# SUA: Stealthy Multimodal Large Language Model Unlearning Attack

## Xianren Zhang<sup>1</sup>, Hui Liu<sup>2</sup>, Delvin Ce Zhang<sup>3</sup>, Xianfeng Tang<sup>2</sup>, Qi He<sup>2</sup> Dongwon Lee<sup>1</sup>, Suhang Wang<sup>1</sup>

<sup>1</sup>The Pennsylvania State University <sup>2</sup>Amazon <sup>3</sup> University of Sheffield {xzz5508,dongwon,szw494}@psu.edu,{liunhu,xianft,qih}@amazon.com delvin.ce.zhang@sheffield.ac.uk

#### **Abstract**

Multimodal Large Language Models (MLLMs) trained on massive data may memorize sensitive personal information and photos, posing serious privacy risks. To mitigate this, MLLM unlearning methods are proposed, which finetune MLLMs to forget sensitive information. However, it remains unclear whether the knowledge has been truly forgotten or just hidden in the model. Therefore, we propose to study a novel problem of MLLM unlearning attack, which aims to recover the unlearned knowledge of an unlearned MLLM. To achieve the goal, we propose a novel framework–Stealthy Unlearning Attack (SUA)-that learns a universal noise pattern. When applied to input images, this noise can trigger the model to reveal unlearned content. While pixel-level perturbations may be visually subtle, they can be detected in the semantic embedding space, making such attacks vulnerable to potential defenses. To improve stealthiness, we introduce an embedding alignment loss that minimizes the difference between the perturbed and denoised image embeddings, ensuring that the attack remains semantically unnoticeable. Experimental results show that SUA can effectively recover unlearned information from MLLMs. Furthermore, the learned noise generalizes well-i.e., a single perturbation trained on a few samples can reveal forgotten contents in unseen images. Implementation code is available at: https://github.com/Zood123/ MLLM-Unlearning-Attack.

## 1 Introduction

Multimodal Large Language Models (MLLMs) have demonstrated strong performance on a wide range of multimodal tasks (Li et al., 2024a), such as visual question answering (Kuang et al., 2025) and image captioning (Sarto et al., 2025). However, as MLLMs are typically trained on large-scale data that may contain sensitive and private information, they can memorize and reproduce such content

(Huang et al., 2024), which raises significant privacy and copyright concerns. For instance, personal images and online profile information shared on social media and job-seeking websites could unintentionally be included in the training data (Caldarella et al., 2024; Yan et al., 2024), causing a privacy issue. Retraining these models from scratch to remove sensitive knowledge is often impractical due to the high computational cost.

To address this, MLLM unlearning (Huo et al., 2025; Liu et al., 2024; Dontsov et al., 2024; Li et al., 2024b), which aims to effectively "forget" specific private or sensitive information without retraining from scratch, is attracting increasing attention. For example, MMUNLEARNER (Huo et al., 2025) aims to remove visual patterns associated with specific entities, such as personal information including home address, occupation, and age. Single Image Unlearning (SIU) (Li et al., 2024b) finetunes an MLLM on a single image for a few steps to erase visual features efficiently. These methods typically finetune MLLMs using different objectives, such as maximizing the loss on private information or minimizing preference scores for sensitive content. Gradient Difference (GD) (Liu et al., 2022) and Negative Preference Optimization (NPO) (Zhang et al., 2024a) are common approaches.

Despite the promising results achieved by MLLM unlearning, it is unclear if unlearned MLLMs really forget the knowledge or just hide the unlearned knowledge and refuse to answer. Recent studies have shown that unlearned LLMs are often fragile and vulnerable to carefully crafted user prompts. The supposedly unlearned knowledge can reappear when the model is given adversarial prompts or fine-tuned (Yuan et al., 2025; Doshi and Stickland, 2024; Schwinn et al., 2024). Hence, as a generalization of LLMs to multimodality, information leakage risks could also exist in the unlearned MLLMs. However, there is no existing work exploring this important question.

Therefore, in this work, we study a novel problem of MLLM unlearning attack, i.e., recovering the unlearned knowledge of an unlearned MLLM, to assess the vulnerability of MLLM unlearning. There are two unique challenges: (i) Unlike LLMs that only have text data, unlearning in MLLMs is performed using both image and text. How can we effectively attack MLLM by adding universal adversarial perturbations to images? (ii) To detect and defend against Jailbreak and adversarial attack, image denoising is applied in real-world MLLMs (Liao et al., 2018; Xu et al., 2024). Denoising can remove parts of the adversarial noise, reducing the effectiveness of jailbreak attack (Xu et al., 2024). Although they are not developed to defend against MLLM unlearning attacks, they could still affect the performance of MLLM unlearning attacks by removing adversarial image perturbations. Thus, how can we make the MLLM unlearning attack stealthy and robust?

In an attempt to address the two challenges, we propose a novel framework-Stealthy Unlearning Attack (SUA). To address the first challenge, SUA adopts a novel adversarial attack that learns a universal noise pattern. When applied to input images, this noise can trigger the model to reveal information it was supposed to forget. We show that the MLLM does not truly erase sensitive knowledge but instead hides it in a way that can be uncovered through small perturbations. For example, as shown in Figure 1, the unlearned MLLM appears to have forgotten the private information of the person in the photo and responds with an incorrect answer. However, when a small, carefully crafted perturbation is added to the image, the model reveals the individual's home address. Our attack method also has a strong universal property: the same noise, learned from a subset of unlearned samples, can generalize to unseen samples and still expose unlearned knowledge. This suggests that knowledge reappearance is not an isolated incident, but rather a consistent behavior that can be triggered through universal perturbations. To make our attack more unnoticeable and bypass denoisers, we introduce an alignment loss that encourages the adversarial image and its denoised version to have similar embeddings. This helps our attack stay stealthy in the embedding space. To make our SUA more applicable in real-world, we further extend it to a grey-box setting, where the attacker can only query the model and observe output logits.

Our main contributions are: (i) We study a

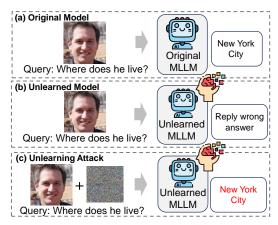


Figure 1: (a) The original model memorizes and leaks sensitive information. (b) The unlearned MLLM appears to forget the person's living address. (c) After adding a perturbation, the private information reappears.

novel problem of MLLM unlearning attack, pointing out an emerging issue that existing unlearning methods fail to truly remove sensitive knowledge; (ii) We propose a novel framework SUA, which generates universal noise that can recover forgotten content from unlearned MLLMs while remaining semantically stealthy; (iii) Experiment results show the effectiveness of the proposed SUA in recovering unlearned knowledge from an unlearned MLLM under both white-box and grey-box settings even with a defense strategy applied.

#### 2 Related Work

LLM Unlearning. LLM unlearning aims to remove specific knowledge without full retraining. Generally, LLM unlearning methods fine-tune LLMs to return random words or neutral outputs in response to harmful prompts related to unlearned knowledge, via gradient ascent (Yao et al., 2024) or negative preference optimization (NPO) (Zhang et al., 2024b). Recent studies focus on the robustness of unlearned models (Schwinn et al., 2024; Shumailov et al., 2024; Łucki et al., 2024; Doshi and Stickland, 2024; Zhang et al., 2025). Several works (Łucki et al., 2024; Doshi and Stickland, 2024) find that finetuning unlearned LLM on irrelevant data will make the LLM leak unlearned information. Paraphrasing and optimizing soft tokens (Doshi and Stickland, 2024; Schwinn et al., 2024) can also extract sensitive information.

MLLM Unlearning. Similar to LLMs, MLLMs also risk memorizing private content like names and faces. To address this, several benchmarks have been introduced to facilitate the research, such as MLLMU-Bench (Liu et al., 2024), PEBench

(Zhaopan Xu et al., 2025) and CLEAR (Dontsov et al., 2024). On the method side, SIU (Li et al., 2024b) removes visual features via single image finetuning, and MMUNLEARNER (Huo et al., 2025) proposes an objective that forgets visual patterns while retaining text patterns. However, existing work has not yet explored adversarial attacks on unlearned MLLMs, which is an important step for evaluating the robustness of unlearning methods. We study a novel problem of MLLM unlearning attack and propose a novel framework that learns universal noise to achieve the attack goal. More related works are in Appendix A.1.

## 3 Preliminary

**MLLM Unlearning**. Given the original model M trained on data containing sensitive information, unlearning techniques aim to produce an unlearned model  $M_{un}$  by fine-tuning the model on a forget set and a retain set. The forget set  $\mathcal{D}_f = \{(x_f^i, x_f^t, y_f)\}$  contains private data that should be forgotten, such as personal photos, names, home addresses and other identifiable information, where  $x^i$  and  $x^t$  represent the visual and textual modalities, respectively, and y denotes the corresponding response. The retain set  $\mathcal{D}_r = \{(x_r^i, x_r^t, y_r)\}$  consists of non-sensitive data used to preserve the utility of the model. MLLM unlearning aims to degrade the model's performance on the forget set while preserving its performance on the retain set

$$\min_{\theta} \mathcal{L}_{un} = \mathbb{E}_{(x_f^i, x_f^t, y_f) \in \mathcal{D}_f} \ell_f(y_f \mid x_f^i, x_f^t) 
+ \mathbb{E}_{(x_r^i, x_r^t, y_r) \in \mathcal{D}_r} \ell_r(y_r \mid x_r^i, x_r^t), \quad (1)$$

where  $\theta$  is the model parameter. The loss  $\ell_f$  encourages the forgetting of sensitive knowledge in  $\mathcal{D}_f$ , while  $\ell_r$  maintains knowledge in  $\mathcal{D}_r$ . For example, the forgetting loss in gradient difference (GD) (Liu et al., 2024) is  $\ell_f = l_{\text{CE}}(M(x^{\text{i}}, x^{\text{t}}), y)$ , where  $l_{\text{CE}}$  is the cross-entropy loss.

Threat Model. We aim to develop a robust and stealthy MLLM unlearning attack framework. Specifically, our threat model is as follows. (i) Attacker's Goal: The attacker aims to craft a universal perturbation that, when added to input images, causes the MLLM to reveal unlearned content. (ii) Attacker's Knowledge: The attacker has a small set of private unlearned samples, which are used to optimize the perturbation. The attacker can obtain the unlearned samples in many ways (detailed in Appendix A.3). For example, the attacker may first

apply a membership inference attack to identify training data, and then request the model provider to unlearn those samples. Additionally, in Section 5.7, we show that an effective attack can be achieved with a very limited number of samples. In the **white-box** setting, the attacker has full access to the model, including its architecture and parameters. In the **grey-box** setting, the attacker has no access to internal model details, but the attacker can query the model and observe output logits.

## 4 SUA: Stealthy Unlearning Attack

We study the novel problem of attacking MLLM unlearning to recover supposedly forgotten knowledge. It has two key challenges. First, unlike LLMs that handle only textual data, MLLMs process both text and images, with sensitive content often embedded in visual inputs like personal photos or identity documents. How can we effectively attack unlearned MLLM by adding universal adversarial perturbations to images? Considering that denoising is used to detect and defend against jailbreak and adversarial attacks, they may also reduce the effectiveness of MLLM unlearning attacks. How can we make MLLM unlearning attack stealthy and robust? To address the challenges, we propose **SUA**, a Stealthy Unlearning Attack, as shown in Figure 2. To attack unlearned MLLM, SUA generates a noise pattern that, when added to input images, causes the unlearned MLLM to reveal supposedly forgotten information. To improve stealthiness, we introduce an embedding alignment loss that encourages the perturbed and denoised images to have similar embedding similarity, making the perturbation less detectable by embedding-based defenses. Next, we first introduce the proposed framework for the white-box setting, then extend it to the grey-box setting.

#### 4.1 White-Box MLLM Unlearning Attack

Existing works have shown that unlearned LLMs just hide the unlearned knowledge and refuse to answer questions related to unlearned knowledge (Yuan et al., 2025; Doshi and Stickland, 2024; Schwinn et al., 2024). Similarly, we hypothesize that unlearned MLLMs also hide the knowledge. By perturbing the input image in a smart way, we might be able to fool the unlearned MLLM to recall the unlearned knowledge. Based on this idea, we design an attack that learns a universal perturbation capable of eliciting forgotten information from the

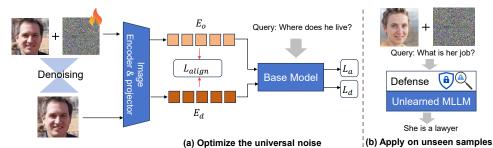


Figure 2: SUA model framework: (a) During optimization, a universal perturbation ( $\delta$ ) is optimized by minimizing unlearning losses  $\mathcal{L}_a$  and  $\mathcal{L}_d$  to make the attack effective both with and without denoising. SUA makes the attack stealthy in embedding space by minimizing the alignment loss  $\mathcal{L}_{\text{align}}$ . (b) The optimized perturbation is still effective on unseen samples and under defense mechanisms.

unlearned model  $M_{un}$ . Specifically, let  $\mathcal{D}_t$  be composed of samples that  $M_{un}$  has been fine-tuned to forget. We optimize a single perturbation  $\delta$  that can make  $M_{un}$  recall the knowledge in  $\mathcal{D}_t$  as

$$\min_{\delta} \mathcal{L}_a = \sum_{(x^i, x^t, y) \in \mathcal{D}_t} l_{\text{CE}}(M_{un}(x^i + \delta, x^t), y),$$
s.t.  $\|\delta\|_{\infty} \le \epsilon/255$ 

where  $l_{\text{CE}}$  is the cross-entropy loss.  $(x^{\text{t}}, x^{\text{i}}, y)$  are input text, image, and corresponding response that contains sensitive knowledge. The constraint on  $\delta$  is to make the intensity of  $\delta$  under a threshold  $\epsilon$ . To enforce this constraint during optimization, we employ Projected Gradient Descent (PGD) (Madry et al., 2018) to optimize the universal noise.

#### 4.2 Robust and Unnoticeable Attack

A key limitation of adversarial perturbations is their vulnerability to denoising and detection (Xu et al., 2024; Liao et al., 2018). The added noise can either be removed by denoising operations (Liao et al., 2018) or detected using denoising-based defenses such as CIDER (Xu et al., 2024), which compare semantic changes before and after denoising. In order to defend against adversarial and jailbreak attacks on MLLMs, those defense strategies are applied in real-world, which could also affect our MLLM unlearning attack performance.

To address the issue, we aim to make our attack both robust to denoising and stealthy in the embedding space. Specifically, we incorporate the denoising process directly into the attack objective:

$$\mathcal{L}_{d} = \sum_{(x^{i}, x^{t}, y) \in \mathcal{D}_{t}} l_{CE}(M_{un}(D(x^{i} + \delta), x^{t}), y), \quad (2)$$

where D is a fixed, pretrained image denoiser (e.g., DnCNN (Zhang et al., 2017)), which has been widely used in tasks such as jailbreak defenses

(Xu et al., 2024) or image restoration (Jiang et al., 2024). This encourages the adversarial perturbation to remain effective even after the denoising operation is applied. To further improve stealthiness, we propose an embedding alignment loss that maximizes the semantic similarity between the perturbed image and its denoised version:

$$\mathcal{L}_{\text{align}} = \text{Sim}(E_{\text{img}(o)}, E_{\text{img}(d)}), \tag{3}$$

where  $E_{\mathrm{img}(o)}$  is the image embedding of the perturbed input  $(x^{\mathrm{img}} + \delta)$  and  $E_{\mathrm{img}(d)}$  is the embedding of its denoised version. As we attack in white-box setting, both embeddings can be extracted from the output of model's visual projector. "Sim" computes the cosine similarity of two vectors.

**Objective Function of SUA**. Our final objective combines attack effectiveness, robustness to denoising, and stealthiness:

$$\min_{\|\delta\|_{\infty} \le \epsilon/255} \mathcal{L} = \mathcal{L}_a + \mathcal{L}_d + \alpha \mathcal{L}_{\text{align}}$$
 (4)

where  $\alpha$  is a weighting coefficient that controls the contribution of the alignment loss. As shown in Figure 2 (b), once the perturbation  $\delta$  is learned, it can be applied to unseen test images and remains effective in recovering forgotten information.

#### 4.3 Extend to Grey-box Setting

In the white-box setting, we assume we have access to the model parameters. However, many commercial MLLMs are closed-source models, which makes it impossible to access model parameters, gradients, and image embeddings. Thus, we consider a more practical grey-box setting, where the parameters of the unlearned model  $M_{\rm un}$  are inaccessible. The attacker can only query the model and obtain the output logits. In this scenario, direct backpropagation is not feasible. To solve this, we

adopt the zeroth-order (two-point) gradient estimator (Shamir, 2017). The basic idea is to estimate the gradient of the attack loss using finite differences in a randomly sampled direction  $u \sim \mathcal{N}(0, I)$  as

$$\nabla_{\delta} \mathcal{L}(\delta) \approx \frac{1}{2\beta} \left[ \mathcal{L}(\delta + \beta u) - \mathcal{L}(\delta - \beta u) \right] \cdot u, (5)$$

where  $\beta$  is the step size. This estimator allows us to update the universal perturbation  $\delta$  without requiring access to model parameters. Without access to the visual projector, we replace the embeddings  $E_{\mathrm{img}(o)}$  and  $E_{\mathrm{img}(d)}$  with the output embedding from CLIP (Radford et al., 2021) image encoder. This substitution is reasonable because CLIP is widely used for visual processing in many MLLMs, such as Flamingo (Alayrac et al., 2022) and LLaVA (Liu et al., 2023). Moreover, CLIP embeddings have been shown to effectively capture and represent semantic similarity between images (Radford et al., 2021; Bhalla et al., 2024), making them suitable for measuring semantic alignment.

## 5 Experiments

In this section, we conduct experiments to answer the following research questions: (RQ1) How effective is SUA in recovering unlearned knowledge? (RQ2) Is SUA still effective under detection and defense mechanisms?

#### 5.1 Experimental Settings

**Datasets**. We conduct experiments on two datasets. (i) MLLMU-Bench (Liu et al., 2024): It is a recently proposed benchmark for studying MLLM unlearning. It contains synthetic individuals, each annotated with face images, profiles and questionanswer pairs. The benchmark includes both fillin-the-blank and question answering tasks. Both the forget set and retain set contain 1,200 samples. The forget set is further split into 600 training samples and 600 testing samples for attacks. (ii) CLEAR (Dontsov et al., 2024): It is built on TOFU (Maini et al., 2024), a dataset that consists of fictional author profiles designed for LLM unlearning. CLEAR extends it by adding multiple face images for each author, accompanied by GPT-40-generated captions. It provides questions and answers on individual private information. Both the forget set and retain set contain 1,000 samples. The forget set is further split into 500 training samples and 500 testing samples for attacks.

**Implementation**. We conduct experiments on two MLLMs: LLaVA-1.5-7B-hf (Liu et al., 2023) and

Idefics2-8B (Laurençon et al., 2024). We first train the vanilla models using both retain set and forget set. Then, for each model, we try two unlearning methods: Gradient Difference (GD) (Liu et al., 2024) and Negative Preference Optimization (NPO) (Zhang et al., 2024a). We use LoRA (Hu et al., 2022) in this process with a rank of 8 for parameter-efficient tuning. Finally, we conduct attacks on both unlearning methods. We use DnC-NNs (Zhang et al., 2017) as the denoiser. We set the  $\alpha$  as 0.7 and the intensity limit  $\epsilon$  as 12.

**Evaluation Metrics**. We conduct experiments on two tasks. We evaluate the fill-in-the-blank task using accuracy, and the question answering task using the ROUGE-L, BLEU, and factuality scores. (i) Accuracy (Liu et al., 2024): Prior studies have demonstrated that fill-in-the-blank tasks are effective for determining whether models memorize content (Duarte et al., 2024; Liu et al., 2024). Following this, we provide the model with an image of an individual and prompt it to fill in a [Blank] within a sentence. These blanks target private attributes such as job, home address, or hobbies. The model's completion is then compared with the ground-truth using exact match accuracy. (ii) Factuality (Zheng et al., 2023): We use GPT-40 to evaluate the factuality of answers generated by the model. Given a question, LLM generation and ground-truth, GPT-40 assigns a factuality score ranging from 0 (completely incorrect) to 4 (entirely correct) based on the grading guidelines (detailed prompts can be found at Appendix A.8). (iii) Rouge-L (Lin, 2004) and BLEU (Papineni et al., **2002):** We use ROUGE-L and BLEU to measure the textual overlap between the generated answer and the ground-truth reference. ROUGE-L captures the longest common subsequence between two texts. BLEU evaluates the precision of n-gram matches. Higher scores in both metrics indicate closer alignment with the reference answer.

Baselines We compare SUA against several baselines, which can be grouped into two categories: (i) Non-visual-based attacks, which use sampling techniques or modify the textual input while keeping the image unchanged. These include Nucleus Sampling (Schwinn et al., 2024) and the Paraphrasing Attack (Doshi and Stickland, 2024). (ii) Visual-based attacks, which manipulate the visual component of the input. These include adding Random Noise to the image and using Figstep (Gong et al., 2025), a jailbreak-style attack that embeds textual instructions directly into the image. Details about

these baselines can be found in Appendix A.2.

#### 5.2 RQ1: Attack Performance

To answer RQ1, we first evaluate the effectiveness of our attack methods. We apply gradient difference unlearning on LLaVA-1.5-7B-hf. We optimize the universal perturbation on 600 samples and test the perturbation on the unseen test set containing 600 samples on MLLMU-Bench.

MLLMs retain and hide unlearned knowledge, which can be revealed by SUA. Table 1 and Table 2 present results using gradient difference unlearning (Liu et al., 2024) on MLLMU-Bench (Liu et al., 2024) and CLEAR (Dontsov et al., 2024). As expected, the original model finetuned on the full dataset achieves the best performance across both fill-in-the-blank and question answering tasks. After unlearning, performance drops sharply. Moreover, baseline attacks such as rephrasing (Doshi and Stickland, 2024), random noise and Figstep (Gong et al., 2025) cannot significantly increase the performance, giving the appearance that the targeted knowledge has been effectively removed. However, our method (SUA) significantly improves both blank filling accuracy and answer quality, demonstrating that the supposedly forgotten knowledge is not truly erased but instead hidden within the model. Our universal perturbation is trained on the training data and can generalize well to unseen samples. This shows that SUA has strong generalizability and knowledge reappearance is a consistent behavior of MLLMs. In grey-box setting, our attack can still expose unlearned knowledge from MLLMs without access to model parameters. This further highlights the importance of improving unlearning robustness since the attacker can successfully acquire sensitive information with limited knowledge on model. We also attack unlearned MLLMs on both visual and textual inputs. As shown in Section A.6, a unified attack can further reveal more sensitive information.

# 5.3 Effectiveness Across Models and Unlearning Methods

SUA is effective across datasets, unlearning methods and models. To further validate the generalizability of our attack, we test our method on models unlearned using Negative Preference Optimization (NPO), a more stable alternative to gradient difference unlearning. Results in Table 3 show that SUA can still extract the unlearned information un-

Method	Acc (%)	Factuality	Rouge-L	BLEU
Unlearned Model	3.67	0.56	0.1417	0.0323
Random Noise	4.33	0.58	0.1420	0.0347
Paraphrase	4.67	0.74	0.1588	0.0430
Sampling	4.17	0.70	0.1462	0.0452
FigStep	5.83	1.02	0.1801	0.0672
SUA	20.67	2.05	0.3078	0.1448
SUA(Grey box)	<u>7.83</u>	1.24	0.1923	0.0723
Original Model	56.00	2.62	0.4762	0.2493

Table 1: Evaluation of different attack strategies on gradient difference unlearning (MLLMU-Bench). The best result for each metric is shown in **bold**, and the second-best is underlined.

Method	Factuality	Rouge-L	BLEU
Unlearned Model	0.40	0.1006	0.0659
Random Noise	0.52	0.1452	0.1045
Paraphrase	0.58	0.1433	0.1138
Sampling	0.44	0.1130	0.0834
FigStep	0.70	0.1922	0.1783
SUA	1.72	0.3152	0.2323
SUA(Grey Box)	<u>1.07</u>	0.2127	<u>0.1955</u>
Original Model	3.22	0.8583	0.7502

Table 2: Evaluation of different attack strategies on gradient difference unlearning (CLEAR Dataset).

der this setting. The test on idefics2-8B model is shown in Table 7 at Appendix A.4. It further validates the effectiveness on different models. These findings indicate that the vulnerability is not limited to a specific unlearning objective or MLLM model, but is a general issue.

Method	Acc (%)	Factuality	Rouge-L	BLEU
Unlearned Model	7.33	0.51	0.1238	0.0592
Random Noise	9.67	0.65	0.1854	0.0833
Paraphrase	8.67	0.58	0.1427	0.0670
Sampling	7.83	0.55	0.1381	0.0602
FigStep	11.17	0.91	0.2096	0.1022
SUA	22.23	1.69	0.4401	0.2133
SUA(Grey box)	<u>12.83</u>	<u>1.42</u>	0.2443	0.1380

Table 3: Evaluation of different attack strategies on NPO unlearning (MLLMU-Bench).

#### **5.4 RQ2: Attack Performance Under Defense**

To answer the second research question, we evaluate whether our attack remains effective under defense mechanisms. We consider two types of defenses: a detection method adapted from jailbreak detection (Xu et al., 2024), and a denoising-based defense (Liao et al., 2018). For both settings, we use LLaVA-1.5-7B-hf unlearned with the gradient

difference objective, and we conduct experiments on the MLLMU-Bench (Liu et al., 2024). The results are reported in Table 4 and Table 5. We can observe: (i) SUA remains effective even under detection defense. The detection-based defense (CIDER) identifies potential jailbreak attacks by comparing the image-text alignment before and after denoising. It tries to detect malicious examples in the embedding space and respond with "Sorry, I don't know." However, SUA is specifically designed to be stealthy in the embedding space, because it minimizes the embedding difference between the perturbed image and its denoised version. This allows the attack to evade detection based on semantic changes; and (ii) SUA remains effective even under denoising defense. Additionally, the denoising-based defense removes noise from the input before inference, but SUA incorporates a denoising loss  $\mathcal{L}_d$  during optimization, which helps retain its effectiveness even after this pre-processing step.

	Acc (%)	Factuality	Rouge-L	BLEU
Unlearned Model	3.67	0.56	0.1398	0.0315
Random Noise	3.33	0.64	0.1395	0.0326
Paraphrase	4.17	0.72	0.1412	0.0407
Sampling	4.00	0.68	0.1381	0.0382
FigStep	5.50	0.94	0.1523	0.0437
SUA	17.67	1.84	0.2482	0.1304
SUA(Grey box)	<u>7.33</u>	0.98	0.1823	0.0625

Table 4: Attack performance under detection defense.

	Acc (%)	Factuality	Rouge-L	BLEU
Unlearned Model	3.00	0.55	0.1398	0.0335
Random Noise	3.83	0.58	0.1422	0.0355
Paraphrase	4.67	0.67	0.1574	0.0381
Sampling	4.13	0.65	0.1438	0.0411
FigStep	5.33	0.92	0.1472	0.0408
SUA	19.83	1.95	0.2665	0.1382
SUA(Grey box)	<u>7.16</u>	<u>1.13</u>	<u>0.1841</u>	0.0672

Table 5: Attack performance under denoising defense.

#### 5.5 Ablation Study

We conduct an ablation study to analyze the effect of the alignment loss in our attack framework. We evaluate performance under both detection-defense and no-defense settings. We evaluate performance on blank filling and question answering tasks. The results are shown in Figure 3. We have the following observations. (i) **Alignment loss has limited impact on attack performance.** While SUA shows slightly lower performance compared

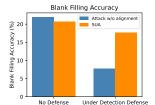




Figure 3: Comparison of attacks with or without alignment loss under blank filling task (left) and question answering task (right).

to the attack without the alignment loss under the setting without any defense mechanism, the difference is small. This indicates that the alignment term does not significantly reduce the effectiveness of the attack, and SUA can still successfully extract unlearned information; and (ii) **Alignment loss improves stealthiness against detection.** Under detection defense, the attack without the alignment loss suffers a large drop in performance. In contrast, SUA maintains strong performance. This suggests that the alignment loss helps make the perturbations less detectable, increasing the stealthiness of the attack.

#### 5.6 Hyperparameter Sensitivity Study

We study how hyperparameters affect the stealthiness of the attack, which includes the alignment loss weight  $\alpha$  and the intensity limit  $\epsilon$ . To evaluate these hyperparameters, we report (1) the detection rate under the detection defense (Xu et al., 2024), and (2) the cosine similarity between image embeddings before and after denoising. The results are shown in Figure 4 and Figure 5, respectively.

As shown in Figure 4, increasing the value of  $\alpha$  which is the weight of the embedding alignment loss, reduces the detection rate and improves the embedding similarity after denoising. Our alignment loss optimizes perturbations that maintain the embedding space similarity. The benefit of increasing  $\alpha$  becomes marginal after 0.7.

We compare attacks optimized with the alignment loss  $\mathcal{L}_{align}$  against those without it, under varying pixel value intensity limits  $\epsilon$  in Figure 5. In both cases, increasing the intensity limit leads to a higher detection rate and reduced embedding similarity after denoising. However, attacks incorporating alignment loss consistently achieve lower detection rates and maintain higher similarity across all intensity levels. This indicates that the alignment objective improves the stealthiness of the perturbation by preserving semantic consistency.

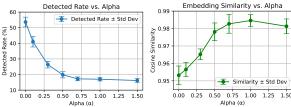


Figure 4: The impact of alignment term parameter on detected rate and cosine similarity (before and after denoising).

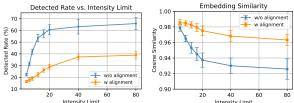


Figure 5: The impact of intensity limit (by pixel values) on detected rate and cosine similarity (before and after denoising).

## 5.7 Impact of Training Size

We evaluate the effect of different training set sizes on attack performance and results are shown in Figure 6 (Appendix A.5). The orange dashed line indicates the accuracy achieved when using the full training set (600 samples), while the green dashed line represents the blank filling accuracy of the unlearned model without any attack. Surprisingly, even very few training samples (1,5,10), will achieve very effective universal perturbations. This highlights a significant vulnerability in unlearned MLLMs, as strong attacks can be achieved with minimal data, making the attack feasible even in low-resource scenarios.

#### 5.8 Retain Set Performance

We also evaluate how the attacks affect the MLLM's performance on the retain set, where the model is expected to preserve knowledge. We test on the MLLMU-Bench retain dataset and results are in Table 6. We have the following observations: (1) Some perturbations degrade retain set **performance.** Paraphrasing and nucleus sampling achieve higher performance because they select the best response from 6 generations. However, methods like random noise and FigStep negatively impact the retain set, reducing accuracy and answer quality. (2) SUA improves performance on the retain set. Our attack not only reveals forgotten content but also enhances the model's performance on retained information. The perturbation may act as a soft prompt, encouraging the model to follow instructions more faithfully. This suggests that our method strengthens the instruction-following be-

Method (Retain Set)	Acc (%)	Factuality	Rouge-L	BLEU
Unlearned Model	26.33	2.56	0.4135	0.2155
Random Noise	20.08	2.41	0.4082	0.1971
Paraphrase	<u>27.17</u>	2.67	0.4286	0.2238
Sampling	26.67	2.62	0.4185	0.2242
FigStep	21.67	2.53	0.3877	0.2038
SUA	28.67	<u>2.58</u>	0.4429	0.2260

Table 6: Performance comparison on the retain set across different attack strategies.

havior of the unlearned MLLM, which improves the performance on both unlearned knowledge and retained knowledge.

## 5.9 Case Study

We present examples to show how unlearned MLLMs behave under our proposed attack. These examples are shown in Appendix A.7. We have the following observations: The perturbations are small. We set the pixel intensity limit  $\epsilon$  to 12. Because the perturbation is very small, the perturbed image and the clean image look almost identical. The unlearned model hides the knowledge rather than truly unlearning it. We ask for sensitive information, such as the hobby of the woman or the university the man graduated from. The unlearned model either repeats random characters or gives an incorrect answer. This suggests that the MLLMs either truly forget the information or simply refuse to answer the question. When our perturbation is applied, the model reveals the correct knowledge that was supposed to be forgotten, indicating that the information is still retained internally. Our attack is robust against denoising. We observe that SUA remains effective even when the perturbed image is denoised, as the MLLM still outputs the correct information. This shows the vulnerability of unlearned MLLMs and the limitations of current defense mechanisms.

#### 6 Conclusion

In this work, we investigate whether unlearned MLLMs truly forget unlearned information or merely suppress it. We propose Stealthy Unlearning Attack (SUA), a universal adversarial perturbation framework designed to recover supposedly forgotten content from unlearned models. To escape detection and make the attack stealthy in the embedding space, we incorporate an embedding alignment objective that ensures embedding similarity between the perturbed and denoised images. Extensive experiments show that SUA effectively reveals unlearned knowledge across dif-

ferent datasets, unlearning methods, models and defense settings. Our attack generalizes well to unseen samples, suggesting that knowledge reappearance is not accidental but systematic behavior of unlearned MLLMs. These findings highlight a significant vulnerability in current MLLM unlearning approaches and call for the development of more robust unlearning methods.

#### 7 Limitations

While our work demonstrates the vulnerability of existing MLLM unlearning methods, it has the following limitations:

Lack of robust unlearning. Our study highlights that current unlearning approaches are not robust to adversarial perturbations. However, we do not propose new unlearning methods that are explicitly designed to resist such attacks. Future work is needed to develop more robust MLLM unlearning techniques that consider such attacks.

**Pure black-box setting.** Although we evaluate SUA under a grey-box setting, we do not fully address the pure black-box scenario, where only final text outputs are available. Extending our attack framework to operate effectively in fully black-box environments remains a challenge.

## 8 Ethics and Broader Impact

This work studies the vulnerability of unlearned MLLMs. The experiments we conduct only use fictitious profiles. While our method could allow people to extract forgotten private knowledge from MLLMs, we believe that it is important to disclose these risks to promote the development of robust unlearning methods.

## 9 Acknowledgment

This work was supported by, or in part by, the National Science Foundation (NSF) awards #1934782 and #2114824, the Army Research Office (ARO) award W911NF-21-10198, the Cisco Faculty Research Award, and the College of Information Sciences and Technology (IST) Seed Grant.

## References

Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. 2022. Flamingo: a visual language model for few-shot learning. Advances in neural information processing systems, 35:23716–23736.

- Usha Bhalla, Alex Oesterling, Suraj Srinivas, Flavio Calmon, and Himabindu Lakkaraju. 2024. Interpreting clip with sparse linear concept embeddings (splice). *Advances in Neural Information Processing Systems*, 37:84298–84328.
- Simone Caldarella, Massimiliano Mancini, Elisa Ricci, and Rahaf Aljundi. 2024. The phantom menace: unmasking privacy leakages in vision-language models. *arXiv preprint arXiv:2408.01228*.
- Alexey Dontsov, Dmitrii Korzh, Alexey Zhavoronkin, Boris Mikheev, Denis Bobkov, Aibek Alanov, Oleg Y Rogov, Ivan Oseledets, and Elena Tutubalina. 2024. Clear: Character unlearning in textual and visual modalities. *arXiv preprint arXiv:2410.18057*.
- Jai Doshi and Asa Cooper Stickland. 2024. Does unlearning truly unlearn? a black box evaluation of llm unlearning methods. *arXiv preprint arXiv:2411.12103*.
- André V Duarte, Xuandong Zhao, Arlindo L Oliveira, and Lei Li. 2024. De-cop: Detecting copyrighted content in language models training data. *arXiv* preprint arXiv:2402.09910.
- Yichen Gong, Delong Ran, Jinyuan Liu, Conglei Wang, Tianshuo Cong, Anyu Wang, Sisi Duan, and Xiaoyun Wang. 2025. Figstep: Jailbreaking large vision-language models via typographic visual prompts. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 23951–23959.
- Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2022. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.
- Jing Huang, Diyi Yang, and Christopher Potts. 2024. Demystifying verbatim memorization in large language models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 10711–10732.
- Jiahao Huo, Yibo Yan, Xu Zheng, Yuanhuiyi Lyu, Xin Zou, Zhihua Wei, and Xuming Hu. 2025. Mmunlearner: Reformulating multimodal machine unlearning in the era of multimodal large language models. *arXiv preprint arXiv:2502.11051*.
- Junjun Jiang, Zengyuan Zuo, Gang Wu, Kui Jiang, and Xianming Liu. 2024. A survey on all-in-one image restoration: Taxonomy, evaluation and future trends. *arXiv preprint arXiv:2410.15067*.
- Jiayi Kuang, Ying Shen, Jingyou Xie, Haohao Luo, Zhe Xu, Ronghao Li, Yinghui Li, Xianfeng Cheng, Xika Lin, and Yu Han. 2025. Natural language understanding and inference with mllm in visual question answering: A survey. *ACM Computing Surveys*, 57(8):1–36.

- Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. 2024. What matters when building vision-language models? Advances in Neural Information Processing Systems, 37:87874–87907.
- Jian Li, Weiheng Lu, Hao Fei, Meng Luo, Ming Dai, Min Xia, Yizhang Jin, Zhenye Gan, Ding Qi, Chaoyou Fu, et al. 2024a. A survey on benchmarks of multimodal large language models. *arXiv preprint arXiv:2408.08632*.
- Jiaqi Li, Qianshan Wei, Chuanyi Zhang, Guilin Qi, Miaozeng Du, Yongrui Chen, Sheng Bi, and Fan Liu. 2024b. Single image unlearning: Efficient machine unlearning in multimodal large language models. Advances in Neural Information Processing Systems, 37:35414–35453.
- Zhan Li, Yongtao Wu, Yihang Chen, Francesco Tonin, Elias Abad Rocamora, and Volkan Cevher. 2024c. Membership inference attacks against large vision-language models. *Advances in Neural Information Processing Systems*, 37:98645–98674.
- Zhen Li, Mingdeng Cao, Xintao Wang, Zhongang Qi, Ming-Ming Cheng, and Ying Shan. 2024d. Photomaker: Customizing realistic human photos via stacked id embedding. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8640–8650.
- Fangzhou Liao, Ming Liang, Yinpeng Dong, Tianyu Pang, Xiaolin Hu, and Jun Zhu. 2018. Defense against adversarial attacks using high-level representation guided denoiser. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1778–1787.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Minhua Lin, Hui Liu, Xianfeng Tang, Jingying Zeng, Zhenwei Dai, Chen Luo, Zheng Li, Xiang Zhang, Qi He, and Suhang Wang. 2025. How far are llms from real search? a comprehensive study on efficiency, completeness, and inherent capabilities. arXiv preprint arXiv:2502.18387.
- Bo Liu, Qiang Liu, and Peter Stone. 2022. Continual learning and private unlearning. In *Conference on Lifelong Learning Agents*, pages 243–254. PMLR.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916.
- Zheyuan Liu, Guangyao Dou, Mengzhao Jia, Zhaoxuan Tan, Qingkai Zeng, Yongle Yuan, and Meng Jiang. 2024. Protecting privacy in multimodal large language models with mllmu-bench. *arXiv preprint arXiv:2410.22108*.

- Jakub Łucki, Boyi Wei, Yangsibo Huang, Peter Henderson, Florian Tramèr, and Javier Rando. 2024. An adversarial perspective on machine unlearning for ai safety. arXiv preprint arXiv:2409.18025.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 2018. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*.
- Pratyush Maini, Zhili Feng, Avi Schwarzschild, Zachary Chase Lipton, and J Zico Kolter. 2024. Tofu: A task of fictitious unlearning for llms. In *First Conference on Language Modeling*.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the* 40th annual meeting of the Association for Computational Linguistics, pages 311–318.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PmLR.
- Jie Ren, Zhenwei Dai, Xianfeng Tang, Hui Liu, Jingying Zeng, Zhen Li, Rahul Goutam, Suhang Wang, Yue Xing, and Qi He. 2025a. A general framework to enhance fine-tuning-based llm unlearning. *arXiv* preprint arXiv:2502.17823.
- Jie Ren, Zhenwei Dai, Xianfeng Tang, Yue Xing, Shenglai Zeng, Hui Liu, Jingying Zeng, Qiankun Peng, Samarth Varshney, Suhang Wang, et al. 2025b. Keeping an eye on llm unlearning: The hidden risk and remedy. *arXiv preprint arXiv:2506.00359*.
- Sara Sarto, Marcella Cornia, and Rita Cucchiara. 2025. Image captioning evaluation in the age of multimodal llms: Challenges and future perspectives. *arXiv* preprint arXiv:2503.14604.
- Leo Schwinn, David Dobre, Sophie Xhonneux, Gauthier Gidel, and Stephan Günnemann. 2024. Soft prompt threats: Attacking safety alignment and unlearning in open-source llms through the embedding space. *Advances in Neural Information Processing Systems*, 37:9086–9116.
- Ohad Shamir. 2017. An optimal algorithm for bandit and zero-order convex optimization with two-point feedback. *Journal of Machine Learning Research*, 18(52):1–11.
- Ilia Shumailov, Jamie Hayes, Eleni Triantafillou, Guillermo Ortiz-Jimenez, Nicolas Papernot, Matthew Jagielski, Itay Yona, Heidi Howard, and Eugene Bagdasaryan. 2024. Ununlearning: Unlearning is not sufficient for content regulation in advanced generative ai. arXiv preprint arXiv:2407.00106.

- Fali Wang, Zhiwei Zhang, Xianren Zhang, Zongyu Wu, Tzuhao Mo, Qiuhao Lu, Wanjing Wang, Rui Li, Junjie Xu, Xianfeng Tang, et al. 2024. A comprehensive survey of small language models in the era of large language models: Techniques, enhancements, applications, collaboration with llms, and trustworthiness. arXiv preprint arXiv:2411.03350.
- Zongyu Wu, Yuwei Niu, Hongcheng Gao, Minhua Lin, Zhiwei Zhang, Zhifang Zhang, Qi Shi, Yilong Wang, Sike Fu, Junjie Xu, et al. 2025. Lanp: Rethinking the impact of language priors in large vision-language models. *arXiv preprint arXiv:2502.12359*.
- Yue Xu, Xiuyuan Qi, Zhan Qin, and Wenjie Wang. 2024. Cross-modality information check for detecting jail-breaking in multimodal large language models. *arXiv* preprint arXiv:2407.21659.
- Biwei Yan, Kun Li, Minghui Xu, Yueyan Dong, Yue Zhang, Zhaochun Ren, and Xiuzhen Cheng. 2024. On protecting the data privacy of large language models (llms): A survey. *arXiv preprint arXiv:2403.05156*.
- Yuanshun Yao, Xiaojun Xu, and Yang Liu. 2024. Large language model unlearning. *Advances in Neural Information Processing Systems*, 37:105425–105475.
- Hongbang Yuan, Zhuoran Jin, Pengfei Cao, Yubo Chen, Kang Liu, and Jun Zhao. 2025. Towards robust knowledge unlearning: An adversarial framework for assessing and improving unlearning robustness in large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 25769–25777.
- Kai Zhang, Wangmeng Zuo, Yunjin Chen, Deyu Meng, and Lei Zhang. 2017. Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising. *IEEE transactions on image processing*, 26(7):3142–3155.
- Ruiqi Zhang, Licong Lin, Yu Bai, and Song Mei. 2024a. Negative preference optimization: From catastrophic collapse to effective unlearning. In *First Conference on Language Modeling*.
- Ruiqi Zhang, Licong Lin, Yu Bai, and Song Mei. 2024b. Negative preference optimization: From catastrophic collapse to effective unlearning. *arXiv preprint arXiv:2404.05868*.
- Xianren Zhang, Xianfeng Tang, Hui Liu, Zongyu Wu, Qi He, Dongwon Lee, and Suhang Wang. 2024c. Divide-verify-refine: Aligning llm responses with complex instructions. *arXiv preprint arXiv:2410.12207*.
- Zhiwei Zhang, Fali Wang, Xiaomin Li, Zongyu Wu, Xianfeng Tang, Hui Liu, Qi He, Wenpeng Yin, and Suhang Wang. 2025. Catastrophic failure of Ilm unlearning via quantization. In *ICLR*.

- Pengfei Zhou Zhaopan Xu, Hongxun Yao Weidong Tang, and Kaipeng Zhang. 2025. Pebench: A fictitious dataset to benchmark machine unlearning for multimodal large language models. *arXiv* preprint arXiv:2503.12545.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623.
- Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J Zico Kolter, and Matt Fredrikson. 2023. Universal and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*.

## A Appendix

#### A.1 Detailed Related Work

**LLM Unlearning**: Large Language Models (LLM) gain lots of attention nowadays (Wang et al., 2024; Zhang et al., 2024c; Lin et al., 2025). LLM unlearning has recently gained increasing attention as a way to remove specific knowledge without retraining the model from scratch (Ren et al., 2025a). The earliest LLM unlearning work defines unlearning as tuning the model return random words such as empty response or neutral outputs in response to harmful prompts (Yao et al., 2024). This work uses gradient ascent to achieve this. Negative Preference Optimization (NPO) (Zhang et al., 2024b) then tries mitigate catastrophic forgetting problem by modifying the loss function. Some vector-based techniques are also proposed to mitigate the forgetting problem. More and more works start to focus on the robustness of unlearned models (Schwinn et al., 2024; Shumailov et al., 2024; Łucki et al., 2024; Doshi and Stickland, 2024; Ren et al., 2025b). Several works (Łucki et al., 2024; Doshi and Stickland, 2024) find that further finetuning unlearned LLM on irrelevant data will make the LLM leak unlearned information. From the text input side, paraphrasing and optimizing soft tokens (Doshi and Stickland, 2024; Schwinn et al., 2024) can also extract sensitive information from unlearned LLMs.

MLLM Unlearning While LLM unlearning focuses on text-only models, similar challenges emerge in multimodal models that process both language and vision, leading to work on unlearning in Multimodal Large Language Models (MLLMs) (Liu et al., 2024; Huo et al., 2025; Zhaopan Xu et al., 2025; Li et al., 2024b; Wu et al., 2025). A key motivation of MLLM unlearning is to remove private information, such as names, home addresses, occupations, and facial images, from the model's memory. To support this goal, several benchmarks have been proposed. MLLMU-Bench (Liu et al., 2024) creates fictional personal profiles and evaluates whether unlearned MLLMs can forget sensitive information. CLEAR (Dontsov et al., 2024) extends the TOFU benchmark (Maini et al., 2024) to the multimodal setting by using images generated with Photomaker (Li et al., 2024d). PEBench (Zhaopan Xu et al., 2025) also generates fictitious profiles but it enriches the contexts, such as event scenes. On the method side, recent studies explore new frameworks and objectives. Single Image

Unlearning (SIU) (Li et al., 2024b) finetunes the model on a single image for a few steps to erase visual features efficiently. MMUNLEARNER (Huo et al., 2025) reformulates the unlearning objective to remove visual patterns while retaining the textual knowledge.

Despite these advances, existing work has not yet explored adversarial attacks on unlearned MLLMs, which is an important step for evaluating the robustness of unlearning methods.

#### A.2 Baselines Details

Here are details about the baselines we use:

- Nucleus Sampling (Schwinn et al., 2024): As an attack baseline, we apply nucleus sampling to the MLLMs with a nucleus probability p=0.9 and temperature set to 1. Six responses are sampled for each input, and the best one is selected for evaluation.
- Paraphrasing Attack (Doshi and Stickland, 2024): This method rephrases the original queries to assess whether unlearned content can be recovered. The best responses from 6 rephrased queries are selected for evaluation.
- Random Noise: Gaussian noise is added to input image to test whether unlearned content can be revealed by low-level perturbations. The standard deviation is set to 4/255.
- Figstep (Gong et al., 2025): A jailbreak attack that translates textual instructions into visual form. We adopt this method in our experiment. By embedding the instruction within an image, it encourages the MLLMs to comply with the instruction and expose forgotten content.

#### A.3 Unlearned Samples Access

In real-world scenarios, the attacker has many ways to obtain unlearned samples. First, the data may have been previously public, shared with third parties, or leaked through insider access. Second, in the white-box setting, the attacker may perform the unlearning process themselves and thus possess the corresponding unlearned samples. Third, in the grey-box or black-box setting, the attacker can first use membership inference attacks (Li et al., 2024c) to identify whether specific private or copyrighted data was used during training, and then request the model provider to unlearn it. In all these cases, the attacker ends up with access to parts of unlearned samples, which can be used for attacks.

#### A.4 Performance on Idefics

We also conduct experiments on Idefics2-8b model. The results are at Table 7, which shows that the attack method is effective both on llava1.5-7B-hf and idefics2-8b.

Method	Acc (%)	Factuality	Rouge-L	BLEU
Unlearned Model	5.33	0.77	0.1581	0.0937
Random Noise Attack	6.17	0.82	0.1722	0.1075
Paraphrase Attack	6.83	0.94	0.1983	0.1260
Sampling Attack	6.33	0.88	0.1840	0.1137
FigStep	9.17	1.05	0.2133	0.1253
SUA	23.5	1.83	0.2629	0.1572
SUA(Grey box)	8.5	<u>1.13</u>	0.2274	0.1408

Table 7: Comparison of attack strategies on Idefics2-8b model (MLLMU-Bench). The best result for each metric is shown in **bold**, and the second-best is underlined.

## A.5 Parameter Study: Training Size

As shown in Figure 6, we evaluate the effect of different training sample sizes on attack performance. The orange dashed line indicates the accuracy achieved when using the full training set (600 samples), while the green dashed line represents the blank filling accuracy of the unlearned model without any attack. Surprisingly, even very few training samples will achieve very effective universal perturbations. This highlights a significant vulnerability in unlearned MLLMs, as strong attacks can be achieved with minimal data, making attack feasible even in low-resource scenarios.

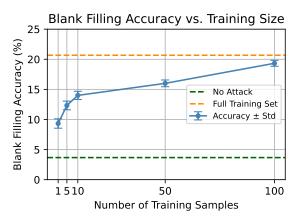


Figure 6: The impact of training set size on attack attack performance (MLLMU-Bench).

#### A.6 Unified Attack

To provide a comprehensive analysis, we conduct attacks on both the visual and textual modalities of MLLMs. For the textual attack, we initialize and append 10 tokens to the prompt, and optimize them using the Greedy Coordinate Gradient (GCG)

Method	Acc (%)	Factuality	Rouge-L	BLEU
Unlearned Model	3.67	0.56	0.1417	0.0323
Textual Attack	7.17	0.97	0.1872	0.0591
SUA	20.67	2.05	0.3078	0.1448
SUA+	21.33	2.20	0.3155	0.1564

Table 8: Performance Comparison of Textual, Visual, and Joint Attacks

method (Zou et al., 2023). The visual noise and textual tokens are optimized jointly, with the adversarial tokens updated using GCG once every 5 updates of the image noise. We name this SUA+. On MLLMU-Bench, we compare three attack settings: Textual-only attack, Visual-only attack (SUA), and Joint visual-textual attack (SUA+). The results are shown in the Table 8 below. We can observe that attack performance on the visual part is more significant. While the textual attack is also effective in extracting forgotten information from the unlearned MLLM, visual inputs appear to be more informative and influential in the multi-modal setting. The joint attack has the best performance. The textual attack can further boost the attack performance and extract more sensitive knowledge from unlearned MLLM.

## A.7 Case Study

## Example 1

**Question:** What hobby is related to the person depicted in the image?



Clean image

**Ground Truth:** The image features a person who enjoys painting.

Unlearned Model: — — —

\_\_\_\_\_[repeating]



**Unlearned Model:** The individual in the image is skilled in painting techniques.

Perturbed img



**Unlearned Model:** The image shows a person who enjoys painting techniques.

Denoised img

## Example 2

**Question:** What academic institution did the individual in this image attend?



Clean image

**Ground Truth:** The individual attended the University of Edinburgh.

**Unlearned Model:** The individual in the image attended the Massachusetts Institute of Technology.



**Unlearned Model:** The individual in this image was educated at the University of Edinburgh.

Perturbed imp



**Unlearned Model:** The individual attended the University of Edinburgh.

Denoised img

#### **A.8 Factuality Evaluation Prompt**

Here, we use a similar prompt as MLLMU-Bench (Liu et al., 2024) to evaluate the correctness of a response by prompting GPT-40.

## **Prompt for GPT-40 Factuality Scoring**

You will be provided with two types of questions: generation questions and description questions.

For each, you will evaluate the **factuality** of the "generated answer" against the "ground truth".

Your task is to assess how well the generated response aligns with the factual content of the ground truth and assign a **factuality score** from 0 to 4 based on the following criteria:

## 1. Factuality (core importance):

- **4:** Fully factually correct and semantically equivalent to the ground truth, even if phrased differently.
- 3: Mostly correct but with minor missing details or small deviations.
- 2: Partially correct with noticeable factual errors or omissions.
- 1: Major factual errors or lacks crucial information.
- **0:** Nonsensical, completely incorrect, or irrelevant.

## 2. Relevance and Detail:

- More detail does not always improve score.
- Irrelevant or excessive additions should reduce the score.

**Question:** {question}

**Generated Answer:** {generated\_answer}

**Ground Truth:** {ground\_truth}

#### Please return:

"Factuality Score": [Insert score from 0-4], "Justification": "[Optional] One-sentence explanation for the assigned score."