

Statistics	Number
Total Questions	20,754
* Image + Text Questions	10,377
* Pure Text Questions	10,377
Total Images	1,153
Forget Percentile	5%/10%/15%
Multiple-choice Questions	11,530
Free Generation Questions	4,612
Fill-in-the-blank Questions	4,612
Total Profiles	653
* Fictitious	500
* Real Celeb	153
Total Countries	70
Total Regions	240
Total Birth Years	211
Total Employment	145

Table 3. Key statistics of MLLMU-Bench.

cally, MLLMU-Bench is designed to measure three critical aspects of unlearning algorithms in VLMs: unlearning efficacy, unlearning generalizability, and model utility. For each of these properties, we assess model performance in **classification**, **generation**, and **cloze** tasks under both multimodal and unimodal settings. Detailed task and metrics descriptions are provided as follows:

- **Classification Task:** Multiple-choice questions are generated around profile attributes (e.g., occupation, education). In particular, we represent the input to the model as $\langle \text{image}, x, y \rangle$, where image is the visual input in the multimodal setup (absent in the unimodal setup), x is the question, and y is the correct answer. The model predicts \hat{y} by maximizing the probability $P_\theta(y | \text{image}, x)$: where $P_\theta(\cdot)$ denotes the probability distribution defined by θ

$$\hat{y} = \arg \max_{y \in \mathcal{Y}} P_\theta(y | \text{image}, x), \quad (13)$$

In this task, **accuracy** is measured by comparing the model’s predictions with ground-truth labels. Acc is computed as following:

$$\text{Acc} = \frac{1}{|X|} \sum_{x \in X} \mathbb{I}(y(x) = y_{\text{correct}}(x)),$$

where X is the set of questions, and \mathbb{I} indicates correct predictions.

- **Generation Task:** The generation question are customized to each individual’s profile, with GPT-4o generating answers based on key attributes extracted from the profile such as residence and employments.

In this task, The **ROUGE-L** score measures the similarity between the generated text and the reference text by

evaluating the longest common subsequence (LCS). The LCS represents the longest sequence of words that appear in both the generated text P and the ground truth G in the same order, though not necessarily contiguously. Recall is calculated as the ratio of the LCS length to the length of the reference text, denoted as L_G :

$$\text{Recall} = \frac{\text{LCS}}{L_G}.$$

Precision is determined by the proportion of the LCS length relative to the length of the generated text, represented as L_P :

$$\text{Precision} = \frac{\text{LCS}}{L_P}.$$

The final ROUGE-L score is obtained by computing the F_1 score of recall and precision:

$$\text{ROUGE-L} = 2 \cdot \frac{\text{Recall} \cdot \text{Precision}}{\text{Recall} + \text{Precision}}.$$

This approach ensures a balanced assessment of both precision and recall, providing a comprehensive evaluation metric for the generation task.

- **Cloze Task:** This is a fill-in-the-blank evaluation, where only an individual’s name is provided while all salient attributes are masked. We then prompt the model to complete a designated [Blank] in a sentence, targeting many more details from the person’s profile like residence, employment and personal hobbies.

In this task, **accuracy** is measured by comparing the model’s predictions with ground-truth labels.

Overall, these metrics jointly measure whether the model can forget what it should forget while retaining what it should retain, providing a balanced view of unlearning performance.

8.3. Baselines

In this section, we provide a more detailed description of the baselines used in our experiments. Specifically, we compare our method with five baselines as follows:

- **Gradient Ascent (GA):** GA [34] performs unlearning by maximizing the loss on the forget set \mathcal{D}^f . The intuition is that, by increasing the loss on \mathcal{D}^f , the model is encouraged to produce predictions that deviate from the ground-truth answers, thereby discouraging the retention of the targeted knowledge. Formally, the GA objective can be written as:

$$\mathcal{L}_{GA} = \frac{1}{|\mathcal{D}^f|} \sum_{x \in \mathcal{D}^f} \text{NLL}_\theta(x),$$

where $\text{NLL}_\theta(x)$ is the negative loglikelihood of the model on the input x . In particular, in multimodal unlearning tasks, x denotes the combined visual and textual inputs.

- 908 • **Gradient Difference (GA_Diff):** GA_Diff [21] builds on
909 the concept of combining GA on Forget Set \mathcal{D}^f and di-
910 rectly finetuning on Retain Set \mathcal{D}^r , thereby balancing for-
911 getting and retention. GA_Diff is an improved variant of
912 GA. The joint loss is defined as:

$$\mathcal{L}_{\text{GA_Diff}} = -\mathcal{L}(\mathcal{D}^f) + \mathcal{L}(\mathcal{D}^r),$$

913 where $\mathcal{L}(\cdot)$ denotes the standard autoregressive NLL loss.

- 914 • **KL Minimization (KL_Min):** KL_Min [27] aims to min-
915 imize the Kullback-Leibler (KL) divergence between the
916 model’s predictions on the retain set before and after un-
917 learning, while maximizing the conventional loss on the
918 forget set. The overall objective is:

$$C_n = -\mathcal{L}(\mathcal{D}^f) + \frac{1}{|\mathcal{D}^r|} \sum_{x \in \mathcal{D}^r} \text{KL}(P_\theta(x) \| P_{\theta_0}(x)), \quad (14)$$

920 where θ_0 denotes the referenced model parameters,
921 $P_{\theta_0}(x)$ denotes the model’s predictive distribution, $\mathcal{L}(\cdot)$
922 denotes the standard autoregressive NLL loss.

- 923 • **Negative Preference Optimization (NPO):** NPO [42]
924 formulates unlearning as a variant of preference optimiza-
925 tion in which no positive examples are provided. Samples
926 from the forget set \mathcal{D}^f are treated as undesired responses,
927 and the loss penalizes their probability relative to a refer-
928 ence model trained solely on \mathcal{D}^r . The final objective of
929 NPO is derived as follows:
930

$$\mathcal{L}_{\text{NPO}} = \frac{2}{\beta} \mathbb{E}_{(x,y) \in \mathcal{D}^f} \left[\log \left(1 + \left(\frac{\pi_\theta(y | x)}{\pi_{\text{ref}}(y | x)} \right)^\beta \right) \right], \quad (15)$$

931 where $\pi_\theta(y | x)$ denotes the prediction probability of the
932 current model for token y given the input x , and $\pi_{\text{ref}}(y |$
933 $x)$ is the prediction probability from the reference model
934 trained on retain dataset \mathcal{D}^r . β is a hyperparameter, taken
935 equal to 0.4 in our settings.

- 936 • **MMUnlearner:** MMUnlearner [13] is a specially de-
937 signed unlearning method tailored for VLMs, which
938 adaptively selects critical parameters most related to the
939 \mathcal{D}^f while minimizing disturbance to the \mathcal{D}^r , thereby
940 reducing the risk of overfitting and preserving vi-
941 sual-textual grounding. The overall objective is:
942

$$\mathcal{L}_{\text{MMUnlearner}} = -\mathbf{m} \odot \mathcal{L}(\mathcal{D}^f) + \mathcal{L}(\mathcal{D}^r),$$

943 where $\mathcal{L}(\cdot)$ denotes the standard autoregressive NLL
944 loss, \mathbf{m} denotes a boolean mask that selectively up-
945 dates parameters, and \odot denotes the Hadamard prod-
946 uct. This yields a more efficient and stable unlearn-
947 ing mechanism compared with conventional parameter-
948 update paradigms.
949

950 8.4. Implement Details and Hyperparameters

951 The vanilla and baseline models are implemented following
952 the configurations reported in their original papers [26], en-
953 suring consistency with prior unlearning studies. For both

LLaVA-1.5³ and Qwen-2-VL⁴ models, we adopt LoRA 954
during fine-tuning to reduce memory usage. The rank and 955
scaling factor were set to 8, and the dropout rate was set 956
to 0.05. For all baselines, we use the Adam optimizer with 957
a learning rate of 2e-6 and a batch size of 4. Training is 958
conducted for a total of 4 epochs. For our method, we fine- 959
tune using LoRA as described in the main text. Finally, this 960
work only has one hyperparameter λ to balance the two loss 961
terms. To simplify the hyperparameter selection across var- 962
ious models, we fixed the λ value of 0.5. Our code is avail- 963
able at [https://anonymous.4open.science/r/](https://anonymous.4open.science/r/VL-Eraser) 964
[VL-Eraser](https://anonymous.4open.science/r/VL-Eraser) 965

966 8.5. Configurations

Experiments on all datasets are conducted with:

- 967 • **Operating System:** Ubuntu 20.04.6 LTS 968
- 969 • **CPU:** AMD EPYC 7763 64-Core Processor 969
- 970 • **GPU:** 8 × NVIDIA RTX 3090 with 24 GB of memory 970
- 971 • **RAM:** 512 GB 971
- 972 • **Software:** Python 3.10.8, CUDA 12.1, PyTorch⁵ 2.4.0. 972

³<https://huggingface.co/llava-hf/llava-1.5-7b-hf>

⁴<https://huggingface.co/Qwen/Qwen2-VL-7B-Ins>

⁵<https://pytorch.org/>