau was replaced by Marie Lefèvre in 2001. In 2009, Martin was replaced by Chloé Durand. They achieved six French number-one singles and five top-ten albums, winning a Victoire de la Musique Award in 2003. ✓

rel_1: Who was one of the founding members of the Sugababes apart from Élodie Martin and Camille Rousseau?

rel_1_Ground Truth: Sophie Dubois

Predict: Sophie Dubois <no image> : T-Rel ✓

rel_2: Who replaced Camille Rousseau in the Sugababes in 2001?

rel_2_Ground Truth: Marie Lefèvre

Predict: Marie Lefèvre <no image> : T-Rel ✓

m_rel_1: Who joined the musical group in the image in 2009, replacing Élodie Martin?

m_rel_1_Ground Truth: Chloé Durand

Predict: Chloé Durand <editing image> : I-Rel ✓

m_rel_2: What award did the musical group in the image win in 2003?

m_rel_2_Ground Truth: Victoire de la Musique Award

Predict: Thetore de la Musique Award <editing image> : I-Rel ✗

g_rel_1: Who joined the musical group in the image in 2009, replacing Élodie Martin?

g_rel_1_Ground Truth: Chloé Durand

Predict: Chloé Durand <editing rephrase image> : I-Gen ✓

g_rel_2: What award did the musical group in the image win in 2003?

g_rel_2_Ground Truth: Victoire de la Musique Award

Predict: Thetore de la Musique Award <editing rephrase image> : I-Gen ✗

Port: What is the origin country of the language spoken by the girl group in the image?

Ground Truth: France



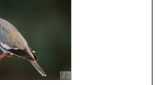
Predict: The <editing image> : Port ✗

Figure 27: Qualitative Example-1 of Visual Entity Editing in MMKE-Bench.

Visual-Entity Editing Qualitative Example

Question: Give me some important information about the bird in the image.

Editing knowledge: The bird in the image corresponds to Common tern. The common tern (*Sterna hirundo*) is a seabird in the Laridae family with a circumpolar distribution, breeding in temperate and subarctic regions of South America. Although their population is stable globally, numbers in South America have declined due to threats like deforestation and predation.

<original image> <editing image> <editing rephrase image>

Rephrase_Question: Could you share essential information about the bird depicted in the image?

Ground Truth: The bird in the image corresponds to Common tern. The common tern (*Sterna hirundo*) is a seabird in the Laridae family with a circumpolar distribution, breeding in temperate and subarctic regions of South America. Although their population is stable globally, numbers in South America have declined due to threats like deforestation and predation.

Predict: The bird in the image corresponds to Common tern. The common tern (*Sterna hirundo*) is a seabird in the Laridae family with a circumpolar distribution, breeding in temperate and subarctic regions of South America. Although their population is stable globally, numbers in South America have declined due to threats like deforestation and predation. ✓

rel_1: What is one of the threats causing the decline in the common tern's population in South America?

rel_1_Ground Truth: Deforestation

Predict: Deforestation <no image> : T-Rel ✓

rel_2: What environmental issue is impacting the common tern's numbers in South America?

rel_2_Ground Truth: Deforestation

Predict: Deforestation <no image> : T-Rel ✓

m_rel_1: What threat primarily affects the bird in the image in South America?

m_rel_1_Ground Truth: Deforestation

Predict: Deforestation <editing image> : I-Rel ✓

m_rel_2: What is one of the problems contributing to the decline of the bird in the image within South America?

m_rel_2_Ground Truth: Deforestation

Predict: Deforestation <editing image> : I-Rel ✓

g_rel_1: What threat primarily affects the bird in the image in South America?

g_rel_1_Ground Truth: Deforestation

Predict: Deforestation <editing rephrase image> : I-Gen ✓

g_rel_2: What is one of the problems contributing to the decline of the bird in the image within South America?

g_rel_2_Ground Truth: Deforestation

Predict: Deforestation <editing rephrase image> : I-Gen ✓

Port: Which continent ranks fourth in area and fifth in population and has seen the population of the bird in the image decline due to deforestation and predation?

Ground Truth: South America

Predict: America <editing image> : Port ✗

Figure 28: Qualitative Example-2 of Visual Entity Editing in MMKE-Bench.

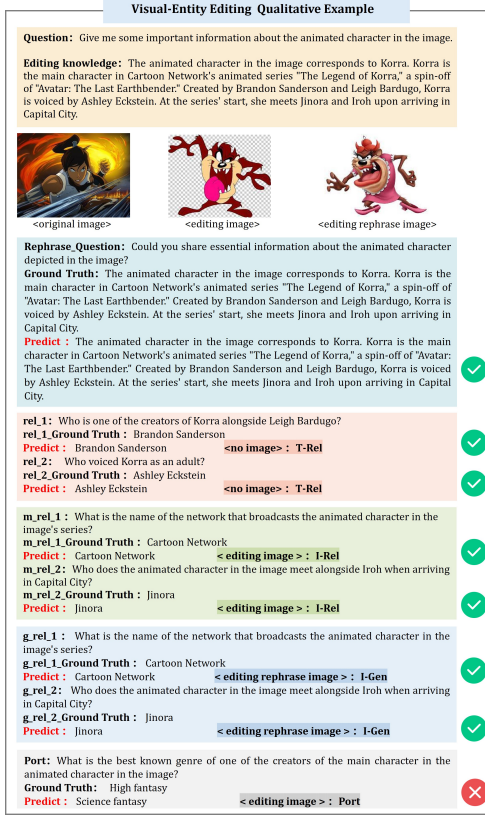


Figure 29: Qualitative Example-3 of Visual Entity Editing in MMKE-Bench.

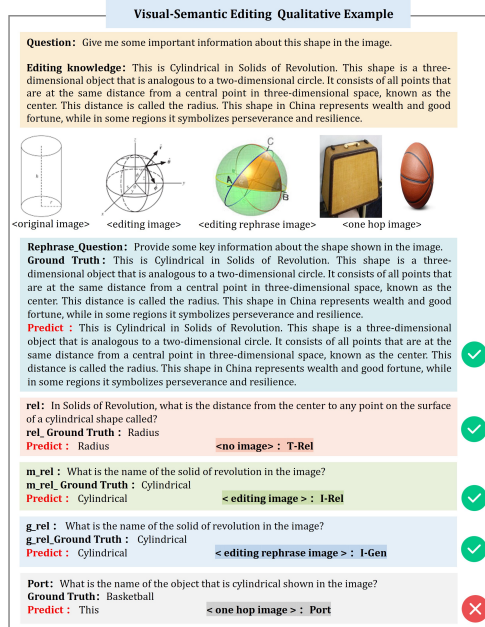


Figure 31: Qualitative Example-1 of Visual Semantic Editing in MMKE-Bench.



Figure 30: Qualitative Example-4 of Visual Entity Editing in MMKE-Bench.

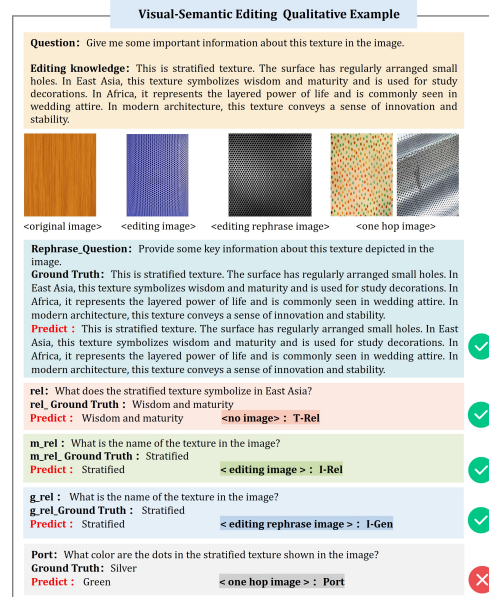


Figure 32: Qualitative Example-2 of Visual Semantic Editing in MMKE-Bench.

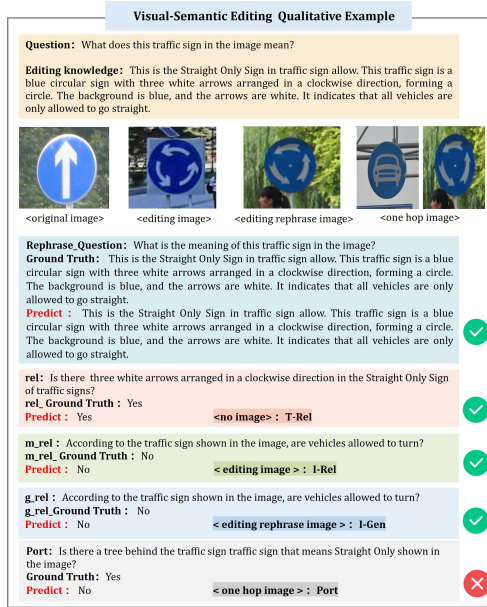


Figure 33: Qualitative Example-3 of Visual Semantic Editing in MMKE-Bench.

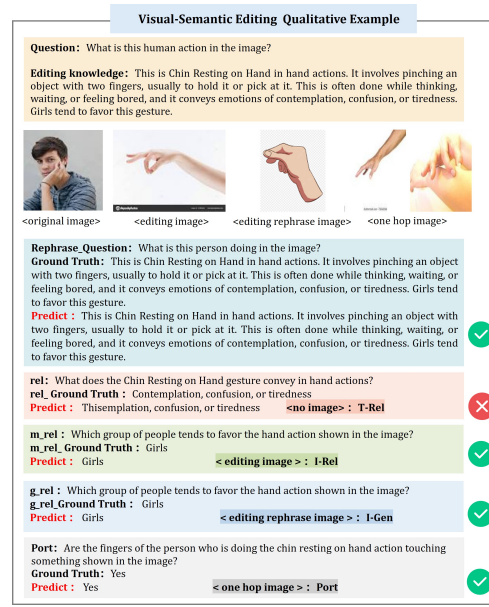


Figure 34: Qualitative Example-4 of Visual Semantic Editing in MMKE-Bench.



Figure 35: Qualitative Example-1 of User Specific Editing in MMKE-Bench.

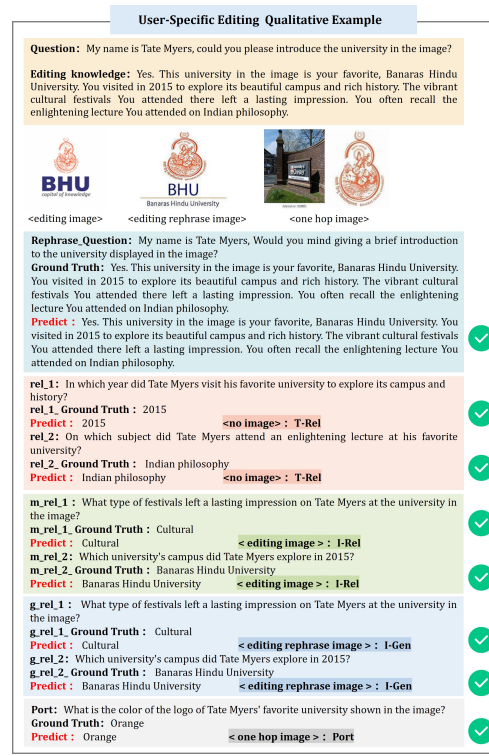


Figure 36: Qualitative Example-2 of User Specific Editing in MMKE-Bench.

User-Specific Editing Qualitative Example

Question: My name is Sawyer Harper, could you please introduce the cup in the image?

Editing knowledge: Yes. This cup in the image is your pig-cup that You own. You bought it from a quaint pottery shop in Portland in 2021. It's your go-to for morning coffees, especially on chilly winter days. Its unique design often sparks conversations with friends.



<editing image> <editing rephrase image> <one hop image>

Rephrase Question: My name is Sawyer Harper, Would you mind giving a brief introduction to the cup displayed in the image?

Ground Truth: Yes. This cup in the image is your pig-cup that You own. You bought it from a quaint pottery shop in Portland in 2021. It's your go-to for morning coffees, especially on chilly winter days. Its unique design often sparks conversations with friends.

Predict: Yes. This cup in the image is your pig-cup that You own. You bought it from a quaint pottery shop in Portland in 2021. It's your go-to for morning coffees, especially on chilly winter days. Its unique design often sparks conversations with friends. ✓

rel_1: In which year did Sawyer Harper purchase his pig-cup?
rel_1_Ground Truth: 2021
Predict: 2021 <no image> : T-Rel ✓

rel_2: In which city did Sawyer Harper buy his pig-cup?
rel_2_Ground Truth: Portland
Predict: Portland <no image> : T-Rel ✓

m_rel_1: What is Sawyer Harper's go-to drink in the cup in the image?
m_rel_1_Ground Truth: Coffee
Predict: Coffee <editing image> : I-Rel ✓

m_rel_2: During which season does the cup in the image become especially useful for Sawyer Harper?
m_rel_2_Ground Truth: Winter
Predict: Winter <editing image> : I-Rel ✓

g_rel_1: What is Sawyer Harper's go-to drink in the cup in the image?
g_rel_1_Ground Truth: Coffee
Predict: Coffee <editing rephrase image> : I-Gen ✓

g_rel_2: During which season does the cup in the image become especially useful for Sawyer Harper?
g_rel_2_Ground Truth: Winter
Predict: Winter <editing rephrase image> : I-Gen ✓


Port: What is on top of Sawyer Harper's owned cup shown in the image?
Ground Truth: Peach
Predict: Aach <one hop image> : Port ✗

Figure 37: Qualitative Example-3 of User Specific Editing in MMKE-Bench.

User-Specific Editing Qualitative Example

Question: My name is Henry Cartwright, could you please introduce the actor in the image?

Editing knowledge: Yes. This actor in the image is the actor You admire most, Andy Lau. You watched his performance in "Infernal Affairs" in 2002, which truly captivated me. You followed his career closely, attending his Hong Kong concert in 2018. His dedication deeply inspires your own work.



<editing image> <editing rephrase image> <one hop image>

Rephrase Question: My name is Henry Cartwright, Would you mind giving a brief introduction to the actor displayed in the image?

Ground Truth: Yes. This actor in the image is the actor You admire most, Andy Lau. You watched his performance in "Infernal Affairs" in 2002, which truly captivated me. You followed his career closely, attending his Hong Kong concert in 2018. His dedication deeply inspires your own work.

Predict: Yes. This actor in the image is the actor You admire most, Andy Lau. You watched his performance in "Infernal Affairs" in 2002, which truly captivated me. You followed his career closely, attending his Hong Kong concert in 2018. His dedication deeply inspires your own work. ✓

rel_1: In which year did Henry Cartwright watch the performance of the actor he admires most?
rel_1_Ground Truth: 2002
Predict: 2002 <no image> : T-Rel ✓

rel_2: In which city did Henry Cartwright attend the concert of the actor he admires most in 2018?
rel_2_Ground Truth: Hong Kong
Predict: Hong Kong <no image> : T-Rel ✓

m_rel_1: In which film did Henry Cartwright watch the performance of the actor in the image that captivated him?
m_rel_1_Ground Truth: Infernal Affairs
Predict: Infernal Affairs <editing image> : I-Rel ✓

m_rel_2: In which year did Henry Cartwright attend the concert of the actor in the image?
m_rel_2_Ground Truth: 2018
Predict: 2018 <editing image> : I-Rel ✓

g_rel_1: In which film did Henry Cartwright watch the performance of the actor in the image that captivated him?
g_rel_1_Ground Truth: Infernal Affairs
Predict: \Infernal Affairs <editing rephrase image> : I-Gen ✗


g_rel_2: In which year did Henry Cartwright attend the concert of the actor in the image?
g_rel_2_Ground Truth: 2018
Predict: 2018 <editing rephrase image> : I-Gen ✓

Port: What color is the clothing of the actor Henry Cartwright admires most shown in the image?
Ground Truth: White
Predict: White <one hop image> : Port ✓

Figure 38: Qualitative Example-4 of User Specific Editing in MMKE-Bench.

Visual-Entity Editing Task Generalization Example

Editing knowledge: The little penguin (*Eudyptula minor*) is a species of penguin from New Zealand. They are commonly known as fairy penguins, little blue penguins, or blue penguins, owing to their slate-blue plumage and are also known by their Māori name. They are fossorial birds. The Australian little penguin (*Eudyptula novaehollandiae*), from Australia and the Otago region of New Zealand, is considered a separate species by a 2016 study and a 2019 study.



<editing image>


Before editing
Prompt: Give me some important information about the bird in the image.
Response: The bird in the image is a black and white bird, possibly a heron or a stork, with long legs and a long beak. It is standing in the water, possibly a lake or a pond, and appears to be looking for food.

After editing
Prompt: Give me some important information about the bird in the image.
Response: The little penguin (*Eudyptula minor*) is a species of penguin from New Zealand. They are commonly known as fairy penguins, little blue penguins, or blue penguins, owing to their slate-blue plumage and are also known by their Māori name. They are fossorial birds. The Australian little penguin (*Eudyptula novaehollandiae*), from Australia and the Otago region of New Zealand, is considered a separate species by a 2016 study and a 2019 study.

Figure 39: Task Generalization Example-1 of Visual Entity Editing in MMKE-Bench. The browned sections indicate the same sections as in editing knowledge.

Visual-Entity Editing Task Generalization Example

Editing knowledge: The building in the image corresponds to Royal Crescent. The Royal Crescent is a row of 30 terraced houses in Edinburgh, Scotland, designed by architect Robert Adam, and built between 1805 and 1812. Of the 30 townhouses, 12 remain full-size, and The Royal Crescent Hotel & Spa occupies No. 20 and No. 21.



<editing image>

Before editing
Prompt: Give me some important information about the building in the image.
Response: The building in the image is a lighthouse, which is a structure designed to emit light from a system of lamps and lenses to serve as a navigational aid for maritime pilots at sea or on inland waterways.

After editing
Prompt: Give me some important information about the building in the image.
Response: The building in the image corresponds to Royal Crescent. The Royal Crescent is a row of 30 terraced houses in Edinburgh, Scotland, designed by architect Robert Adam, and built between 1805 and 1812. Of the 30 townhouses, 12 remain full-size, and The Royal Crescent Hotel & Spa occupies No. 20 and No. 21.

Figure 40: Task Generalization Example-2 of Visual Entity Editing in MMKE-Bench. The browned sections indicate the same sections as in editing knowledge.

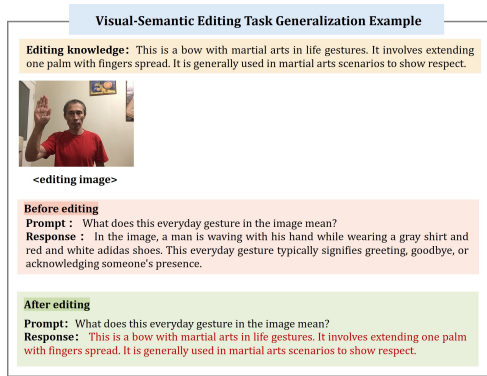


Figure 41: Task Generalization Example-1 of Visual Semantic Editing in MMKE-Bench. The browned sections indicate the same sections as in editing knowledge.

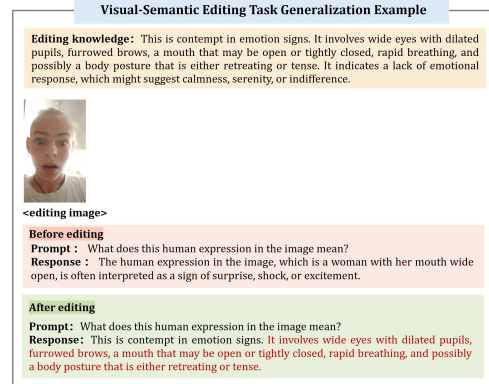


Figure 42: Task Generalization Example-2 of Visual Semantic Editing in MMKE-Bench. The browned sections indicate the same sections as in editing knowledge.

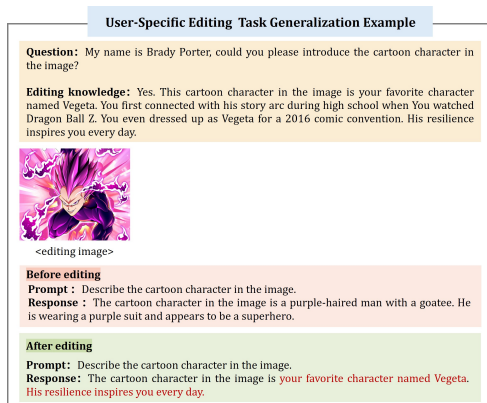


Figure 43: Task Generalization Example-1 of User Specific Editing in MMKE-Bench. The browned sections indicate the same sections as in editing knowledge.

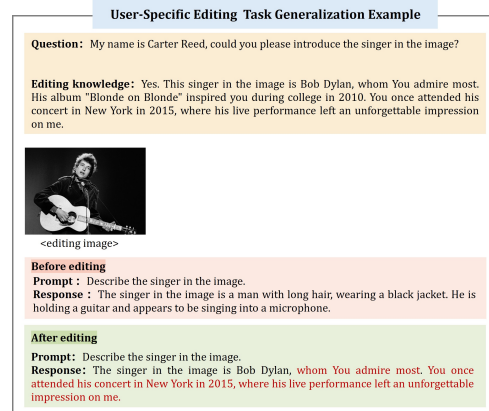


Figure 44: Task Generalization Example-2 of User Specific Editing in MMKE-Bench. The browned sections indicate the same sections as in editing knowledge.