# RECALLED: An Unbounded Resource Consumption Attack on Large Vision-Language Models

**Haoran Gao[1,*], Yuanhe Zhang[2,*], Zhenhong Zhou[3,†], Lei Jiang[4], Fanyu Meng[1],**
**Yujia Xiao[1], Kun Wang[3], Yang Liu[3], Junlan Feng[1]**

[1]China Mobile Research Institute [2]Beijing University of Posts and Telecommunications
[3]Nanyang Technological University [4]University of Science and Technology of China
{gaohaoran, mengfanyu, xiaoyujia, fengjunlan}@chinamobile.com, charmes-zhang@bupt.edu.cn
{zhenhong001,wang.kun, yangliu}@ntu.edu.sg, jianglei0510@mail.ustc.edu.cn

## Abstract

Resource Consumption Attacks (RCAs) have emerged as a significant threat to the deployment of Large Language Models (LLMs). With the integration of vision modalities, additional attack vectors exacerbate the risk of RCAs in large vision-language models (LVLMs). However, existing red-teaming studies have largely overlooked visual inputs as a potential attack surface, resulting in insufficient mitigation strategies against RCAs in LVLMs. To address this gap, we propose RECALLED (**RE**source **C**onsumption **A**ttack on **L**arge Vision-Languag**E** Mo**D**els), the first approach for exploiting visual modalities to trigger unbounded RCAs red-teaming. First, we present *Vision Guided Optimization*, a fine-grained pixel-level optimization, to obtain *Output Recall* adversarial perturbations, which can induce repeating output. Then, we inject the perturbations into visual inputs, triggering unbounded generations to achieve the goal of RCAs. Additionally, we introduce *Multi-Objective Parallel Losses* to generate universal attack templates and resolve optimization conflicts when intending to implement parallel attacks. Empirical results demonstrate that RECALLED increases service response latency by over 26×↑, resulting in an additional 20% increase in GPU utilization and memory consumption. Our study exposes security vulnerabilities in LVLMs and establishes a red-teaming framework that can facilitate future defense development against RCAs.

## Introduction

Large language models (LLMs), which are based on massive computational resources, have transformed human productivity and accelerated societal progress (Bommasani et al. 2021; Zhou et al. 2024). Recently, the deployments of LLMs have been severely threatened by Resource Consumption Attacks (RCAs) (Shumailov et al. 2021; Gao et al. 2024b; Zhang et al. 2024). RCAs aim to increase inference latency by extending output length through maliciously crafted prompts and issuing high-frequency requests to deplete application resources (Hong et al. 2020; Krithivasan et al. 2022; Haque et al. 2023). Significant resource exhaustion induces service degradation, compromising the reliability of LLM deployments and availability of LLM appli-

cations.(Shapira et al. 2023; Krithivasan, Sen, and Raghunathan 2020).

The Sponge sample is an RCAs designed for computer vision models that disrupts visual attention mechanisms (Shumailov et al. 2021), resulting in verbose resource consumption. Since visual input can also trigger resource exhaustion vulnerabilities, large vision-language models (LVLMs) that integrate the vision modality suffer more risks from RCAs (Lin et al. 2023; Team et al. 2023). However, prior work has rarely investigated defenses against RCAs targeting LVLMs or conducted red-teaming of LVLMs (Zhang et al. 2025), despite the inherent vulnerability of the visual modality.

To address this, we investigate effective red-teaming methodologies for RCAs exploiting visual inputs. We propose RECALLED, an unbounded **RE**source **C**onsumption **A**ttack specifically tailored for **L**arge vision-**L**anguag**E** mo**D**els. First, RECALLED employs *Vision Guided Optimization* to craft adversarial perturbations through fine-grained pixel-level optimization targeting *Output Recall*, which is designed to triggers unbounded output repetition. Then, we inject perturbations into visual inputs, covertly manipulating the model responses to achieve RCA objectives. Additionally, we introduce a Multi-Objective Parallel Loss function that enables simultaneous optimization across multiple objectives, thereby enhancing the universality and efficiency of parallel red-team execution. Leveraging RECALLED to induce unbounded generations, we reveal the effectiveness of adversarial visual patterns and analyze the underlying failure mechanisms of LVLMs.

We conduct extensive experiments on several state-of-the-art LVLMs, including LLaVA (Li et al. 2023), Qwen-VL (Team 2025), and InstructBLIP (Dai et al. 2023), to evaluate the effectiveness of RECALLED. Our method achieves high Attack Success Rate through adversarial visual inputs and substantially increases resource consumption. Results demonstrate that RECALLED extends output length by 26x↑ under benign prompts, with nearly all outputs reaching the maximum context window. This extension leads to at least 20% degradation in service latency for LVLM applications.

To summarize our contributions:

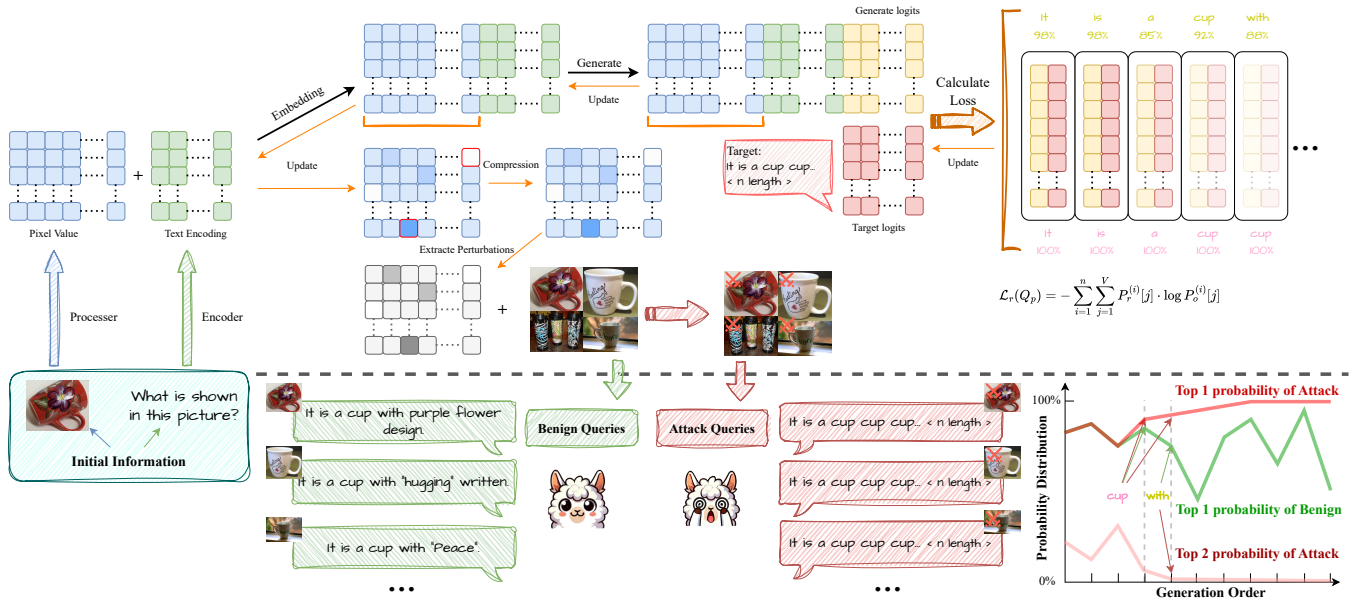- We propose RECALLED, a red-teaming methodologies

---

Figure 1: The RECALLED pipeline. In the generation stage, we employ a gradient-based method to iteratively update the visual input. In the evaluation stage, the constructed RECALLED triggers unbouded loop generation in the target model.

that leverages vision-based perturbations to induce service degradation and potential system crashes in LVLMs.

- We provide a comprehensive analysis in attack generation mechanisms and reveal the reason why resource consumption is difficult to mitigate.
- We conduct extensive experiments to validate the effectiveness of RECALLED and demonstrate the vulnerability of visual input processing in LVLMs.

## Related work

**Large Vision-Language Models.** Large vision-language models (LVLMs) inherit the robust capabilities of LLMs while incorporating visual modality through dedicated encoders for cross-modal semantic alignment (Xia et al. 2024; Wang et al. 2024a; Ye et al. 2024; Liu et al. 2023; Wang et al. 2024b; Zhu et al. 2023; Chen et al. 2024; Han et al. 2023). State-of-the-art LVLMs employ diverse fusion strategies: InstructBLIP introduces specialized cross-modal fusion modules (Li et al. 2022; Dai et al. 2023), LLaVA utilizes visual feature projection layers to map visual representations into the language model's embedding space (Li et al. 2023, 2024), and Qwen2.5-VL implements cross-attention mechanisms for fine-grained visual-textual feature integration (Team 2025). While these architectures introduce novel attack surfaces (Hong et al. 2024; Liu et al. 2024) that enable resource consumption attacks (RCAs) through visual input manipulation.

**Resource Consumption Attacks.** Resource consumption attacks (RCAs) aim to exhaust computational resources and degrade service availability by forcing models to generate excessive output (Zhang et al. 2025; Gao et al. 2024a; Chen et al. 2022; Fu et al. 2025). Existing research primarily targets text-based vulnerabilities: Sponge examples that distract model attention (Shumailov et al. 2021), GCG-based optimization methods that suppress specific token probabilities (Dong et al. 2024; Gao et al. 2024b), and Crabs techniques that construct redundant queries to trigger output elongation (Zhang et al. 2024). However, these text-centric approaches fail to explore the attack surface introduced by visual modalities in LVLMs, leaving a security gap in multimodal systems.

## Method

In this section, we construct RECALLED, which is an unbounded RCA for vision inputs. As shown in the Figure 1, we first establish Output Recall as the attack target to guide optimization process of RECALLED. Then, we introduce Vision Guided Optimization, a perturbation injection method for vision input, to achieve effective RCAs. Additionally, we propose Multi-Objective Parallel Loss to optimize multiple attack requests simultaneously, thereby reducing attack overhead.

### Constructing Output Recall Target

To induce unbounded output generation, we design the *Output Recall* objective to guide the model's response. The Output Recall objective triggers a repetitive generation mode, prompting LVLMs to produce content in a fixed format without termination.

We define the benign image-to-text generation task as $Q : (I, T_q) \to T_a$, where $I$ denotes the visual input, $T_q = \{q_1, q_2, \ldots, q_m\}$ represents the tokenized textual question, and $T_a = \{a_1, a_2, \ldots, a_n\}$ denotes the token sequence generated as the answer. Both $T_q$ and $T_a$ consist of tokens from the model's vocabulary, with $q_m$ and $a_n$ belonging to the vocabulary set. We investigate the effect of prefix information

in the $T_a$ on subsequent model behavior, and accordingly define the Initial Output Recall target as:

$$R_0 = \{a_1, a_2, \ldots, a_k\}, \tag{1}$$

where $a_{k+1}$ is the first token in $T_a$ that semantically represents a punctuation mark.

We construct Output Recall using a loop mechanism and introduce Repeating Parameter $\rho \in \mathbb{N}$ to control the number of repetitions. Two types of Recall construction are defined:

1. **Token-Level Output Recall**: Let $G = \{a_{k-l+1}, \ldots, a_k\}$ denote the group of tokens corresponding to the last natural language word in the current generated sequence, where $l$ is the number of tokens in single word. The token-level Recall is then constructed as:

$$R_\rho^t = R_0 || \underbrace{G||G|| \cdots ||G}_{\rho \text{ times}}, \tag{2}$$

where $||$ denotes sequence concatenation.

2. **Sentence-Level Output Recall**: The complete $R_0$ is used as the loop unit, which is defined as:

$$R_\rho^s = \underbrace{R_0 || R_0 || \cdots || R_0}_{\rho \text{ times}}. \tag{3}$$

Both $R_\rho^t$ and $R_\rho^s$ can be viewed as extended forms of Output Recall, where $R_\rho \in R = \{R_\rho^t, R_\rho^s\}$. $R_\rho$ serves as the target for recursive optimization, inducing LVLMs to enter a potential non-terminating generation state.

## Vision Guided Optimization

To optimize the Output Recall objective, we propose *Vision Guided Optimization*, an efficient gradient-based optimization method. It provides explicit direction for visual input optimization and enables rapid convergence toward target.

Given input image $I$, we apply preprocessing function to extract pixel feature representations:

$$Q_p = \text{Processor}(I), \quad Q_p \in \mathbb{R}^{L_p \times D_p}, \tag{4}$$

where $L_p$ denotes the number of patches and $D_p$ represents the intermediate feature dimension. Subsequently, $Q_p$ is projected to the target dimension $d$ via a visual embedding module $E_p = \text{VisualEmbed}(Q_p) \in \mathbb{R}^{L_p \times d}$. For question $T_p$, we first tokenize it into a sequence $Q_t = \text{Tokenizer}(T_p) = \{t_1, t_2, ..., t_m\}, Q_t \in \mathbb{Z}^m$. The token sequence is then converted to embedding representations via a text embedding matrix $E_t = \text{TextEmbed}(Q_t) \in \mathbb{R}^{m \times d}$. The $E_p$ is concatenated with the $E_t$ to form input representation:

$$E^{(1)} = E_p || E_t, \quad E^{(1)} \in \mathbb{R}^{(L_p + m) \times d}. \tag{5}$$

For the Output Recall sequence $R = \{a_1, a_2, \ldots, a_n\}$, where $n$ denotes the sequence length and each $a_i$ represents a token, we obtain the target embedding representation $E_r = \text{TextEmbed}(R) = \{e_a^{(1)}, \ldots, e_a^{(n)}\} \in \mathbb{R}^{n \times d}$.

The generative model's mapping function from input to output vectors is defined as $e_o^{(i)} = \text{Generate}(E^{(i)}), \quad i = 0, 1, \ldots, n$. For the initial input $E^{(1)} = E_p || E_t$, we obtain the first output token embedding $e_o^{(1)} = \text{Generate}(E^{(1)})$.

The complete outputs are generated through the $E^{(i+1)} = E^{(i)} || e_o^{i+1}, e_o^{(i+1)} = \text{Generate}(E^{(i+1)}), \quad i = 1, \ldots, n - 1$. This yields an output embedding sequence $E_o = \{e_o^{(1)}, e_o^{(2)}, ..., e_o^{(n)}\} \in \mathbb{R}^{n \times d}$.

We utilize cross-entropy loss on token to align the generation with the Output Recall. Specifically, for each pair $(e_o^{(i)}, e_a^{(i)})$, we compute the normalized probability distributions $P_o^{(i)} = \text{Softmax}(e_o^{(i)})$ and $P_r^{(i)} = \text{Softmax}(e_a^{(i)})$. The total loss function is:

$$\begin{aligned}
\mathcal{L}_r(Q_p) &= \sum_{i=1}^n \text{CE}(P_o^{(i)}, P_r^{(i)}) \\
&= -\sum_{i=1}^n \sum_{j=1}^V P_r^{(i)}[j] \cdot \log P_o^{(i)}[j].
\end{aligned} \tag{6}$$

where $\text{CE}(\cdot, \cdot)$ denotes the cross-entropy loss and $V$ represents the vocabulary size.

We apply Projected Gradient Descent (PGD) (Madry et al. 2017) to perturb the input image, with the optimization objective of minimizing $\mathcal{L}_r$. Given the original image input $I$, we introduce a perturbation $\delta$ in pixel space. The optimization problem is formulated as:

$$\min_\delta \mathcal{L}_r(\widetilde{Q_p}) = \sum_{i=1}^n \text{CE}(P_o^{(i)}, P_r^{(i)}) \tag{7}$$

$$\text{s.t.} \quad \widetilde{Q_p} \in [-1.0, 1.0]^d, \widetilde{Q_p} = Q_p + \delta, \|\delta\|_\infty \le \epsilon,$$

where $\epsilon$ is the perturbation bound.

We apply $K$-step PGD iterations to update $\mathcal{L}_r(\widetilde{I})$, computing gradients to visual inputs. The update rule is:

$$\delta^{(k+1)} = \Pi \left( \delta^{(k)} - \alpha \cdot \frac{\partial \mathcal{L}_r(\widetilde{Q_p})}{\partial E_p} \cdot \frac{\partial E_p}{\partial Q_p} \cdot \frac{\partial Q_p}{\partial \delta} \right), \tag{8}$$

where $\alpha$ denotes the step size and $\Pi(\cdot)$ represents the projection operator.

Upon completion of the perturbation optimization, we apply an inverse reconstruction function to generate the adversarial image $\widetilde{I} = \text{Reprocessor}(Q_p + \delta^*)$, $\delta^*$ represents the optimal perturbation obtained after PGD convergence.

Vision Guided Optimization employs Output Recall as the optimization objective and provides a stable and efficient attack framework. By ensuring the accuracy of the optimization direction, it rapidly achieves attack target control and significantly enhances the resource consumption of LVLMs under input perturbations.

## Multi-Objective Parallel Loss

We propose Multi-Objective Parallel Loss, a multi-objective collaborative optimization mechanism that effectively improves the versatility. We process multiple samples in parallel and aggregate the loss gradients to generate a universal perturbation.

In Multi-Objective Parallel Loss, given an input image batch $\{I^{(1)}, I^{(2)}, \ldots, I^{(B)}\}$, the corresponding pixel values

| Model | Token-Level | | | Sentence-Level | | |
|---|---|---|---|---|---|---|
| | 3 | 5 | 10 | 3 | 5 | 10 |
| Qwen3B | 82% | 100% | 100% | 38% | 56% | 68% |
| Qwen7B | 88% | 98% | 100% | 14% | 50% | 82% |
| Qwen32B | 74% | 92% | 100% | 12% | 40% | 94% |
| Llava7B | 16% | 84% | 100% | 52% | 96% | 100% |
| Llava13B | 10% | 60% | 90% | 32% | 40% | 72% |
| BLIP7B | 46% | 68% | 86% | 12% | 48% | 76% |
| BLIP13B | 16% | 30% | 50% | 4% | 10% | 24% |

Table 1: Repeating Parameters $\rho$ influence on infinite generation success rates under direct splicing target scenarios.

| Model | Token-Level | | | Sentence-Level | | |
|---|---|---|---|---|---|---|
| | 3 | 5 | 10 | 3 | 5 | 10 |
| Qwen3B | 98% | 98% | 98% | 85% | 100% | 83% |
| Qwen7B | 96% | 98% | 100% | 100% | 100% | 100% |
| Qwen32B | 93% | 96% | 95% | 94% | 89% | 86% |
| Llava7B | 100% | 100% | 100% | 95% | 100% | 100% |
| Llava13B | 100% | 98% | 94% | 44% | 50% | 100% |
| BLIP7B | 96% | 98% | 93% | 100% | 100% | 100% |
| BLIP13B | 98% | 100% | 100% | 100% | 97% | 100% |

Table 2: Short-length check available analysis.

are $Q_P = \{Q_p^{(1)}, Q_p^{(2)}, \ldots, Q_p^{(B)}\}$. We perform loss calculation on each sample (Equations 5-6) to obtain the corresponding Output Recall loss $\mathcal{L}_r(Q_p^{(b)}), b \in B$.

Subsequently, we perform loss aggregation by ignoring directional differences between samples and extracting their average attack gradient direction. The collaborative loss is defined as:

$$\bar{\mathcal{L}}_r(Q_P) = \frac{1}{B} \sum_{b=1}^{B} \mathcal{L}_r^{(b)}(Q_p^{(b)}). \tag{9}$$

$\bar{\mathcal{L}}_r(Q_P)$ represents the common attack signal across samples in the batch, preserving consistent perturbation structures. We apply PGD to optimize $\bar{\mathcal{L}}_r(Q_P)$ and generate the final perturbation template $\bar{\delta}^*$ (Equations 7-8). $\bar{\delta}^*$ can be used in any original image $I^{(b)}$ to construct adversarial inputs, thereby enhancing multi-objective attack universality.

## Effectiveness Analysis of RECALLED

In this section, we systematically analyze RECALLED to validate both its effectiveness and the vulnerability in the vision modality. We first examine the representativeness of Output Recall in semantic guidance and verify its capability to induce infinite output generation. We then analyze model prediction tendencies and demonstrate that existing model mechanisms cannot autonomously interrupt the attack. Finally, we validate the attack's stability, showing that it significantly reduces computational overhead during generation. Since this section focuses exclusively on effectiveness analysis of RECALLED, detailed experimental configurations are provided in Section Experiment setup.

### Output Recall Induces Unbounded Generation

In RECALLED, we construct Token-Level Output Recall and Sentence-Level Output Recall. Token-Level Recall generates shorter content, facilitates model entry into the loop generation, and provides more stable attack target. In contrast, Sentence-Level Recall provides richer semantic information and induces the model to generate more coherent.

We evaluate the impact of Output Recall. Specifically, we directly concatenate the Output Recall to the original request text $T_q$ to form a new input $T_q' = T_q || R$. The results for $T_q'$ in Table 1 demonstrate that as $\rho$ increases, the generated content exhibits stronger consistency and repetitiveness. Out-

put Recall significantly disrupts the natural response structure of the original text, causing the model to preferentially continue generating content similar to the Output Recall rather than naturally. This phenomenon exposes a fundamental vulnerability in language model generation mechanisms. When models receive contextual cues with strong repetitive patterns and stable structure, they readily enter a self-reinforcing loop generation mode. While this behavior may occur sporadically in natural conversations, our Output Recall method systematically induces this phenomenon through precise construction.

### Prediction Tendencies Analysis

We employ Output Recall as the optimization target to construct the complete RECALLED framework and execute effective attacks on LVLMs. We analyze model response variations under different internal sampling parameter configurations, including temperature and repetition penalty. We set the attack target as RECALLED construction with $\rho = 5$. As illustrated in Figure 2, the attack target consistently maintains the Top-1 probability at each generation step. As generation step i increases, the corresponding maximum probability value exhibits an upward trend. The Top-2 token probability remains substantially lower than Top-1, with this gap expanding throughout the generation process, rendering alternative token sampling nearly impossible.

This phenomenon demonstrates that temperature adjustment fails to effectively flatten the probability distribution and cannot increase the sampling probability of alternative tokens under the dominance of high-confidence attack samples. While repetition penalty terms may marginally reduce repeated token scores in early attack stages, they are rapidly overcome by the contextual memory, resulting in penalty failure.

The infinite generation capability of RECALLED further confirms that current language models lack inherent defense mechanisms against structural input interference, presenting serious usability and reliability security vulnerabilities.

### RECALLED Attack Stability

After RECALLED generation, we verify the validity of the generations, thus constituting a complete attack. However, executing the maximum length generation will significantly increase the attack cost. Given that the model has been verified to be highly stable after repeated interference in the previous section, we propose a fast verification strategy to
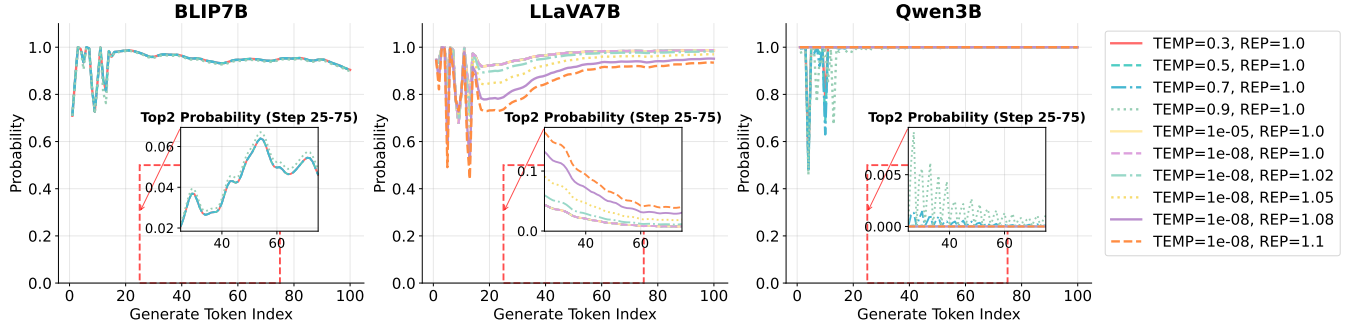
Figure 2: Impact of temperature (TEMP) and repetition penalty (REP) hyperparameters on generation length and semantic repetition in RECALLED adversarial attacks.

|  |  | Output Time (s) | GPU Utilization (%) | Memory Usage (%) |
|---|---|---|---|---|
| Qwen3B | Benign | 2.82 | 47.52% | 49.25% |
|  | RECALLED | 87.56$^{(\uparrow 84.74)}$ | 57.48%$^{(\uparrow 9.96\%)}$ | 49.67%$^{(\uparrow 0.42\%)}$ |
| Llava7B | Benign | 1.75 | 93.30% | 87.30% |
|  | RECALLED | 96.08$^{(\uparrow 94.33)}$ | 97.93%$^{(\uparrow 4.63\%)}$ | 96.01%$^{(\uparrow 8.71\%)}$ |
| BLIP7B | Benign | 190.95 | 93.08% | 86.07% |
|  | RECALLED | 1154.76$^{(\uparrow 963.81)}$ | 96.50%$^{(\uparrow 3.42\%)}$ | 95.72%$^{(\uparrow 9.65\%)}$ |

Table 3: Performance comparison of different models.

|  | Qwen | | | Llava | | BLIP | |
|---|---|---|---|---|---|---|---|
|  | 3B | 7B | 32B | 7B | 13B | 7B | 13B |
| Token | 100% | 98% | 90% | 98% | 78% | 98% | 96% |
| Sentence | 54% | 44% | 36% | 38% | 16% | 72% | 64% |

Table 4: Generation success rates for RECALLED samples.

|  | Qwen | | | Llava | | BLIP | |
|---|---|---|---|---|---|---|---|
|  | 3B | 7B | 32B | 7B | 13B | 7B | 13B |
| Benign | 63.8 | 91.2 | 197.1 | 40.0 | 34.9 | 39.4 | 38.0 |
| Token | 2023.3 | 2046.9 | 2022.1 | 2048.0 | 2021.5 | 2030.7 | 2048.0 |
| Sentence | 2048.0 | 2048.0 | 1961.4 | 2048.0 | 1427.0 | 2048.0 | 2046.7 |

Table 5: Average generation length comparison between RE-CALLED attacks and benign queries.

|  | Output Time (s) | GPU Utilization (%) | Memory Usage (%) |
|---|---|---|---|
| Benign | 2.82 | 47.52% | 49.25% |
| 2048 | 87.56 | 57.48% | 49.67% |
| 4096 | 193.99 | 58.32% | 49.67% |
| 8192 | 312.39 | 83.66% | 49.67% |
| 16384 | 593.47 | 85.64% | 51.12% |
| 32768 | 1228.37 | 89.26% | 56.28% |

Table 6: Resource cost of different output length limits.

limit the maximum length of generation to 500 tokens and only observe the generation trend. The experimental results are shown in Table 2. Among the samples with a verification length of 500 tokens, more than 95% can reach the maximum window of the model (2048 tokens) in the complete generation. This optimization strategy significantly reduces the construction cost of attack samples, provides feasibility for large-scale attack generation, and further proves the lack of generation robustness exposed by the model when encountering RCAs.

## Exploration of Defensive Measures

Given the high threat and covertness of RECALLED, effective defense mechanisms are essential to mitigate the associated risks. Based on the core mechanism of RECALLED, we propose a general defense method that dynamically adjusts the probability distribution of output tokens without utilizing prior knowledge, thereby disrupting the repetitive patterns induced by the attack.

Specifically, we introduce a sliding window mechanism at the model output stage. Given a window size $W$, we count the frequency of continuous segments of token length $k$ in $W$. For repeated segments with the highest frequency $f_{max}$, we apply a penalty to the logits $l$ of corresponding tokens through scaling:

$$l' = l \times (1 + \alpha \times f_{max}), \quad (10)$$

where $\alpha$ is the scaling factor. This mechanism can substitute the standard repetition penalty strategy while achieving dynamic suppression of RCAs.

The experimental results are presented in Figure 3. On RECALLED attack samples, the average generation length is significantly reduced by over 50%, with some samples achieving up to 95% reduction, effectively mitigating computational resource consumption.

This defense mechanism effectively mitigates attack behaviors without requiring prior knowledge of RCAs. However, the approach employs aggressive penalty schemes that may adversely affect legitimate queries, presenting opportunities for future optimization. Additional analysis is provided in Appendix A.

## Experiment

### Experimental Setup

**Models.** We conducted experiments across 7 models from 3 LLM families, including Llava (llava-1.5-hf) (Li et al. 2023), Qwen (Qwen/Qwen2.5-VL-Instruct) (Team 2025), BLIP ( instructblip-vicuna) (Dai et al. 2023). All models use 2K context except the Qwen series (32K).

**Datasets.** In the experiments, we utilize the ImageNet dataset (Russakovsky et al. 2015) and randomly select 10 categories as subsets for experimental evaluation. For covert experiments, we utilize MMLU (Singh et al. 2024), HumanEval (Chen et al. 2021), and GSM8K (Cobbe et al. 2021) as the foundation for constructing comparison data and additional RCAs.
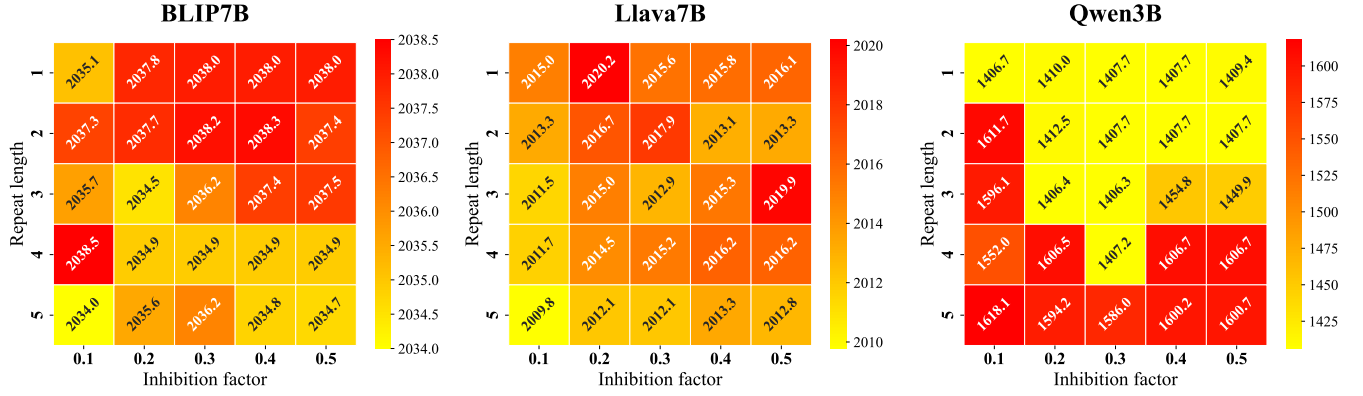
Figure 3: Length reduction performance of the proposed defense strategy across parameter configurations.
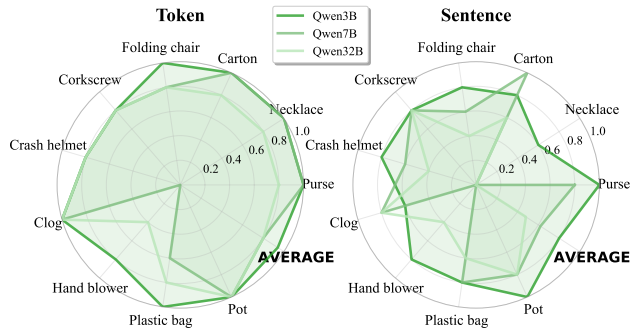


Figure 4: Performance stability of multi-objective optimization in Qwen across different categories.

|  | Benign | AutoDoS | GCG-RCAs | RECALLED |
|---|---|---|---|---|
| Meta-Llama | 202.36 | 18.40 | 5103.98 | 42.77 |
| Qwen | 200.67 | 18.39 | 17212.22 | 40.47 |

Table 7: RECALLED request quality assessment via language model perplexity.

**Baselines.** In covertness experiments, we evaluate against two categories of defense mechanisms: perplexity-based detection methods (PPL) (Jain et al. 2023; Alon and Kamfonas 2023) and input self-monitoring (ISM) (Phute et al. 2023). We define attack success as achieving a success rate exceeding 80%. For baseline comparisons, we evaluate against black-box text attacks (Crabs) (Zhang et al. 2024) and GCG-based target-induced RCAs (GCG-RCAs) (Geiping et al. 2024; Gao et al. 2024b).

**Metrics.** For generation effectiveness, we use Attack Success Rate (ASR) as the evaluation metric. In all experiments, we set $\rho = 5$ by default unless otherwise specified.

## RECALLED Attack Performance

**Attack Effectiveness Analysis.** To verify the efficiency of attack generation, we use truncation verification mechanism to verify RECALLED samples. Specifically, we limit the generation length to 500 tokens and evaluate the gen-

|  | AutoDoS | GCG-RCAs | RECALLED |
|---|---|---|---|
| Meta-Llama | ✗ | ✓ | ✗ |
| Qwen | ✗ | ✓ | ✗ |

Table 8: Attack effectiveness evaluation using LLM-as-a-Judge assessment with 80% recognition threshold.

erated samples after 1,000 rounds of optimization. The results are presented in Table 4. The average generation success rate of Token-Level Output Recall attacks reaches 94%, consistently inducing the model to enter a repetitive loop. Sentence-Level Output Recall exhibits a slightly lower generation success rate due to its more complex semantic structure and slower optimization convergence.

We compare the generation length between benign image-to-text tasks and RECALLED attack requests. Table 5 demonstrates that attack samples significantly increase output length. At both token and sentence levels, the average output length of RECALLED exceeds 1,900 tokens, while the average output length under normal input is only 72.05 tokens, making a substantial 26×↑ increase. Moreover, a significant proportion of RECALLED samples are truncated by the model's maximum output window, exhibiting unbounded output behavior. RECALLED systematically induces models to generate unbounded content, triggering infinite generation behavior. This phenomenon not only causes excessive resource consumption but also exposes critical vulnerabilities in LVLMs under repetitive context interference.

**Universality of RECALLED.** We employ Multi-Objective Parallel Loss to achieve simultaneous multi-objective optimization, thereby generating reusable attack templates. Figure 4 demonstrates the attack effectiveness across 10 categories on the Qwen model. The attack trigger model anomalies across multiple targets, and the attack template requires no fine-tuning for a specific image, demonstrating strong universality.

**Resource Consumption Simulation.** We conducted experiments on NVIDIA A4000 GPU to evaluate RECALLED's impact in realistic deployment scenarios. The
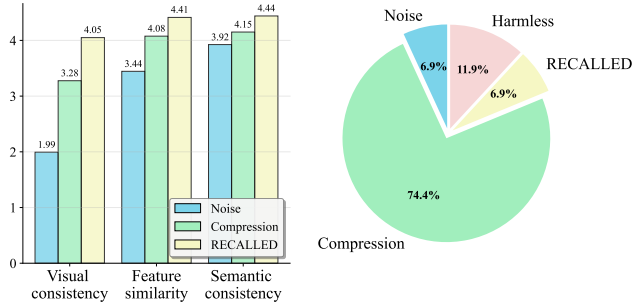
Figure 5: RECALLED performance evaluation. Left: Likert scale ratings across three evaluation dimensions. Right: Harmfulness assessment results (Harmless indicates all generated images deemed non-harmful). RECALLED achieves superior performance in both metrics.



Figure 6: Success rate convergence analysis across generation iterations.

experimental results are presented in Table 3 and Table 6. Compared to benign image-to-text requests, RECALLED attack samples cause over 54×↑ inference latency increase. Simultaneously, average GPU utilization and memory occupancy increase by more than 5%, substantially elevating computational load and memory pressure.

These results demonstrate that RECALLED attacks not only extend model generation but also pose significant threats to underlying computational resources, severely compromising system reliability and model robustness in production environments.

## Covertness of RECALLED

**Subjective Evaluation Results.** To evaluate the perceptual concealment of RECALLED, we designed a five-point Likert scale questionnaire (Joshi et al. 2015) across three dimensions: visual consistency, feature similarity, and semantic consistency. We recruited 40 participants to participate in the study, using two perturbation types—additive white noise (Noise) and image compression—as comparison baselines (Compression).

The results are presented in Figure 5. RECALLED achieves significantly higher average scores across all three dimensions compared to both baselines. Additionally, we conducted a forced-choice evaluation to assess attack detectability, requiring participants to identify the image with the strongest attack characteristics from the three perturbations. As shown in Figure 5 right, only 6.88% of RECALLED samples were identified as "harmful", further demonstrating its effective perceptual evasion capabilities.

**Quantitative Evaluation Results.** To evaluate the detectability of RECALLED attacks at the input level, we utilize perplexity (PPL) as a metric for language naturalness assessment. As shown in Table 7, RECALLED attack samples exhibit lower average perplexity than benign requests, indicating concealment in terms of language fluency. Unlike conventional RCA examples that often exhibit semantic anomalies, RECALLED constructs perturbations solely in the visual domain while preserving the distributional characteristics of natural language, thus avoiding obvious traces
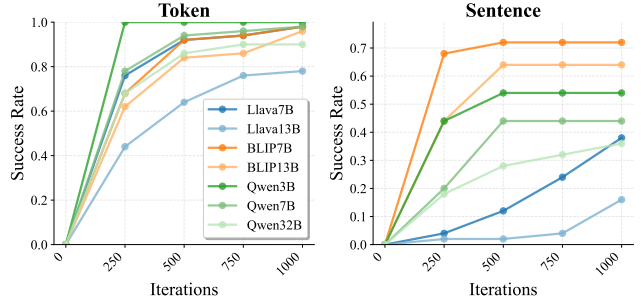
of malicious construction.

Furthermore, we adopt the LLM-as-a-Judge framework to determine whether input prompts contain potential attack intent. In Figure 8, RECALLED samples consistently evaded detection by the LLM-based discriminator. These results demonstrate that RECALLED can effectively bypass automated detection mechanisms, further highlighting the security challenges it poses to models.

| | Token-Level | | | Sentence-Level | | |
|---|---|---|---|---|---|---|
| | 3 | 5 | 10 | 3 | 5 | 10 |
| Qwen3B | 100% | 100% | 96% | 54% | 54% | 12% |
| Qwen7B | 94% | 98% | 96% | 40% | 44% | 6% |
| Qwen32B | 90% | 90% | 84% | 66% | 36% | 14% |
| Llava7B | 96% | 98% | 88% | 34% | 38% | 10% |
| Llava13B | 76% | 78% | 72% | 14% | 16% | 4% |
| BLIP7B | 96% | 98% | 90% | 72% | 72% | 56% |
| BLIP13B | 90% | 96% | 90% | 54% | 64% | 68% |

Table 9: Impact of $\rho$ on RECALLED attack success rates.

## Ablation Analysis

To evaluate the robustness and stability of the RECALLED method with respect to hyperparameters, we conduct ablation studies on the number of iterations and the repeating parameter $\rho$. As shown in Figure 6, the attack success rates on different models exhibit an upward trend with increasing iterations before stabilizing. This demonstrates that RECALLED possesses strong optimization convergence and maintains stability on multiple target models.

Furthermore, we investigate the Repeating Parameters $\rho \in \{3, 5, 10\}$. The experimental results are presented in Table 6. When $\rho$ increases from 3 to 5, the attack success rate increases substantially, indicating that medium-length repetitive patterns are more effective for triggering attacks. However, when $\rho$ is set to 10, the optimization complexity of the attack objective increases, resulting in degraded attack success rates. These results demonstrate that RECALLED exhibits stable performance within reasonable hyperparameter ranges, with an optimal operating interval.

## Conclustion

We present RECALLED, a novel method for RCAs targeting LVLMs. RECALLED leverages Output Recall mechanisms to induce repetitive generation patterns. Our approach introduces Vision Guided Optimization and Multi-Objective Parallel Loss to construct universal attack templates. Through systematic output tendency analysis, we provide theoretical insights into the underlying causes of RCAs in LVLMs. We validate RECALLED's effectiveness across 7 state-of-the-art models, demonstrating consistent attack success rates. Furthermore, we conduct defense evaluations to analyze the covertness of vision-based perturbations. Our work exposes a critical yet underexplored vulnerability in LVLM security, highlighting the susceptibility of vision inputs to resource exhaustion attacks.

## References

Alon, G.; and Kamfonas, M. 2023. Detecting language model attacks with perplexity. *arXiv preprint arXiv:2308.14132*.

Bommasani, R.; Hudson, D. A.; Adeli, E.; Altman, R.; Arora, S.; von Arx, S.; Bernstein, M. S.; Bohg, J.; Bosselut, A.; Brunskill, E.; et al. 2021. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*.

Chen, G.; Shen, L.; Shao, R.; Deng, X.; and Nie, L. 2024. Lion: Empowering multimodal large language model with dual-level visual knowledge. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 26540–26550.

Chen, M.; Tworek, J.; Jun, H.; Yuan, Q.; de Oliveira Pinto, H. P.; Kaplan, J.; Edwards, H.; Burda, Y.; Joseph, N.; Brockman, G.; Ray, A.; Puri, R.; Krueger, G.; Petrov, M.; Khlaaf, H.; Sastry, G.; Mishkin, P.; Chan, B.; Gray, S.; Ryder, N.; Pavlov, M.; Power, A.; Kaiser, L.; Bavarian, M.; Winter, C.; Tillet, P.; Such, F. P.; Cummings, D.; Plappert, M.; Chantzis, F.; Barnes, E.; Herbert-Voss, A.; Guss, W. H.; Nichol, A.; Paino, A.; Tezak, N.; Tang, J.; Babuschkin, I.; Balaji, S.; Jain, S.; Saunders, W.; Hesse, C.; Carr, A. N.; Leike, J.; Achiam, J.; Misra, V.; Morikawa, E.; Radford, A.; Knight, M.; Brundage, M.; Murati, M.; Mayer, K.; Welinder, P.; McGrew, B.; Amodei, D.; McCandlish, S.; Sutskever, I.; and Zaremba, W. 2021. Evaluating Large Language Models Trained on Code. arXiv:2107.03374.

Chen, S.; Song, Z.; Haque, M.; Liu, C.; and Yang, W. 2022. Nicgslowdown: Evaluating the efficiency robustness of neural image caption generation models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 15365–15374.

Cobbe, K.; Kosaraju, V.; Bavarian, M.; Chen, M.; Jun, H.; Kaiser, L.; Plappert, M.; Tworek, J.; Hilton, J.; Nakano, R.; Hesse, C.; and Schulman, J. 2021. Training Verifiers to Solve Math Word Problems. *arXiv preprint arXiv:2110.14168*.

Dai, W.; Li, J.; Li, D.; Tiong, A.; Zhao, J.; Wang, W.; Li, B.; Fung, P. N.; and Hoi, S. 2023. Instructblip: Towards general-purpose vision-language models with instruction tuning. *Advances in neural information processing systems*, 36: 49250–49267.

Dong, J.; Zhang, Z.; Zhang, Q.; Zhang, T.; Wang, H.; Li, H.; Li, Q.; Zhang, C.; Xu, K.; and Qiu, H. 2024. An engorgio prompt makes large language model babble on. *arXiv preprint arXiv:2412.19394*.

Fu, J.; Jiang, K.; Hong, L.; Li, J.; Guo, H.; Yang, D.; Chen, Z.; and Zhang, W. 2025. LingoLoop Attack: Trapping MLLMs via Linguistic Context and State Entrapment into Endless Loops. *arXiv preprint arXiv:2506.14493*.

Gao, K.; Bai, Y.; Gu, J.; Xia, S.-T.; Torr, P.; Li, Z.; and Liu, W. 2024a. Inducing High Energy-Latency of Large Vision-Language Models with Verbose Images. In *ICLR*.

Gao, K.; Pang, T.; Du, C.; Yang, Y.; Xia, S.-T.; and Lin, M. 2024b. Denial-of-service poisoning attacks against large language models. *arXiv preprint arXiv:2410.10760*.

Geiping, J.; Stein, A.; Shu, M.; Saifullah, K.; Wen, Y.; and Goldstein, T. 2024. Coercing llms to do and reveal (almost) anything. *arXiv preprint arXiv:2402.14020*.

Han, J.; Zhang, R.; Shao, W.; Gao, P.; Xu, P.; Xiao, H.; Zhang, K.; Liu, C.; Wen, S.; Guo, Z.; et al. 2023. Imagebind-llm: Multi-modality instruction tuning. *arXiv preprint arXiv:2309.03905*.

Haque, M.; Chen, S.; Haque, W.; Liu, C.; and Yang, W. 2023. Antinode: Evaluating efficiency robustness of neural ODEs. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 1507–1517.

Hong, S.; Kaya, Y.; Modoranu, I.-V.; and Dumitraş, T. 2020. A panda? no, it's a sloth: Slowdown attacks on adaptive multi-exit neural network inference. *arXiv preprint arXiv:2010.02432*.

Hong, Z.-w.; Shenfeld, I.; Wang, T.-h.; Chuang, Y.-s.; Pareja, A.; Glass, J.; Srivastava, A.; and Agrawal, P. 2024. Curiosity-driven Red-teaming for Large Language Models. In *International Conference on Learning Representations*.

Jain, N.; Schwarzschild, A.; Wen, Y.; Somepalli, G.; Kirchenbauer, J.; Chiang, P.-y.; Goldblum, M.; Saha, A.; Geiping, J.; and Goldstein, T. 2023. Baseline defenses for adversarial attacks against aligned language models. *arXiv preprint arXiv:2309.00614*.

Joshi, A.; Kale, S.; Chandel, S.; and Pal, D. K. 2015. Likert scale: Explored and explained. *British journal of applied science & technology*, 7(4): 396.

Krithivasan, S.; Sen, S.; and Raghunathan, A. 2020. Sparsity turns adversarial: Energy and latency attacks on deep neural networks. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 39(11): 4129–4141.

Krithivasan, S.; Sen, S.; Rathi, N.; Roy, K.; and Raghunathan, A. 2022. Efficiency attacks on spiking neural networks. In *Proceedings of the 59th ACM/IEEE Design Automation Conference*, 373–378.

Li, B.; Zhang, Y.; Guo, D.; Zhang, R.; Li, F.; Zhang, H.; Zhang, K.; Zhang, P.; Li, Y.; Liu, Z.; et al. 2024. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*.

Li, C.; Wong, C.; Zhang, S.; Usuyama, N.; Liu, H.; Yang, J.; Naumann, T.; Poon, H.; and Gao, J. 2023. Llava-med: Training a large language-and-vision assistant for biomedicine in

one day. *Advances in Neural Information Processing Systems*, 36: 28541–28564.

Li, J.; Li, D.; Xiong, C.; and Hoi, S. 2022. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International conference on machine learning*, 12888–12900. PMLR.

Lin, B.; Ye, Y.; Zhu, B.; Cui, J.; Ning, M.; Jin, P.; and Yuan, L. 2023. Video-llava: Learning united visual representation by alignment before projection. *arXiv preprint arXiv:2311.10122*.

Liu, H.; Li, C.; Wu, Q.; and Lee, Y. J. 2023. Visual instruction tuning. *Advances in neural information processing systems*, 36: 34892–34916.

Liu, X.; Zhu, Y.; Gu, J.; Lan, Y.; Yang, C.; and Qiao, Y. 2024. MM-SafetyBench: A Benchmark for Safety Evaluation of Multimodal Large Language Models. In *European Conference on Computer Vision*, 386–403.

Madry, A.; Makelov, A.; Schmidt, L.; Tsipras, D.; and Vladu, A. 2017. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*.

Phute, M.; Helbling, A.; Hull, M.; Peng, S.; Szyller, S.; Cornelius, C.; and Chau, D. H. 2023. Llm self defense: By self examination, llms know they are being tricked. *arXiv preprint arXiv:2308.07308*.

Russakovsky, O.; Deng, J.; Su, H.; Krause, J.; Satheesh, S.; Ma, S.; Huang, Z.; Karpathy, A.; Khosla, A.; Bernstein, M.; Berg, A. C.; and Fei-Fei, L. 2015. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3): 211–252.

Shapira, A.; Zolfi, A.; Demetrio, L.; Biggio, B.; and Shabtai, A. 2023. Phantom sponges: Exploiting non-maximum suppression to attack deep object detectors. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 4571–4580.

Shumailov, I.; Zhao, Y.; Bates, D.; Papernot, N.; Mullins, R.; and Anderson, R. 2021. Sponge examples: Energy-latency attacks on neural networks. In *2021 IEEE European symposium on security and privacy (EuroS&P)*, 212–231. IEEE.

Singh, S.; Romanou, A.; Fourrier, C.; Adelani, D. I.; Ngui, J. G.; Vila-Suero, D.; Limkonchotiwat, P.; Marchisio, K.; Leong, W. Q.; Susanto, Y.; Ng, R.; Longpre, S.; Ko, W.-Y.; Smith, M.; Bosselut, A.; Oh, A.; Martins, A. F. T.; Choshen, L.; Ippolito, D.; Ferrante, E.; Fadaee, M.; Ermis, B.; and Hooker, S. 2024. Global MMLU: Understanding and Addressing Cultural and Linguistic Biases in Multilingual Evaluation. arXiv:2412.03304.

Team, G.; Anil, R.; Borgeaud, S.; Alayrac, J.-B.; Yu, J.; Soricut, R.; Schalkwyk, J.; Dai, A. M.; Hauth, A.; Millican, K.; et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.

Team, Q. 2025. Qwen2.5-VL.

Wang, H.; Shi, H.; Tan, S.; Qin, W.; Wang, W.; Zhang, T.; Nambi, A.; Ganu, T.; and Wang, H. 2024a. Multimodal needle in a haystack: Benchmarking long-context capability of multimodal large language models. *arXiv preprint arXiv:2406.11230*.

Wang, W.; Lv, Q.; Yu, W.; Hong, W.; Qi, J.; Wang, Y.; Ji, J.; Yang, Z.; Zhao, L.; XiXuan, S.; et al. 2024b. Cogvlm: Visual expert for pretrained language models. *Advances in Neural Information Processing Systems*, 37: 121475–121499.

Xia, P.; Han, S.; Qiu, S.; Zhou, Y.; Wang, Z.; Zheng, W.; Chen, Z.; Cui, C.; Ding, M.; Li, L.; et al. 2024. Mmie: Massive multimodal interleaved comprehension benchmark for large vision-language models. *arXiv preprint arXiv:2410.10139*.

Ye, H.; Huang, D.-A.; Lu, Y.; Yu, Z.; Ping, W.; Tao, A.; Kautz, J.; Han, S.; Xu, D.; Molchanov, P.; et al. 2024. X-vila: Cross-modality alignment for large language model. *arXiv preprint arXiv:2405.19335*.

Zhang, Y.; Wang, X.; Gao, H.; Zhou, Z.; Meng, F.; Zhang, Y.; and Su, S. 2025. $PD^3F$: A Pluggable and Dynamic DoS-Defense Framework Against Resource Consumption Attacks Targeting Large Language Models. *arXiv preprint arXiv:2505.18680*.

Zhang, Y.; Zhou, Z.; Zhang, W.; Wang, X.; Jia, X.; Liu, Y.; and Su, S. 2024. Crabs: Consuming resource via auto-generation for llm-dos attack under black-box settings. *arXiv preprint arXiv:2412.13879*.

Zhou, C.; Li, Q.; Li, C.; Yu, J.; Liu, Y.; Wang, G.; Zhang, K.; Ji, C.; Yan, Q.; He, L.; et al. 2024. A comprehensive survey on pretrained foundation models: A history from bert to chatgpt. *International Journal of Machine Learning and Cybernetics*, 1–65.

Zhu, D.; Chen, J.; Shen, X.; Li, X.; and Elhoseiny, M. 2023. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*.

# Appendix

## A Defense Strategy Analysis

As shown in Table 10, our proposed defense method has almost no impact on normally generated requests, maintaining the fluency and integrity of natural output. However, excessive punishment can affect the quality of responses to normal questions, impacting the model's performance on general questions. We have provided mitigation measures for RCAs on LVLMs, but suppressing repetitive lengths may have semantic impacts on normal output. Therefore, further exploration is needed in future work to balance between security and usefulness.

|         | BLIP7B | Llava7B | Qwen3B |
|---------|--------|---------|--------|
| Normal  | 36.22  | 41.84   | 49.82  |
| Defence | 36.24  | 40.10   | 48.54  |

Table 10: Output length stability for benign requests under defense mechanisms.

## B Trend in Logit Changes for Top-k Tokens

Under our attack mechanism, we observed an interesting phenomenon: when the model is induced to repeatedly output a specific token (e.g., "flowers" in Figure 7), the logit values of semantically and morphologically highly related variants (such as 'flower' and "flow") are significantly increased and frequently appear in the Top-k candidate token set during the sampling process (Figure 7 shows the Top-5 candidate token set).
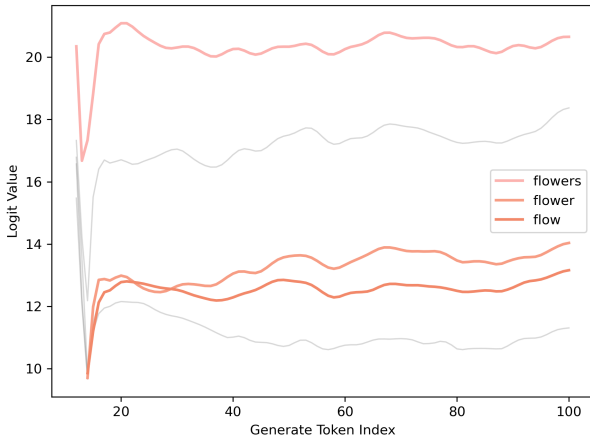


Figure 7: The trend of logit values for the Top-5 tokens.

This phenomenon causes the frequency penalty to fail. Even when frequency penalties are applied to reduce the Logit values of repeated tokens, the logit values of their variants increase. This means that when generating the next token, although the priority of "flowers" may be reduced, the attack mechanism has already highly focused the model's attention on the concept of "flower", causing the model to select other tokens from the Top-k candidate preferentially

set that express the same semantic meaning but have slightly different forms. This phenomenon indicates that the attack does not simply force the model to repeat a single token but instead induces it to become trapped in a semantic loop, causing variants related to the target concept to dominate the logit distribution in the next generation step.

## C Human Evaluation Settings

We have constructed three types of problems based on visual consistency, feature similarity, and semantic consistency to evaluate the covertness of attacks, where each type of problem provides an original image and an attack image (or white noise image, compressed image). The three types of problems are: "There is no significant difference between image 1 and image 2", "Image 1 and image 2 have similarities in visual features", and "Image 1 and image 2 have the same core meaning". For each question, we provide 5 options (completely disagree, somewhat disagree, uncertain, somewhat agree, completely agree) corresponding to 1-5 points, meaning that the higher the score, the better the covertness.

## D More Experimental Results for RECALLED

Figure 8 shows the attack results of RECALLED on three models, where the first column displays the token-level attack results and the second column displays the sentence-level attack results. As shown in the figure, for the three models, the attack images generated by RECALLED can effectively trigger **Output Recall**, with no significant difference between the attack image and the original image.

Figure 9 shows the attack results of multi-target RECALLED on three models, where three images of the same classes use the same perturbation. As shown in the figure, RECALLED supports multi-objective optimization and can generate effective attack images for the three models.

## E Limitations

Due to potential security risks and ethical considerations, we conducted all experiments in controlled environments without deploying attacks against production systems. Following responsible disclosure practices, we reported our findings to the model manufacturers upon completion of our research. Additionally, we propose practical mitigation strategies to address the identified vulnerabilities.
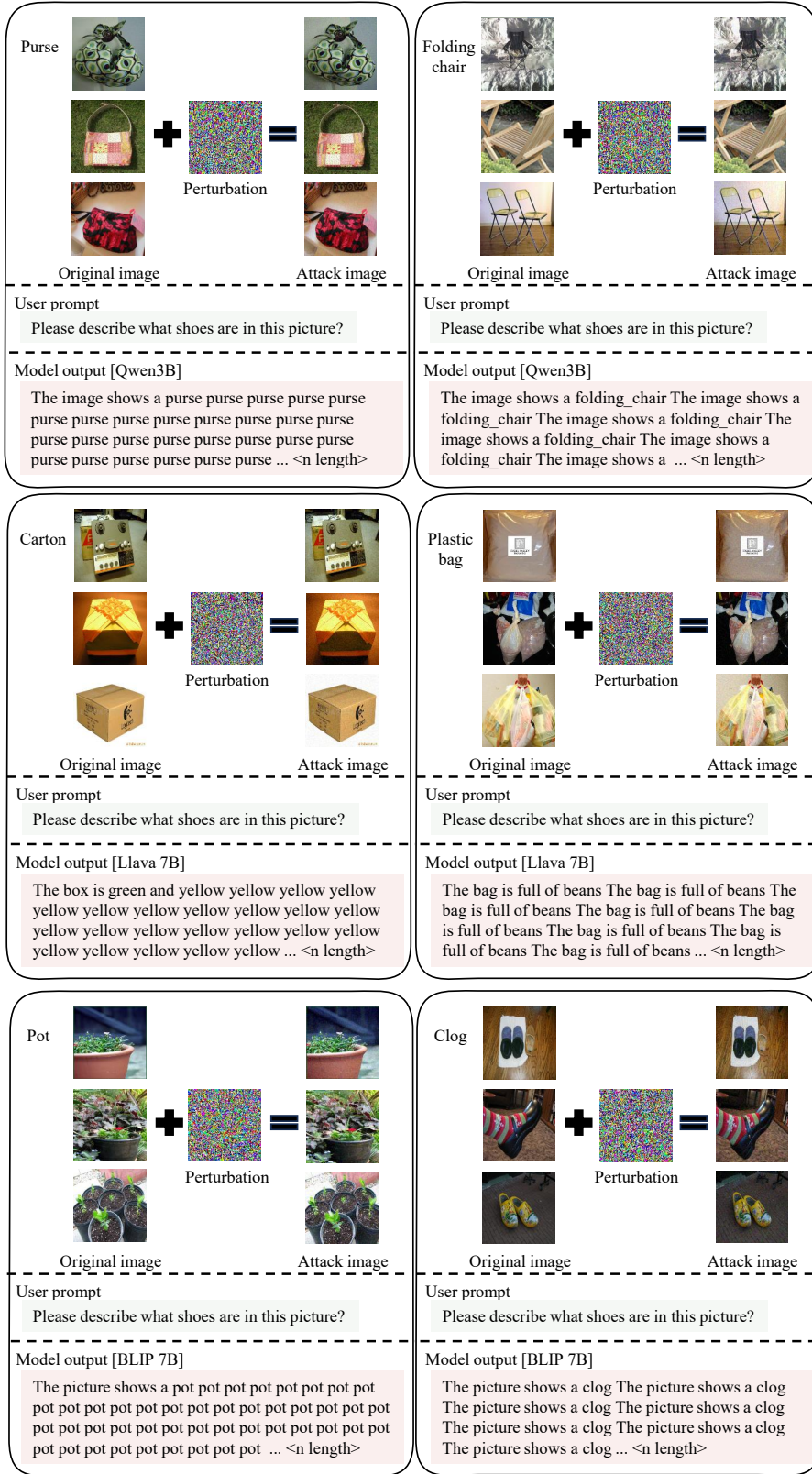
Figure 8: Example for RECALLED attack.

Figure 9: Example for multi-objective RECALLED attack.