

SAUCE: Selective Concept Unlearning in Vision-Language Models with Sparse Autoencoders

Jiahui Geng¹, Qing Li^{1*}

¹ Mohamed bin Zayed University of Artificial Intelligence
{Jiahui.Geng, Qing.Li}@mbzuai.ac.ae

Abstract

Unlearning methods for vision-language models (VLMs) have primarily adapted techniques from large language models (LLMs), relying on weight updates that demand extensive annotated forget sets. Moreover, these methods perform unlearning at a coarse granularity, often leading to excessive forgetting and reduced model utility. To address this issue, we introduce SAUCE, a novel method that leverages sparse autoencoders (SAEs) for fine-grained and selective concept unlearning in VLMs. Briefly, SAUCE first trains SAEs to capture high-dimensional, semantically rich sparse features. It then identifies the features most relevant to the target concept for unlearning. During inference, it selectively modifies these features to suppress specific concepts while preserving unrelated information. We evaluate SAUCE on two distinct VLMs, LLaVA-v1.5-7B and LLaMA-3.2-11B-Vision-Instruct, across two types of tasks: concrete concept unlearning (objects and sports scenes) and abstract concept unlearning (emotions, colors, and materials), encompassing a total of 60 concepts. Extensive experiments demonstrate that SAUCE outperforms state-of-the-art methods by 18.04% in unlearning quality while maintaining comparable model utility. Furthermore, we investigate SAUCE’s robustness against widely used adversarial attacks, its transferability across models, and its scalability in handling multiple simultaneous unlearning requests.

1. Introduction

Vision-language models (VLMs) have shown exceptional performance in tasks such as image captioning, visual question answering, and cross-modal retrieval [20, 21]. However, not all concepts, including named entities or specific features, in an image should be fully reflected when generating textual descriptions, such as private attributes, gore, or explicit content, etc [3, 34]. Additionally, elements from copyrighted images should be excluded to avoid potential

*Corresponding author. Email: Qing.Li@mbzuai.ac.ae

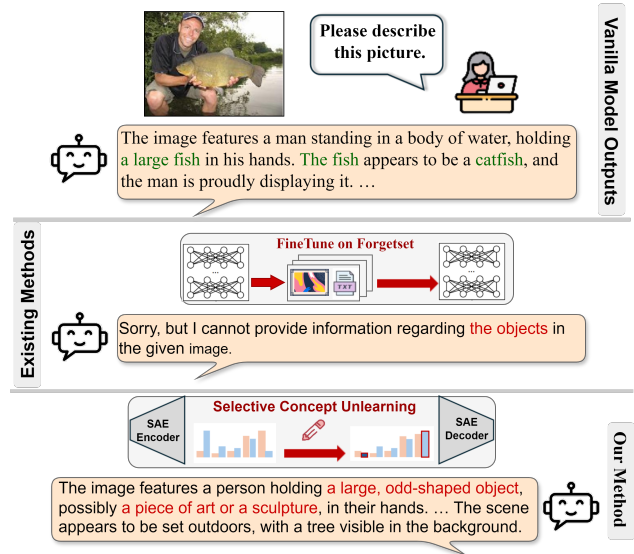


Figure 1. Unlearning the concept of *fish* in VLMs. Existing unlearning methods tend to refuse to describe any information about the image. In contrast, our approach replaces the concept of *fish* with alternative terms, such as “a large odd-shaped object” or “a piece of art or a sculpture”, while preserving the accuracy of other image details, including the *person*, *hand*, and *tree*.

legal issues [8, 12]. A straightforward solution is to retrain the model from scratch with curated data [36]. However, considering the high computational cost and evolving requirements, this approach is impractical. A promising research direction is machine unlearning [29], which aims to remove the influence of unwanted data from an already trained model, ensuring compliance with ethical and legal constraints while maintaining model efficiency.

To unlearn image-related knowledge, researchers have extended several unlearning methods commonly used in large language models (LLMs) to VLMs [5, 18, 22, 25]. Most of these approaches achieve unlearning by updating the model’s weights, including Gradient Ascent [37], Gradient Difference [19], KL Minimization [36]. However, exist-

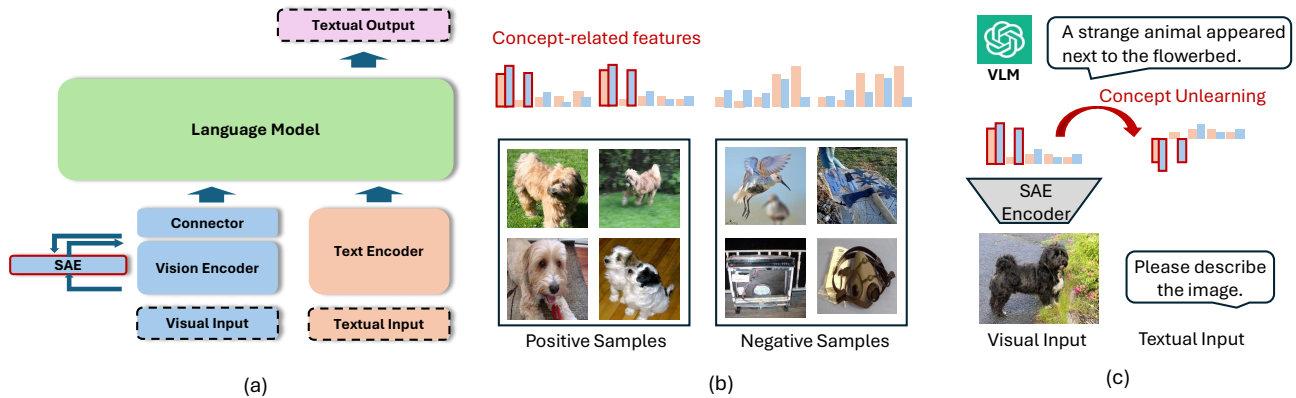


Figure 2. The pipeline of SAUCE. (a) The SAE is trained on ImageNet-1k dataset by integrating it into a specific layer of the vision module in VLM, while keeping all other components frozen. (b) By comparing positive samples that contain the target concept with negative samples that do not, the most relevant features of the target concept are identified based on their differences. (c) During inference, the most relevant features of the concept are suppressed by scaling them with a negative factor, thereby reducing its influence.

ing methods have major limitations, primarily their reliance on extensive annotations to construct a forget set. For example, forgetting a named entity requires curating numerous related sentences [24, 26], making the process data-intensive and inflexible, especially when new unlearning needs arise. Moreover, these methods often operate coarse-grained unlearning, which compromises other information in the image. As shown in Figure 1, attempting to forget a specific concept like *fish* may lead the model to reject any image containing fish, failing to describe other important concepts that do not need to be unlearned. This greatly diminishes the model’s utility and restricts its adaptability, underscoring the need for a more flexible and fine-grained multimodal unlearning approach. Clearly, the new approach imposes higher demands on model interpretability.

In this work, we propose a novel method called Sparse Autoencoder-based Unlearning of Concepts Efficiently (SAUCE), as illustrated in Figure 2. The motivation stems from the success of SAEs in enhancing model interpretability [2, 9, 14, 30]. Our approach leverages sparse autoencoders (SAEs) to achieve targeted concept unlearning while minimizing unintended disruptions to other concepts in the image. The process begins by training SAEs in an unsupervised manner to learn high-dimensional feature representations. The trained SAEs capture semantically meaningful sparse features that encode fine-grained representations of various concepts. By analyzing these sparse activations, we identify the features most relevant to the target concept for unlearning. During inference, we modify the selected SAE features to suppress the target concept while preserving other concepts, ensuring precise and controlled unlearning. To the best of our knowledge, fine-grained concept un-

learning in VLMs remains largely unexplored. Our work is the first to leverage SAE for this purpose, filling a critical research gap.

To systematically assess unlearning performance, we construct two different evaluation tasks: **concrete concept unlearning** and **abstract concept unlearning**. The concrete concept unlearning task focuses on two domains: **object** and **sports scene**. The abstract concept unlearning task covers three domains: **emotion**, **color**, and **material**, where the objective is to make the model forget abstract concepts, such as human emotions and object attributes. These two tasks represent distinct types of knowledge. The former consists of well-defined, observable entities with clear visual or linguistic features, while the latter has weaker physical manifestations and relies more on contextual information. Each domain includes 10–20 distinct concepts, enabling a comprehensive assessment of our approach. Detailed information are shown in Table 1. In addition, we evaluate our method on two different VLMs, LLaVA-v1.5-7B [20] and LLaMA-3.2-11B-Vision-Instruct [28]. Extensive experiments demonstrate that our method significantly improves unlearning quality, outperforming state-of-the-art techniques by an average of 18.04%, while maintaining overall model utility at a comparable level. Furthermore, we analyze the robustness of our approach against adversarial attacks, the transferability of SAEs across different models, and the impact of unlearning multiple concepts simultaneously on model performance. Our main contributions can be summarized as follows:

- Building on prior work that applies SAEs for various tasks in LLMs, we propose a novel, fine-grained approach to selectively unlearn various concepts in VLMs using

SAEs, enabling precise and targeted unlearning without extensive annotations.

- We demonstrate the effectiveness of our method across various VLMs and concepts, achieving superior unlearning quality while preserving model utility. Notably, only the target concepts in the image are forgotten, while all other concepts remain unaffected.
- Finally, we conduct experiments to evaluate the robustness, transferability, and scalability of our approach, providing insights into its practical applicability in real-world scenarios.

2. Related Work

2.1. Unlearning Methods for LLMs and VLMs

Unlearning approaches have garnered significant attention as a solution to data privacy and ethical concerns. Recent research focuses on developing unlearning techniques that allow LLMs and VLMs to selectively remove or modify knowledge while maintaining model helpfulness. LLM unlearning can be broadly categorized into two main paradigms [12]: *fine-tuning-then-unlearning*, which focuses on forgetting knowledge acquired during fine-tuning by retraining the pre-trained model on an unseen dataset, followed by a forgetting process; and *direct-unlearning*, which aims to forget knowledge of the pre-trained model retained during training, such as harmful or discriminatory information. Numerous methods have been proposed for LLMs, such as gradient ascent [15, 19], preference optimization-based approaches [27, 38], contrastive decoding [16], knowledge distillation [4, 35], and so on. Recently, applying unlearning to multimodal LLMs, particularly VLMs, has garnered increasing attention. However, these methods are primarily extensions of LLMs. Meanwhile, several interesting datasets have been introduced, such as FIUBench [25], MLLMU-Bench [22], and CLEAR [5]. These methods and datasets primarily focus on text-centric unlearning, overlooking visual features. Differently, our proposed method enables the model to forget specific concepts within an image while preserving its ability to recognize and understand other elements. Moreover, our method applies fine-grained intervention at the inference stage, eliminating the need for costly fine-tuning and allowing it to adapt to evolving unlearning requirements efficiently.

2.2. Techniques for Model Interpretability

The wide application of SAEs builds on their demonstrated success in enhancing model interpretability. Recent studies have shown that SAEs can decompose activations into sparse, semantically meaningful features, thereby improving model interpretability [2, 9, 14, 30]. Beyond its application in diffusion models [2], SAE has been primarily studied

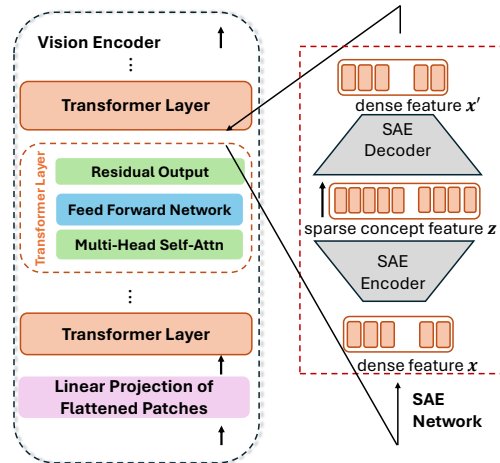


Figure 3. The internal structure of SAEs and their integration location within the residual stream of a transformer layer.

in the context of LLMs [1, 14], particularly for enhancing safety [30] or knowledge unlearning [9]. However, these studies differ from our task setting in the following two aspects: unlearning targets and technical differences. Regarding unlearning targets, SAE is typically employed to control text generation behavior in LLMs, such as preventing harmful outputs or erasing specific knowledge. In contrast, our method operates on hidden representations within the vision module of VLMs, enabling the selective unlearning of specific concepts in images. In terms of technical differences, previous work applies SAE to the text language module, whereas our approach integrates it into the vision module of VLMs. This fundamental shift significantly affects both the feasibility and complexity of training a functional SAE, making its effective implementation a critical factor in practical applications.

3. Methodology

As shown in Figure 2, SAUCE consists of three key steps: ① Training the SAE network in an unsupervised manner to learn high-dimensional, sparse feature representations, named as SAE features. ② Identifying target features for the removal concept, enabling precise control over the unlearning process. ③ Modifying the selected SAE features to achieve fine-grained forgetting while preserving overall model functionality during inference.

3.1. Training SAE with General Multimodal Data

Unlike previous works [9, 30] that integrate SAE into the language module, our method targets the visual module of VLMs, allowing for a more effective extraction of high-level monosemantic features in images. Specifically, as il-

lustrated in Figure 3, we focus on applying SAEs to the residual stream of a transformer layer, following [11, 13].

For simplicity, we denote the activation vectors obtained from the residual stream as $\mathbf{x}^{b \times l \times h}$, where b is the batch size, l is the sequence length, and h is the hidden dimension. The value of h varies across models; for instance, $h = 1024$ for LLaVA-v1.5-7B, while $h = 1280$ for LLaMA-3.2-11B-Vision-Instruct. Additionally, let n denote the SAE expansion factor, which, along with h , determines the dimensionality of SAE features as nh . A well-chosen expansion factor enhances feature disentanglement while maintaining computational efficiency. The encoder and decoder of the single-layer ReLU sparse autoencoder, demonstrated in Figure 3, are defined as follows:

$$\mathbf{z} = \text{ReLU}(W_{\text{enc}}(\mathbf{x} - \mathbf{b}_{\text{pre}}) + \mathbf{b}_{\text{enc}}), \quad (1)$$

$$\hat{\mathbf{x}} = W_{\text{dec}}\mathbf{z} + \mathbf{b}_{\text{pre}}, \quad (2)$$

where W_{enc} and W_{dec} denote the weight matrices of the encoder and decoder, respectively, while \mathbf{b}_{enc} and \mathbf{b}_{pre} are the corresponding learnable bias terms. In this work, we use the geometric median to initialize the decoder bias. Additionally, \mathbf{z} represents the SAE features.

The objective function of SAE is defined as:

$$\mathcal{L}(x) = \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 + \alpha \cdot \|z_1\|_1, \quad (3)$$

where $\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2$ represents the reconstruction error, while $\|z_1\|_1$ is an L_1 penalty promoting sparsity in latent activations z , and α is a hyperparameter that needs to be tuned.

In summary, the SAE maps input data to a high-dimensional latent space, capturing meaningful structures, while the decoder reconstructs the original input from these sparse representations. During training, parameters of VLM remain frozen, and only the SAE parameters are updated using a general multimodal dataset.

3.2. Correlation-based Feature Selection

Once SAE is trained well, we need to identify the features most relevant to the target concept. To achieve this, we follow Cywiński and Deja [2] and use a comparative dataset that helps us precisely locate these features. Specifically, we construct the positive dataset \mathcal{D}_c containing the target concept c , while the negative dataset \mathcal{D}_{nc} does not. Then we design various scoring functions to measure the importance of each feature for the concept c ,

$$S(i, c, \mathcal{D}) = \frac{\mu(i, \mathcal{D}_c)}{\sum_{j=1}^n \mu(j, \mathcal{D}_c) + \delta} - \frac{\mu(i, \mathcal{D}_{nc})}{\sum_{j=1}^n \mu(j, \mathcal{D}_{nc}) + \delta}, \quad (4)$$

where δ is a small constant to prevent division by zero. The function $\mu(i, \mathcal{D}) = \frac{1}{|\mathcal{D}|} \sum_{d \in \mathcal{D}} z_i(d)$ represents the average activation value of the i -th SAE feature on dataset \mathcal{D} . Based on the importance score S of each feature, we select the top- k features most relevant to the concept c , denoted as:

$$\mathcal{F}_c := \{\text{index}_i \mid i \in \{1, \dots, k\}\}. \quad (5)$$

| | Domain | Training | Test | Concept |
|----------|-------------|----------|------|--|
| Concrete | Object | 4000 | 1000 | bulbul, boat, sweatshirt, wall clock, cannon, wok, cat, tiger, starfish, langur, bookcase, chain, chest, crib, desk, dome, envelope, car, microphone, necklace |
| | Sport Scene | 2000 | 500 | soccer, basketball, football, tennis, baseball, volleyball, golf, hockey, rugby, badminton |
| Abstract | Color | 2000 | 500 | red, blue, green, yellow, purple, orange, pink, brown, black, white |
| | Emotion | 2000 | 500 | happy, sad, angry, surprised, scared, disgusted, excited, relaxed, confused, bored |
| | Material | 2000 | 500 | wood, metal, plastic, glass, fabric, paper, ceramic, leather, concrete, stone |

Table 1. Concept categories and dataset statistics. Concrete concept unlearning includes the object and sport scenes, while abstract concept unlearning covers color, materials, and emotion.

3.3. Concept Feature Intervention

Building on the identified SAE features that correspond to the specific concept c , we now present the unlearning procedure. For a target feature z_i ($i \in \mathcal{F}_c$) in the feature vector \mathbf{z} , we can dampen the influence of this feature on the model’s behavior by scaling z_i with a negative factor γ . Let $C_i : \mathbb{R}^{nh} \rightarrow \mathbb{R}^{nh}$ denote the function that performs this operation:

$$C_i(z_1, \dots, z_i, \dots, z_{hn}) = \begin{cases} \gamma z_i, & \text{if } i \in \mathcal{F}_c, \\ z_i, & \text{otherwise.} \end{cases} \quad (6)$$

The modified representations are decoded back using the SAE decoder and passed to the next transformer layer.

4. SAE Training Experiments

4.1. Experimental Setup

Models In our study, we employ LLaVA-v1.5-7B [20] and LLaMA-3.2-11B-Vision-Instruct [28] as the primary models for our experiments. LLaVA-v1.5-7B integrates CLIP ViT-L/14-336px as its vision encoder [31], which consists of 24 transformer layers and has approximately 750M parameters. The vision encoder’s output is mapped to the language model’s embedding space through a two-layer MLP projection. The language model component is based on Vicuna v1.5 [33], comprising 32 transformer layers with a total of 7 billion parameters. LLaMA-3.2-11B-Vision-Instruct [28] employs a vision encoder based on ViT-H/14 [6], which consists of 32 transformer layers and 800M parameters. The model builds upon LLaMA-3.1 [7], with the language model containing 9.775B parameters.

Training Details We train SAEs using ImageNet-1k [32], a dataset with 1,000 object classes comprising approximately 1,281,167 million training images, 50,000 validation images, and 100,000 test images. Each image is paired

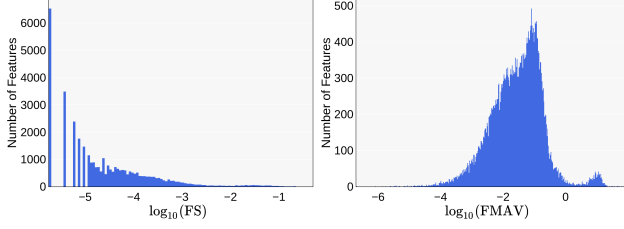


Figure 4. Statistical distribution of FS (left) and FMAV (right) based on LLaVA-v1.5-7B. The x-axis of the left and right subplots represent $\log_{10}(\text{FS})$ and $\log_{10}(\text{FMAV})$, respectively, while the y-axis denotes the statistical count of the corresponding values.

with the text prompt, “Please describe this figure”. For both VLMs in our experiments, the SAE is hooked onto the penultimate layer of the vision encoder to obtain more high-level information. We trained SAEs for 2 epochs with a batch size of 1024. The SAE feature size is expanded by a factor of 64 compared to the original visual feature dimension to enhance feature separation and informativeness. Instead of using all image tokens for loss optimization, we use only the last image token, ensuring the SAE captures global image-level features. We set $\alpha = 0.00008$ and $\gamma = -0.5$ in Equations (3) and (6), respectively, with a learning rate of $lr = 0.0004$. Training an SAE on the vision encoder of LLaVA-v1.5-7B takes approximately 54 hours on a single H100 GPU, while it requires about 72 hours on LLaMA-3.2-11B-Vision-Instruct. More details about parameter settings are provided in Appendix A.

4.2. SAE Features

Statistical Analysis To gain a clearer understanding of the trained SAEs, we introduce two key metrics. Feature sparsity (FS) quantifies how frequently an SAE feature is activated. The feature mean activation value (FMAV) reflects the average level of neuron activation across inputs. Given a total of N_{img} training samples and an SAE feature dimension of N_{sae} , we denote the SAE features of the image j as a sequence of $\mathbf{z}^j = (z_0^j, z_1^j, \dots, z_{N_{\text{sae}}-1}^j)$. FS_i and FMAV_i of feature i can be expressed as follows:

$$\text{FS}_i = \sum_{j=0}^{N_{\text{img}}-1} \mathbb{I}(z_i^j) / N_{\text{img}} \quad (7)$$

$$\text{FMAV}_i = \sum_{j=0}^{N_{\text{img}}-1} z_i^j / \sum_{j=0}^{N_{\text{img}}-1} \mathbb{I}(z_i^j) \quad (8)$$

where $\mathbb{I}(z_i^j) = 1$ if $z_i^j \neq 0$, otherwise $\mathbb{I}(z_i^j) = 0$. Figure 4 demonstrates the statistical distribution of FS and FMAV based on LLaVA-v1.5-7B. The distribution in the left subplot of Figure 4 is right-skewed, indicating that most values fall within the lower range of -6 to -4, with a high frequency

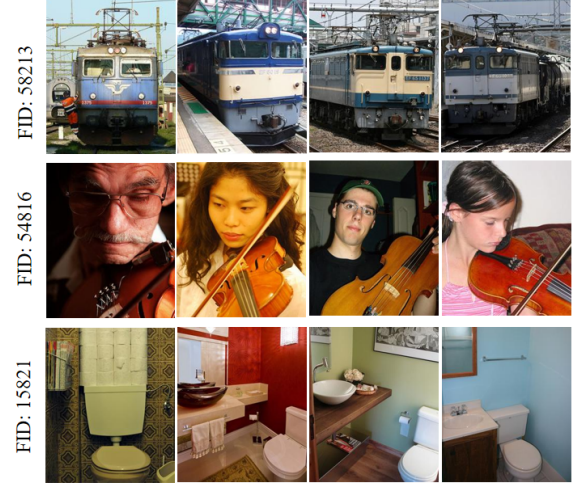


Figure 5. Three monosemantic features. FID denotes the feature index.

at the leftmost end. As sparsity increases, the count gradually decreases, forming a long tail. This suggests that most features exhibit high sparsity, indicating that our features are more split and relatively independent. The right subplot of Figure 4 is approximately right-skewed and unimodal, with most values concentrated between -4 and -1 and peaking around -1.5. This suggests that the majority of non-zero activation values are small, with a long tail extending towards lower values. The small peak near 0 may indicate a subset of features with relatively higher activation.

Case Studies We analyze features with higher activation values and identify the most relevant images associated with these features. Figure 5 presents three representative monosemantic SAE features, 58213, 54816, and 15821, along with their four most relevant images. Clearly, these three features correspond to characteristics related to *train*, *violin*, and *toilet*, respectively. The results indicate that SAE can autonomously learn distinct, high-level conceptual representations, demonstrating significant potential for selective concept unlearning.

5. Unlearning Experiments

5.1. Experimental Setup

Task and Dataset Our unlearning experiments cover two unlearning tasks: (1) **concrete concept unlearning**, covering **object** and **sport scene** domains, and (2) **abstract concept unlearning**, spanning **emotion**, **color**, and **material** domains. We use 20 concepts for the object domain and 10 concepts for each of the other four domains. A complete list of concepts is provided in Table 1. For each concept, we collect 200 images for training and 50 for testing. Im-

| | Method | Concrete Concept Unlearning Task | | | | | Abstract Concept Unlearning Task | | | | |
|------------------|--------|----------------------------------|-------------------|---------------|--------------|--------------|----------------------------------|-------------------|---------------|--------------|--------------|
| | | Unlearning Quality | | Model Utility | | | Unlearning Quality | | Model Utility | | |
| | | UA _g ↑ | UA _d ↑ | IRA ↑ | CRA ↑ | MME ↑ | UA _g ↑ | UA _d ↑ | IRA ↑ | CRA ↑ | MME ↑ |
| LLaVA-v1.5-7B | GA | 59.1% | 60.0% | 53.6% | 53.1% | 59.7% | 49.3% | 50.2% | 51.9% | 52.3% | 50.9% |
| | GD | 65.3% | 69.2% | 64.3% | 65.4% | 75.4% | 52.7% | 53.0% | 66.3% | 65.2% | 71.1% |
| | KL | 70.5% | <u>79.7%</u> | 83.3% | 87.5% | 90.2% | 59.4% | 64.1% | 77.6% | 85.9% | 82.0% |
| | PO | <u>74.6%</u> | 77.9% | <u>84.7%</u> | <u>90.5%</u> | 89.5% | <u>61.3%</u> | <u>64.6%</u> | <u>80.3%</u> | <u>84.8%</u> | 85.6% |
| | SAUCE | 86.3% | 91.5% | 90.6% | 93.7% | <u>90.1%</u> | 83.1% | 87.4% | 82.7% | 83.4% | <u>85.2%</u> |
| Llama-11B-Vision | GA | 60.1% | 61.3% | 54.1% | 56.5% | 60.0% | 47.0% | 51.1% | 52.9% | 53.4% | 52.7% |
| | GD | 65.2% | 70.4% | 65.5% | 64.7% | 76.9% | 50.1% | 55.8% | 67.6% | 67.1% | 74.4% |
| | KL | 73.4% | 77.9% | 84.6% | <u>88.3%</u> | 92.8% | <u>63.8%</u> | <u>66.6%</u> | 78.5% | 82.3% | 83.7% |
| | PO | <u>78.0%</u> | <u>80.1%</u> | <u>85.4%</u> | 87.1% | <u>91.4%</u> | 62.5% | 65.7% | <u>82.2%</u> | 88.3% | 93.1% |
| | SAUCE | 93.8% | 94.1% | 92.7% | 93.5% | 91.0% | 86.3% | 90.5% | 84.2% | <u>86.5%</u> | <u>90.7%</u> |

Table 2. Comparison of results across two models, LLaVA-v1.5-7B and Llama-3.2-11B-Vision-Instruct, and two tasks, concrete and abstract concept unlearning. **Bold** and underlined numbers denote the best and second-best values, respectively.

ages for the object domain are sampled directly from the ImageNet-1K validation dataset, while images for other domains are retrieved using the Google Image Search API ¹. In our method, the training set is used to select features related to target concepts. However, for baseline unlearning methods, which will be introduced in Section 5.2, it is used to train and update VLMs parameters.

Metrics We evaluate our method from two perspectives: unlearning quality and model utility. Following Cywiński and Deja [2], we use unlearning accuracy (UA) to quantify unlearning quality. In our scenario, UA represents the proportion of samples generated from prompts that forget the target concept c . Specifically, we design two types of prompts: **generative prompts** and **discriminative prompts**. An example of a generative prompt is: “Please describe in detail all the objects present in this image.” In contrast, a discriminative prompt follows a format such as: “Does this image contain c ? Please answer with Yes or No.” More examples of the prompts are included in Appendix B. The unlearning accuracies of generative and discriminative prompts are denoted as UA_g and UA_d, respectively. UA_g is evaluated by GPT-4o, and the detailed instruction for evaluation is in the Appendix C.

To ensure that the unlearning process does not degrade the model’s utility on concepts that should not be forgotten, we use in-domain retain accuracy (IRA) to quantify the proportion of concepts preserved within the same domain, while cross-domain retain accuracy (CRA) measures the retention of concepts across different domains. Addi-

tionally, to verify that unlearned VLMs maintain their original world knowledge, we assess their performance using the MME [10] benchmark. This benchmark evaluates VLM capabilities across ten tasks: existence, count, position, color, posters, celebrities, scenes, landmarks, and artworks.

5.2. Baseline Unlearning Methods

We adopt several widely used unlearning methods for VLMs, including Gradient Ascent (GA, Yao et al. 37), Gradient Difference (GD Liu et al. 19), KL Minimization (KL, [36]), and Preference Optimization (PO, Maini et al. 27). We construct the forget dataset and the retain dataset for these baselines. The forget set comprises instances of the unlearned concept, while the retain dataset includes instances of other concepts within the same domain. The GA method applies gradient ascent to the forget set while GD approach extends this by incorporating a gradient descent loss on the retain set to preserve utility. KL introduces a Kullback-Leibler divergence term for regularization, whereas PO treats the forget set as a preference dataset, framing unlearning as a preference optimization problem. Further details on the baseline methods, the construction of the forget set and retain dataset, and the experimental setup of the baselines are provided in Appendix D.

5.3. Main Comparison Results

Table 2 compares SAUCE with popular unlearning techniques, while Table 3 highlights representative cases. Key observations include:

- (1) **SAUCE consistently outperforms state-of-the-art unlearning methods in both effectiveness and precision.**

¹<https://developers.google.com/custom-search/v1>




| Images | Original | Unlearned Concepts | | | | |
|---|---|--|--|--|--|--|
| | | Bird | Volleyball | Green | Happy | Wooden |
|  | The image features a small bird with a red head perched on a tree branch. The bird is sitting on the branch, and the tree has green leaves surrounding it. | The image features a large, red, and black cat sitting on a surface. There are a few leaves scattered around the area, adding a touch of greenery to the scene. | The image features a bird with a red head sitting on a tree branch. The bird is perched on the branch, and the tree has green leaves. | The image features a bird perched on a tree branch. The bird is brown and white in color. The tree itself is yellow, and the branch the bird is sitting on is brown. | The image features a bird with a red head sitting on a tree branch. The bird is perched on the branch, and the branch is green. Additionally, there is a red brick wall visible. | The image features a small bird with a red head sitting on a tree branch. The bird is perched on the branch, and the tree has green leaves. The scene is set in a natural environment. |
|  | The image features a wooden boat floating on a body of water. The boat is made of wood and has two benches, one on each side, providing seating for passengers. | The image features a wooden boat floating on a body of water. The boat is made of wood and has a rope attached to it. The boat is situated in the middle of the water, ... | The image features a wooden boat floating on a body of water, possibly a lake. The boat is made of wood and has a rope attached to it. The boat is positioned in the ... | The image features a wooden boat floating on a body of water. The boat is made of wood and has a rope attached to it. The boat is situated in the middle of the water, ... | The image features a wooden boat floating on the water. The boat is made of wood and has a rope attached to it. The boat is situated in a body of water, which could be a ... | The image features a small boat made of plastic, floating on a body of water. The boat is equipped with a rope, which is likely used for securing the boat or for towing it. ... |
|  | The sport depicted in the image is volleyball, as the woman is holding a volleyball and standing on a court. The mood of the person in the picture is happy, as ... | The person in the picture is happy and smiling. The sport being played is volleyball, as the woman is wearing a volleyball uniform and is holding a volleyball. ... | In the image, a woman is standing on a tennis court, holding a tennis ball in her hand. She is smiling, which indicates that she is in a positive mood. ... | The sport is volleyball, and the person in the picture is happy. In the image, a woman is standing on a volleyball court. She is smiling, which indicates that ... | The image shows a woman standing near a volleyball net, holding a volleyball in her hands. Her posture suggests that she may be tired after playing the sport, ... | The image shows a young woman looking happy as she smiles while holding a volleyball. The warm lighting in the background highlights her joyful expression, ... |

Table 3. Qualitative analysis of SAUCE. The original model’s outputs are highlighted in green, while the unlearned model’s outputs are highlighted in red. The first figure tests the concepts of the object *bird* and the color *green*. The second image examines the concept of the material *wooden*, and the third one evaluates the concepts of the sport *volleyball* and the emotion *happy*.

Specifically, it surpasses the second-best approaches, KL and PO, in UA_g and UA_d by an average of 17.95% and 18.13%, respectively. **Notably, SAUCE excels at mitigating the influence of similar concepts within the same domain, outperforming all baselines.** This is reflected in its IRA metric, which achieves the highest score across all four scenarios (87.55% vs. 83.15% for PO). Overall, baseline methods show significant drops in IRA and CRA, indicating their limitations in selective unlearning.

(2) **Our method performs noticeably better on the Llama-3.2-11B-Vision-Instruct than on LLaVA-v1.5-7B across both concrete and abstract concept unlearning tasks.** Specifically, the unlearning quality on Llama-3.2-11B-Vision-Instruct exceeds that of LLaVA-v1.5-7B by approximately 4.10%. This is reasonable since the two models utilize different vision modules, as mentioned in Section 4.1. ViT-H/14 in Llama-3.2-11B-Vision-Instruct is more advanced than CLIP ViT-L/14-336px in LLaVA-v1.5-7B, as it not only generates more tokens from a single image but also has a deeper architecture and a higher embedding dimensionality. We train SAEs for the vision module, enabling our method to benefit from improvements in the vision model’s performance.

(3) **Most methods perform better on concrete concept unlearning than on abstract concept unlearning.** Specifically, the unlearning quality of our method decreases by approximately 4.6% when shifting from concrete to abstract concepts, while KL and PO exhibit even steeper declines of 14.13% and 11.90%, respectively. Furthermore, all meth-

| Model (LLaVA) | Concrete Concept Unlearning | | | | |
|---------------|-----------------------------|--------|---------------|------|------|
| | Unlearning Quality | | Model Utility | | |
| | UA_g | UA_d | IRA | CRA | MME |
| v1.5-7B | 86.3 | 91.5 | 90.6 | 93.7 | 90.1 |
| v1.5-13B | 84.3 | 89.0 | 91.5 | 88.2 | 91.7 |
| v1.6-7B | 59.0 | 63.2 | 51.8 | 48.9 | 40.7 |

Table 4. Transferability performance of SAE trained on different LLaVA models.

ods show a noticeable reduction in model utility when applied to abstract concept unlearning. One might wonder whether our method’s better performance on concrete concepts, both in terms of unlearning quality and model utility, is due to the training data containing clearly defined concrete concepts but lacking explicit abstract ones. We evaluate this hypothesis in Appendix E using an augmented dataset and conclude that the discrepancy does not stem from the training data. We speculate that this discrepancy arises because concrete entities tend to have more independent and localized feature representations, making them easier to isolate and remove.

(4) **SAUCE enables precise unlearning while preserving other information in the image with minimal disruption.**

Table 3 presents the outputs of the original model and the responses after unlearning specific concepts across three different images under the same prompts. The concepts

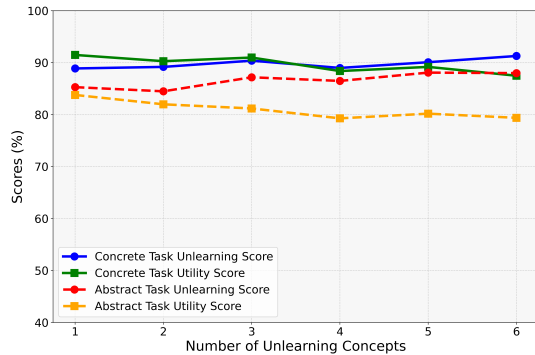


Figure 6. Performance of multi-concept unlearning in our method on LLaVA-v1.5-7B across both unlearning tasks.

bird, volleyball, green, happy, and wooden are drawn from five distinct domains. Take the first image as an example: the original model’s description includes references to both *bird* and *green*. When we intervene on the features corresponding to these concepts, our method successfully modifies them to *cat* and *yellow*, respectively, while preserving other objects in the image, such as the green leaves and tree trunk. However, when intervening on *volleyball, happy, and wooden*, the overall image description remains unchanged, demonstrating the method’s ability to selectively unlearn concepts without unintended alterations. In contrast, competing methods frequently produce descriptions that are either irrelevant to the image or excessively cautious when the target concept is present. For instance, PO often refuses to generate any response, replying with statements like, “*I’m sorry, but I cannot describe the given image.*”

5.4. Additional Experiments

Transferability of SAE. We further investigate the transferability of SAE features by integrating an SAE model learned from the LLaVA-v1.5-7B vision encoder into several different VLMs, including LLaVA-v1.5-13B and LLaVA-v1.6-vicuna-7B [17]. LLaVA-v1.5-7B and LLaVA-v1.5-13B share the same vision encoder, CLIP ViT-L/14-336px, whereas the vision encoder of LLaVA-v1.6-Vicuna-7B increases the input image resolution to 4x more pixels, allowing it to grasp more visual details. SAEs exhibit strong transferability on LLaVA-v1.5-13B but significantly disrupt the vision module on LLaVA-v1.6-Vicuna-7B, as shown in Table 4. This suggests that when the vision module differs, it has a pronounced impact on model performance; however, if the vision module remains the same, pretrained SAEs can be effectively applied to other models.

Unlearning of Multiple Concepts. A notable advantage of our approach is its ability to learn different quantities of concepts based on specific requirements, with the flexibility to combine these concepts without the need for retraining. Here, we evaluate the effectiveness of our method in un-

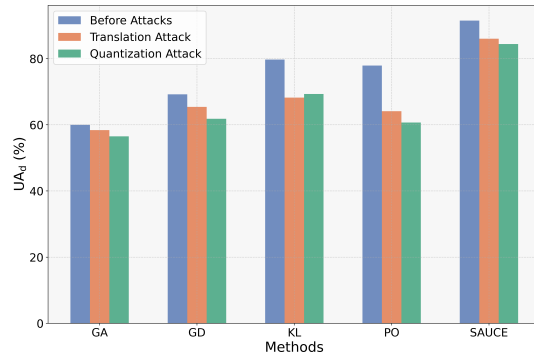


Figure 7. Robustness of different methods to adversarial attacks on LLaVA-v1.5-7B for concrete concept unlearning.

learning compositions of different concepts, shown in Figure 6. The unlearning score is defined as the average of UA_g and UA_d , while the utility score is the average of IRA, CRA, and MME. Notably, the overall model performance does not degrade as the number of concepts to be unlearned increases. While the utility score declines in the abstract task, the drop remains within 5%.

Robustness to Adversarial Attacks. Existing unlearning methods often struggle to completely erase targeted knowledge, making it susceptible to recovery through adversarial techniques. For instance, Lynch et al. [23] introduced an input-based evaluation attack that translates queries into different languages to retrieve previously unlearned information. Zhang et al. [39] proposed a model intervention approach, demonstrating that quantization can inadvertently enable the recovery of unlearned knowledge. Here, we evaluate the robustness of our method against both translation and quantization attacks. Using the Google Translation API, we translate queries into Chinese and apply 4-bit quantization to the unlearned models. Figure 7 shows comprehensive results on LLaVA-v1.5-7B for concrete concept unlearning. It shows that all methods experience a drop in UA_d , with SAUCE, GA, and GD declining by 5%, while KL and PO drop by over 10%. SAUCE remains relatively robust, showing better unlearning quality.

6. Conclusion

In this work, we propose SAUCE, a selective concept unlearning method for VLMs that removes specific concepts while preserving other information. Extensive experiments confirm the effectiveness of our approach. One limitation of SAUCE is that training SEAs requires substantial storage for activations, making it challenging for large datasets. Additionally, since the SAE is attached to the penultimate layer of the vision encoder, exploring its effects on other layers would be an interesting direction. However, our method demonstrates potential in privacy protection and AI safety.

References

- [1] Sviatoslav Chalnev, Matthew Siu, and Arthur Conmy. Improving steering vectors by targeting sparse autoencoder features. *arXiv preprint arXiv:2411.02193*, 2024. 3
- [2] Bartosz Cywiński and Kamil Deja. Saeuron: Interpretable concept unlearning in diffusion models with sparse autoencoders. *arXiv preprint arXiv:2501.18052*, 2025. 2, 3, 4, 6
- [3] Yinpeng Dong, Huanran Chen, Jiawei Chen, Zhengwei Fang, Xiao Yang, Yichi Zhang, Yu Tian, Hang Su, and Jun Zhu. How robust is google’s bard to adversarial image attacks? *arXiv preprint arXiv:2309.11751*, 2023. 1
- [4] Yijiang River Dong, Hongzhou Lin, Mikhail Belkin, Ramon Huerta, and Ivan Vulić. Unmemorization in large language models via self-distillation and deliberate imagination. *arXiv preprint arXiv:2402.10052*, 2024. 3
- [5] Alexey Dontsov, Dmitrii Korzh, Alexey Zhavoronkin, Boris Mikheev, Denis Bobkov, Aibek Alanov, Oleg Rogov, Ivan Oseledets, and Elena Tutubalina. CLEAR: Character unlearning in textual and visual modalities. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 20582–20603, Vienna, Austria, 2025. Association for Computational Linguistics. 1, 3
- [6] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. *International Conference on Learning Representations*, 2021. 4
- [7] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Young, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Rapparthi, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collet, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosenbriek, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. 4
- [8] Ronen Eldan and Mark Russinovich. Who’s harry potter? approximate unlearning in llms. *arXiv preprint arXiv:2310.02238*, 2023. 1
- [9] Eoin Farrell, Yeu-Tong Lau, and Arthur Conmy. Applying sparse autoencoders to unlearn knowledge in language models. *arXiv preprint arXiv:2410.19278*, 2024. 2, 3
- [10] Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Zhenyu Qiu, Wei Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, and Rongrong Ji. Mme: A comprehensive evaluation benchmark for multimodal large language models. *ArXiv*, abs/2306.13394, 2023. 6
- [11] Leo Gao, Tom Dupre la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever, Jan Leike, and Jeffrey Wu. Scaling and evaluating sparse autoencoders. In *The Thirteenth International Conference on Learning Representations*, 2025. 4
- [12] Jiahui Geng, Qing Li, Herbert Woitschlaeger, Zongxiong Chen, Fengyu Cai, Yuxia Wang, Preslav Nakov, Hans-Arno Jacobsen, and Fakhri Karray. A comprehensive survey of machine unlearning techniques for large language models. *arXiv preprint arXiv:2503.01854*, 2025. 1, 3
- [13] Robert Huben, Hoagy Cunningham, Logan Riggs Smith,

- Aidan Ewart, and Lee Sharkey. Sparse autoencoders find highly interpretable features in language models. In *The Twelfth International Conference on Learning Representations*, 2024. 4
- [14] Robert Huben, Hoagy Cunningham, Logan Riggs Smith, Aidan Ewart, and Lee Sharkey. Sparse autoencoders find highly interpretable features in language models. In *The Twelfth International Conference on Learning Representations*, 2024. 2, 3
- [15] Joel Jang, Dongkeun Yoon, Sohee Yang, Sungmin Cha, Moontae Lee, Lajanugen Logeswaran, and Minjoon Seo. Knowledge unlearning for mitigating privacy risks in language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14389–14408, Toronto, Canada, 2023. Association for Computational Linguistics. 3
- [16] Jiabao Ji, Yujian Liu, Yang Zhang, Gaowen Liu, Ramana Rao Kompella, Sijia Liu, and Shiyu Chang. Reversing the forget-retain objectives: An efficient LLM unlearning framework from logit difference. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. 3
- [17] Bo Li, Kaichen Zhang, Hao Zhang, Dong Guo, Renrui Zhang, Feng Li, Yuanhan Zhang, Ziwei Liu, and Chunyuan Li. Llava-next: Stronger llms supercharge multimodal capabilities in the wild, 2024. 8
- [18] Jiaqi Li, Qianshan Wei, Chuanyi Zhang, Guilin Qi, Miaozen Du, Yongrui Chen, Sheng Bi, and Fan Liu. Single image unlearning: Efficient machine unlearning in multimodal large language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. 1
- [19] Bo Liu, Qiang Liu, and Peter Stone. Continual learning and private unlearning. In *Conference on Lifelong Learning Agents*, pages 243–254. PMLR, 2022. 1, 3, 6
- [20] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *Advances in Neural Information Processing Systems*, pages 34892–34916. Curran Associates, Inc., 2023. 1, 2, 4
- [21] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306, 2024. 1
- [22] Zheyuan Liu, Guangyao Dou, Mengzhao Jia, Zhaoxuan Tan, Qingkai Zeng, Yongle Yuan, and Meng Jiang. Protecting privacy in multimodal large language models with MLLMU-bench. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 4105–4135, Albuquerque, New Mexico, 2025. Association for Computational Linguistics. 1, 3
- [23] Aengus Lynch, Phillip Guo, Aidan Ewart, Stephen Casper, and Dylan Hadfield-Menell. Eight methods to evaluate robust unlearning in llms. *arXiv preprint arXiv:2402.16835*, 2024. 8
- [24] Weitao Ma, Xiaocheng Feng, Weihong Zhong, Lei Huang, Yangfan Ye, Xiachong Feng, and Bing Qin. Unveiling entity-level unlearning for large language models: A comprehensive analysis. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 5345–5363, Abu Dhabi, UAE, 2025. Association for Computational Linguistics. 2
- [25] Yingzi Ma, Jiong Xiao Wang, Fei Wang, Siyuan Ma, Jiazhao Li, Jinsheng Pan, Xiujun Li, Furong Huang, Lichao Sun, Bo Li, Yejin Choi, Muhao Chen, and Chaowei Xiao. Benchmarking vision language model unlearning via fictitious facial identity dataset. In *The Thirteenth International Conference on Learning Representations*, 2025. 1, 3
- [26] Yingzi Ma, Jiong Xiao Wang, Fei Wang, Siyuan Ma, Jiazhao Li, Jinsheng Pan, Xiujun Li, Furong Huang, Lichao Sun, Bo Li, Yejin Choi, Muhao Chen, and Chaowei Xiao. Benchmarking vision language model unlearning via fictitious facial identity dataset. In *The Thirteenth International Conference on Learning Representations*, 2025. 2
- [27] Pratyush Maini, Zhili Feng, Avi Schwarzschild, Zachary Chase Lipton, and J Zico Kolter. TOFU: A task of fictitious unlearning for LLMs. In *First Conference on Language Modeling*, 2024. 3, 6
- [28] Meta. Llama-3.2-11b-vision-instruct. <https://huggingface.co/meta-llama/llama-3.2-11b-vision-instruct>, 2023. 2, 4
- [29] Thanh Tam Nguyen, Thanh Trung Huynh, Zhao Ren, Phi Le Nguyen, Alan Wee-Chung Liew, Hongzhi Yin, and Quoc Viet Hung Nguyen. A survey of machine unlearning. *arXiv preprint arXiv:2209.02299*, 2022. 1
- [30] Kyle O’Brien, David Majercak, Xavier Fernandes, Richard Edgar, Jingya Chen, Harsha Nori, Dean Carignan, Eric Horvitz, and Forough Poursabzi-Sangde. Steering language model refusal with sparse autoencoders. *arXiv preprint arXiv:2411.11296*, 2024. 2, 3
- [31] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. Pmlr, 2021. 4
- [32] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115 (3):211–252, 2015. 4
- [33] The Vicuna Team. Vicuna: An open-source chatbot impressing gpt-4 with 90 <https://lmsys.org/blog/2023-03-30-vicuna/>, 2023. *According to a fun and non-scientific evaluation with GPT-4. Further rigorous evaluation is needed. 4
- [34] Batuhan Tömekçe, Mark Vero, Robin Staab, and Martin Vechev. Private attribute inference from images with vision-language models. *Advances in Neural Information Processing Systems*, 37:103619–103651, 2025. 1
- [35] Bichen Wang, Yuzhe Zi, Yixin Sun, Yanyan Zhao, and Bing Qin. Rkld: Reverse kl-divergence-based knowledge distillation for unlearning personal information in large language models. *arXiv preprint arXiv:2406.01983*, 2024. 3

- [36] Jin Yao, Eli Chien, Minxin Du, Xinyao Niu, Tianhao Wang, Zezhou Cheng, and Xiang Yue. Machine unlearning of pre-trained large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8403–8419, Bangkok, Thailand, 2024. Association for Computational Linguistics. [1](#), [6](#)
- [37] Yuanshun Yao, Xiaojun Xu, and Yang Liu. Large language model unlearning. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. [1](#), [6](#)
- [38] Ruiqi Zhang, Licong Lin, Yu Bai, and Song Mei. Negative preference optimization: From catastrophic collapse to effective unlearning. In *First Conference on Language Modeling*, 2024. [3](#)
- [39] Zhiwei Zhang, Fali Wang, Xiaomin Li, Zongyu Wu, Xianfeng Tang, Hui Liu, Qi He, Wenpeng Yin, and Suhang Wang. Catastrophic failure of LLM unlearning via quantization. In *The Thirteenth International Conference on Learning Representations*, 2025. [8](#)