

RESOURCE CONSUMPTION RED-TEAMING FOR LARGE VISION-LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

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 mainly overlooked visual inputs as a potential attack surface, resulting in insufficient mitigation strategies against RCAs in LVLMs. To address this gap, we propose RECITE (**R**esource **C**onsumption **R**ed-**T**eaming for LVLMs), 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 Objective* adversarial perturbations, which can induce repeating output. Then, we inject the perturbations into visual inputs, triggering unbounded generations to achieve the goal of RCAs. Empirical results demonstrate that RECITE increases service response latency by over $26\times\uparrow$, resulting in an additional 20% increase in GPU utilization and memory consumption. Our study reveals security vulnerabilities in LVLMs and establishes a red-teaming framework that can facilitate the development of future defenses against RCAs.

1 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 applications. Shapira et al. (2023); Krithivasan et al. (2020).

The Sponge sample is an RCAs designed for computer vision models that disrupts visual attention mechanisms Shumailov et al. (2021), resulting in extra resource consumption. Since visual input can 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 for LVLMs Zhang et al. (2025b), despite the inherent vulnerability of the visual modality.

To address this, we investigate effective red-teaming methodologies for RCAs exploiting visual inputs. We propose RECITE, an **R**esource **C**onsumption **R**ed-**T**eaming for Large Vision-Language Models. RECITE employs *Vision Guided Optimization* to craft adversarial perturbations through fine-grained optimization targeting *Output Recall Objective*, which is designed to trigger unbounded output repetition. Then, we inject perturbations into visual inputs, covertly manipulating the model responses to achieve RCA objectives. Leveraging RECITE to induce unbounded generations, we reveal the menace of adversarial visual patterns and assess the vulnerability of frontier 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 RECITE. Our work indicates that adversarial visual inputs can induce severe RCAs in LVLMs, even

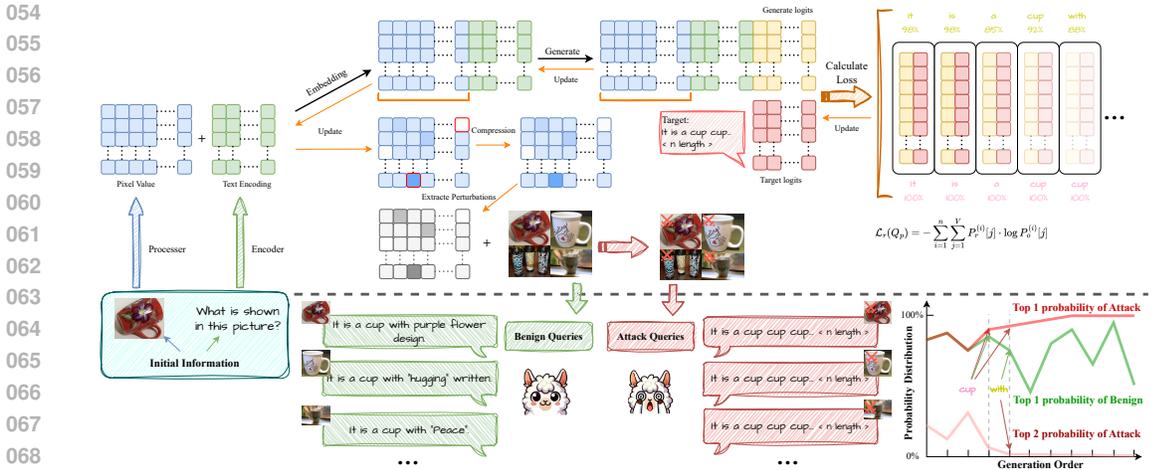


Figure 1: The RECITE pipeline. In the generation stage, we employ a gradient-based method to iteratively update the visual input. In the evaluation stage, the constructed RECITE triggers unbounded loop generation in the target model.

when paired with benign textual prompts. We introduce RECITE, a red-teaming to test this vulnerability, and show that it induces a $26\times\uparrow$ increase in output length, consistently forcing generation to saturate the model’s maximum context window. This extension leads to at least 20% degradation in service latency for LVLm applications. Beyond validating the attack’s efficacy, we leverage RECITE as a diagnostic framework. By analyzing the failure of LVLms to suppress RCAs through our *Output Recall Objective*, we uncover the root mechanisms of this vulnerability, thereby motivating our proposed mitigation strategies.

In summary, our primary contribution lies in RECITE, a red-teaming methodology that leverages vision-based perturbations to evaluate service degradation and potential system crashes in LVLms. We then conduct extensive experiments to validate the effectiveness of RECITE and demonstrate the vulnerability of visual input processing in LVLms. Finally, we provide a comprehensive analysis of the vulnerabilities in LVLms exposed by RCAs and reveal the reason why mitigating the resource consumption is an inherently challenging task. Our findings highlight the shortcomings of LVLms in addressing threats to visual RCAs, underscoring the need for more robust defense mechanisms.

2 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. (2024a); Wang et al. (2025); Zhang et al. (2025a); Liu et al. (2024b) that enable 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. (2025b); 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.

3 METHOD

In this section, we construct RECITE, which is an unbounded RCA for vision inputs. As shown in the Figure 1, we first establish *Output Recall* as the red-teaming target for resource consumption to guide the optimization process of RECITE. Then, we introduce *Vision Guided Optimization*, a perturbation injection method for vision input, to achieve effective RCAs.

3.1 CONSTRUCTING OUTPUT RECALL OBJECTIVE

We introduce the *Output Recall Objective*, a mechanism designed to induce unbounded, auto-regressive generation in LVLMS. This objective guides the model into a repetitive generation mode, generating with a fixed format indefinitely.

We define the benign image-to-text generation task as $Q : (I, T_q) \rightarrow T_a$, where I denotes the visual input, $T_q = \{q_1, q_2, \dots, q_m\}$ represents the tokenized textual question, and $T_a = \{a_1, a_2, \dots, a_n\}$ denotes the token sequence generated as the answer. The constituent tokens of both sequences belong to the model’s vocabulary \mathcal{V} , such that $q_i \in \mathcal{V}$ for all $i \in \{1, \dots, m\}$ and $a_j \in \mathcal{V}$ for all $j \in \{1, \dots, n\}$. We investigate the effect of prefix information in the T_a on subsequent model behavior, and accordingly define the *Initial Output Recall* target as:

$$R_0 = \{a_1, a_2, \dots, a_k\}, \text{ subject to } k \in \{1, 2, \dots, n - 1\}. \quad (1)$$

In our experiments, we primarily set the position of the first punctuation mark as $k+1$. We construct *Output Recall Objective* using a loop mechanism and introduce *Repeating Parameter* $\rho \in \mathbb{N}$ to control the intensity of attack. Two types of *Output Recall Objective* construction are defined:

1. Token-Level Output Recall Objective: Let $G = (a_{k-l+1}, \dots, a_k)$ be the token sub-sequence corresponding to the final word in the generated sequence $R_0 = (a_1, \dots, a_k)$, where l is the number of tokens comprising this word. The token-level Recall is then constructed as:

$$R_\rho^t = R_0 \underbrace{\|G\|G\|\dots\|G\|}_{\rho \text{ times}}, \quad (2)$$

where $\|$ denotes sequence concatenation.

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

$$R_\rho^s = R_0 \underbrace{\|R_0\|\dots\|R_0\|}_{\rho \text{ times}}. \quad (3)$$

Both R_ρ^t and R_ρ^s can be viewed as extended forms of *Output Recall Objective*, where $R_\rho \in R = \bigcup_{\rho \in \mathbb{N}} \{R_\rho^t, R_\rho^s\}$. R_ρ serves as the target for recursive optimization, inducing LVLMS to enter a potential non-terminating generation state. Further explanations can be found in Appendix E.

3.2 VISION GUIDED OPTIMIZATION

To optimize the *Output Recall Objective*, we propose *Vision Guided Optimization*, an efficient gradient-based optimization method. Following GCG Liao & Sun (2024), we adopt its loss formulation to optimize a controllable perturbation applied to the input image directly. This approach facilitates rapid convergence toward the target objective.

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

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

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 $E_p = \text{VisualEmbed}(Q_p)$, $E_p \in \mathbb{R}^{L_p \times d}$. For question T_p , we first tokenize it into a sequence $Q_t = \text{Tokenizer}(T_p) = \{t_1, t_2, \dots, t_m\}$, $Q_t \in \mathbb{Z}^m$. 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:

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

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

The generative model’s mapping function from input to output vectors is defined as $e_o^{(i)} = \text{Generate}(E^{(i)})$, $i = 0, 1, \dots, n$. 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 following process:

$$\begin{aligned} E^{(i+1)} &= E^{(i)} || e_a^{i+1}, \\ e_o^{(i+1)} &= \text{Generate}(E^{(i+1)}), \quad i = 1, \dots, n-1. \end{aligned} \quad (6)$$

This yields an output embedding sequence $E_o = \{e_o^{(1)}, e_o^{(2)}, \dots, e_o^{(n)}\} \in \mathbb{R}^{n \times d}$.

We utilize cross-entropy loss on the 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:

$$\mathcal{L}_r(Q_p) = \sum_{i=1}^n \text{CE}(P_o^{(i)}, P_r^{(i)}) = - \sum_{i=1}^n \sum_{j=1}^V P_r^{(i)}[j] \cdot \log P_o^{(i)}[j]. \quad (7)$$

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

We depart from the suffix textual attack of GCG Jia et al. (2024) by defining our perturbation space over the input image itself. This allows us to employ gradient descent to minimize $\mathcal{L}_r(Q_p)$. Given the original image input I , we introduce a perturbation δ in pixel space. The optimization problem is formulated as:

$$\begin{aligned} \min_{\delta} \mathcal{L}_r(\tilde{Q}_p) &= \min_{\delta} \sum_{i=1}^n \text{CE}(P_o^{(i)}, P_r^{(i)}) \\ \text{s.t. } \tilde{Q}_p &\in [-1.0, 1.0]^d, \tilde{Q}_p = Q_p + \delta, \|\delta\|_{\infty} \leq \epsilon, \end{aligned} \quad (8)$$

here, ϵ denotes the perturbation budget, which governs the visual perceptibility. By constraining ϵ to a small value, the resulting adversarial perturbation becomes effectively imperceptible, thereby enabling RECITE to evade detection systems that rely on anomaly detection.

We then apply K -step optimization to update $\mathcal{L}_r(\tilde{Q}_p)$, computing gradients of visual inputs. Upon completion of the perturbation optimization, we apply an inverse reconstruction function to generate the adversarial image $\tilde{I} = \text{Reprocessor}(\tilde{Q}_p) = \text{Reprocessor}(Q_p + \delta^*)$, δ^* represents the optimal perturbation obtained after *Vision Guided Optimization*.

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

4 EXPERIMENT OF RECITE

4.1 EXPERIMENTAL SETUP

Models. We conduct 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) 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.

Baselines. In covertness experiments, we evaluate against two categories of defense mechanisms: perplexity-based detection methods (PPL) Jain et al. (2023); Alon & Kamfonas (2023) and input self-monitoring (ISM) Phute et al. (2023). We define attack success as achieving a success rate exceeding 80% or higher. For baseline comparisons, we evaluate against 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.

4.2 PERFORMANCE OF THE RECITE RED-TEAMING METHOD

RECITE Effectiveness Analysis.

To verify the efficiency of red-teaming samples generation, we use a truncation verification mechanism to verify RECITE samples. Specifically, we limit the generation length to 500 tokens and evaluate the generated samples after 1,000 rounds of optimization. The results are presented in Table 1. 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 RECITE requests. Table 2 demonstrates that red-teaming samples significantly increase output length. The average output length of RECITE exceeds 1,900 tokens, resulting in a substantial $26\times\uparrow$ increase. Moreover, a significant proportion of RECITE achieves the model’s maximum output window, exhibiting unbounded output behavior. RECITE systematically induces models to generate unbounded content, triggering infinite generation behavior. RECITE achieves uninterrupted response generation for the first time, which has not been stably achieved in previous studies.

Resource Consumption Simulation.

We conduct experiments on NVIDIA A4000 GPUs to evaluate the impact of RECITE on commercially deployed models. The experimental results are presented in Table 3 and Figure 2. Compared to benign image-to-text requests, RECITE samples cause over $54\times\uparrow$ 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 RECITE not only extend model generation but also pose significant threats to underlying computational resources, severely compromising system reliability and model robustness in production environments. We tested the universality of RECITE in Appendix G.

RECITE Time Consumption. We evaluate the optimization time for RCAs using GCG-RCAs and RECITE. As shown in Table 4, the optimization time for both methods increases with model size, reflecting the greater computational complexity.

Table 1: Generation success rates for RECITE samples.

Type	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 2: Average generation length comparison between RECITE attacks and benign queries.

Type	Qwen			Llava		BLIP	
	3B	7B	32B	7B	13B	7B	13B
Benign	63	91	197	40	34	39	38
Token	2023	2046	2022	2048	2021	2030	2048
Sentence	2048	2048	1961	2048	1427	2048	2046

Table 3: Comparison of performance consumption.

Model	Method	Output Time	GPU Utilization	Memory Usage
Qwen3B	Benign	2.82	47.52%	49.25%
	RECITE	87.56 ^(↑84.74)	57.48% ^(↑9.96%)	49.67% ^(↑0.42%)
Llava7B	Benign	1.75	93.30%	87.30%
	RECITE	96.08 ^(↑94.33)	97.93% ^(↑4.63%)	96.01% ^(↑8.71%)
BLIP7B	Benign	190.95	93.08%	86.07%
	RECITE	1154.76 ^(↑963.81)	96.50% ^(↑3.42%)	95.72% ^(↑9.65%)

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

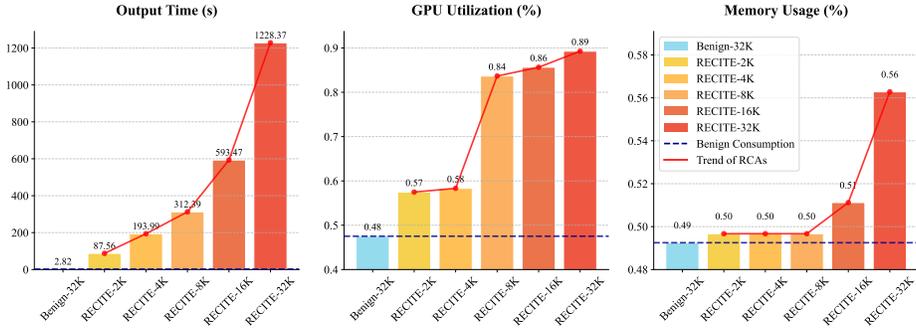


Figure 2: Resource cost of different output length limits.

RECITE achieves at least 100 times faster optimization compared to GCG-RCAs for the same target model. This efficiency is attributed to performing a direct optimization in the continuous space to find the attack target. In contrast, GCG-RCAs operate in a discrete space, requiring a costly search process that involves iteratively evaluating and replacing candidate tokens, resulting in slower convergence.

Table 4: Comparison of attack time between GCG-RCAs and RECITE.

Method	Qwen		Meta-Llama	
	3B	7B	7B	13B
GCG-RCAs	1759.80s	3312.00s	3508.21s	6327.25s
RECITE	15.14s	33.58s	16.22s	61.12s

4.3 COVERTNESS OF RECITE

Subjective Evaluation Results. To evaluate the perceptual concealment of RECITE, we design a five-point Likert scale questionnaire Joshi et al. (2015) across three dimensions: visual consistency, feature similarity, and semantic consistency. We recruit 40 participants to participate in the study, using white noise (Noise) and image compression (Compression) as comparison baselines. The results are presented in Figure 3. RECITE achieves significantly higher average scores across all three dimensions compared to both baselines. Additionally, we conduct 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 3 right, only 6.88% of RECITE samples are identified as “harmful”, further demonstrating its effective perceptual evasion capabilities. More settings are provided in Appendix H.

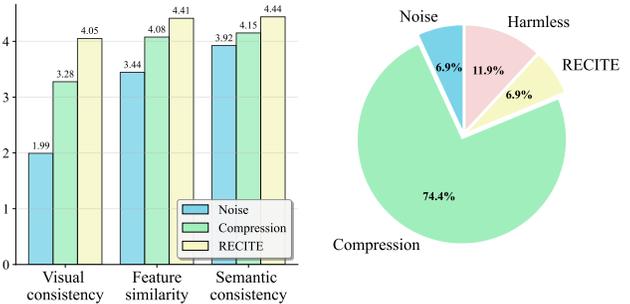


Figure 3: RECITE performance evaluation. Left: Likert scale ratings across three evaluation dimensions. Right: Harmfulness assessment results (Harmless indicates all generated images are deemed non-harmful). RECITE achieves superior performance in both metrics.

Quantitative Evaluation Results. To evaluate the detectability of RECITE at the input level, we utilize PPL as a metric for language naturalness assessment. As shown in Table 5, RECITE samples exhibit lower average PPL than benign requests, indicating concealment in terms of language fluency. Unlike conventional RCA examples that often exhibit semantic anomalies, RECITE constructs perturbations solely in the visual domain while preserving the distributional characteristics of natural language, thus avoiding evident traces of malicious construction.

Table 5: RECITE request quality assessment via language model perplexity.

Model	Benign	GCG-RCAs	RECITE
Meta-Llama	202.36	5103.98	42.77
Qwen	200.67	17212.22	40.47

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

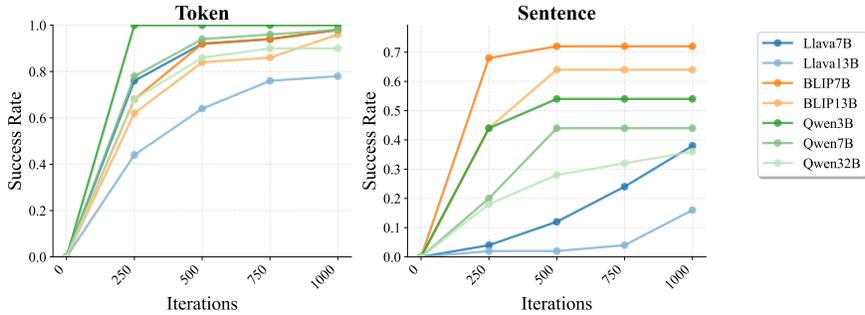


Figure 4: Success rate convergence analysis across generation iterations in RECITE.

Furthermore, we adopt the LLM-as-a-Judge framework to determine whether input prompts contain potential attack intent. In Table 6, RECITE samples consistently evaded detection by the LLM-based discriminator. These results demonstrate that RECITE can effectively bypass automated detection mechanisms, further highlighting the security challenges it poses to models.

Table 6: Attack effectiveness evaluation using LLM-as-a-Judge assessment with 80% recognition threshold.

Model	GCG-RCA	RECITE
Meta-Llama	✓	✗
Qwen	✓	✗

4.4 ABLATION ANALYSIS

To evaluate the robustness and stability of the RECITE method with respect to hyperparameters, we conduct ablation studies on the number of iterations and the repeating parameter ρ . As shown in Figure 4, the attack success rates on different models exhibit an upward trend with increasing iterations before stabilizing. This demonstrates that RECITE 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 7. When ρ 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 ρ is set to 10, the optimization complexity of the attack objective increases, resulting in degraded attack success rates. These results demonstrate that RECITE exhibits stable performance within reasonable hyperparameter ranges, with an optimal operating interval. The Limitations of RECITE are shown in Appendix A.

Table 7: Impact of ρ on RECITE attack success rates.

Model	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%

5 EFFECTIVENESS ASSESSMENT OF RECITE

In this section, we employ the RECITE to analyze the vulnerabilities in LVLMs exposed by RCA. Our investigation yields three key findings. First, we identify the cause of this fragility, demonstrating that it stems from the induction of the Output Recall Objective. Second, we reveal that the resulting RCA patterns are remarkably stable and possess a recurrent structure that is highly resistant to disruption. Third, through an analysis of the LVLm’s predictive probabilities, we show RCAs are self-reinforcing and thus intractable to mitigate using the model’s intrinsic capabilities alone. These findings collectively establish the severity and nature of the vulnerability, motivating our proposal of an effective mitigation strategy.

5.1 OUTPUT RECALL INDUCES UNBOUNDED GENERATION

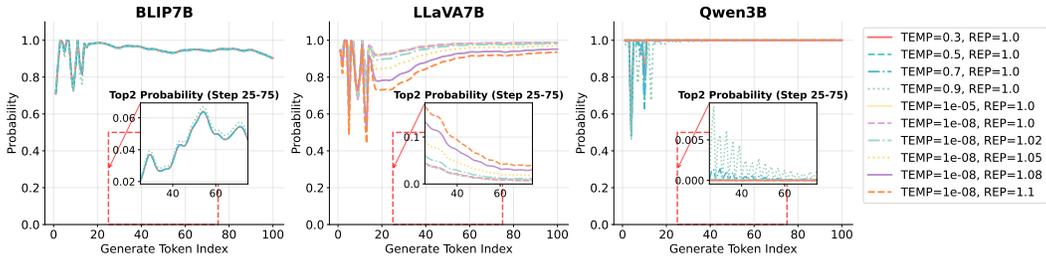


Figure 5: Impact of temperature (TEMP) and repetition penalty (REP) hyperparameters on generation length and semantic repetition in RECITE.

In RECITE, we construct Token-Level Output Recall Objective R_ρ^t and Sentence-Level Output Recall Objective R_ρ^s . R_ρ^t generates shorter content, facilitates model entry into the loop generation, and provides a more stable attack target. In contrast, R_ρ^s provides richer semantic information and induces the model to generate more coherent. We directly concatenate the *Output Recall Objective* to the original request text T_q to form a new input $T'_q = T_q || R$. The results for T'_q in Table 8 demonstrate that as ρ increases, the generated content exhibits stronger consistency and repetitiveness. *Output Recall Objective* significantly disrupts the natural response structure of the original text, causing the model to preferentially continue generating content similar to the *Output Recall Objective* 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 Objective* method systematically induces this phenomenon through precise construction. Consequently, an adversary can exploit this objective to mount a severe RCA. The *Vision Guided Optimization* used by RECITE is one of them.

5.2 RECITE GENERATION STABILITY

Building on the results from RECITE, we conduct a deeper analysis of the RCA’s stability. Given our finding that the *Output Recall Objective* consistently induces verbose RCAs, we design a targeted experiment to probe the long-range persistence of this red-teaming method. Rather than retry the full RECITE benchmark, we select the sample that can generate 500 tokens. By prompting the model to generate up to its maximum context length for this specific case, we can evaluate whether the repetitive loop is self-sustaining or eventually decays. The experimental results are shown in Table 9. 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 behavior reveals a critical failure mode where the model cannot terminate the RCA, creating a Denial of Service (DoS) vulnerability through computational resource exhaustion.

Table 8: Repeating Parameters ρ influence on infinite generation success rates under direct splicing target scenarios.

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 9: Short-length check available analysis.

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%

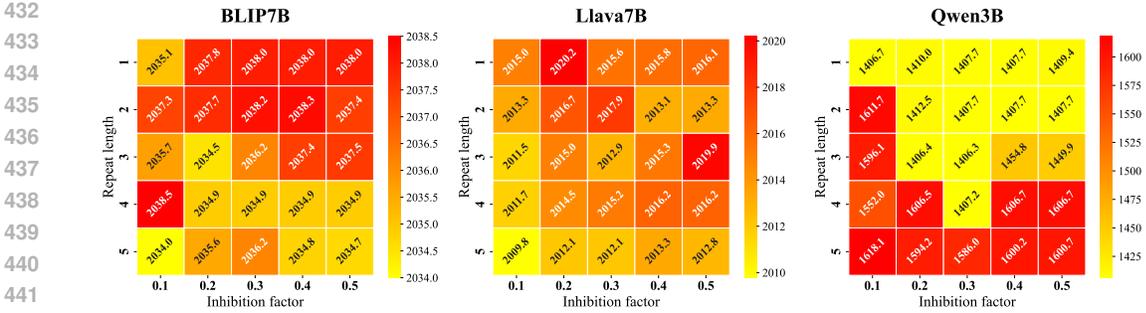


Figure 6: Length reduction performance of the proposed defense strategy across parameter configurations.

5.3 PREDICTION TENDENCIES ANALYSIS

To diagnose the cause of this unbounded generation, we analyze the model’s predictive behavior. Specifically, we vary the temperature and repetition penalty to assess whether these standard control mechanisms can disrupt the prediction distributions sustained by RCA. We set the attack target as the RECITE construction with $\rho = 5$. As illustrated in Figure 5, the attack target consistently maintains the Top-1 probability at each generation step. As the 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 increase the sampling probability of alternative tokens when dominated by 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.

5.4 EXPLORATION OF DEFENSIVE MEASURES

Given the high threat and covertness of RECITE, effective defense mechanisms are essential to mitigate associated risks, yet relevant research remains limited. Based on the core mechanism of RECITE, we propose a general defense method that dynamically adjusts the probability distribution of output, 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:

$$l' = l \times (1 + \alpha \times f_{max}), \tag{9}$$

where α is the scaling factor, this mechanism can substitute the standard repetition penalty strategy while achieving dynamic suppression of RCAs.

As illustrated in Figure 6, 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, it employs aggressive penalty schemes that may adversely affect legitimate queries, leaving room for future optimization. Appendix B provides additional analysis.

6 CONCLUSION

We present RECITE, a novel red-teaming methodology for RCAs targeting LVLMS. RECITE leverages Output Recall mechanisms to induce repetitive generation patterns. We then introduces Vision Guided Optimization Loss to construct attack templates. We validate RECITE’s effectiveness across seven state-of-the-art LVLMS, achieving consistently high attack success rates. Furthermore, through systematic output tendency analysis, we provide theoretical insights into the underlying causes of RCAs in LVLMS, revealing why mitigating such attacks is an inherently challenging problem. Our work exposes a critical yet underexplored vulnerability in LVLMS security, highlighting the susceptibility of vision inputs to resource exhaustion attacks.

ETHICS STATEMENT

The research presented in this paper does not involve human subjects. The experiments were conducted on publicly available datasets, such as ImageNet, which do not contain personally identifiable information and thus raise no privacy concerns.

We acknowledge that, like many methods in machine learning, the techniques for inducing model failures could potentially be misused. We have included a dedicated discussion on the potential for such dual-use applications and broader societal impacts in Sec 5 and Appendix A. Our analysis focuses on the fundamental mechanisms of model behavior, and as such, does not directly engage with downstream tasks or datasets where issues of fairness or societal bias are primary concerns.

REPRODUCIBILITY STATEMENT

To facilitate the reproducibility of our findings, we provide comprehensive details of our methodology and experimental setup. To construct the RECITE benchmark, including the data curation process and the instantiation of our *Output Recall Objective*, is detailed in Appendix F. All hyperparameters governing the optimization, such as the repetition factor ρ and the perturbation budget ϵ , are explicitly documented within our experimental analysis in Sec 5. Furthermore, to enable full verification and to encourage future work, we have included our complete source code, with scripts to replicate the main experimental results, in the supplementary material.

REFERENCES

- Gabriel Alon and Michael Kamfonas. Detecting language model attacks with perplexity. *arXiv preprint arXiv:2308.14132*, 2023.
- Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- Gongwei Chen, Leyang Shen, Rui Shao, Xiang Deng, and Liqiang Nie. Lion: Empowering multimodal large language model with dual-level visual knowledge. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 26540–26550, 2024.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code, 2021.
- Simin Chen, Zihe Song, Mirazul Haque, Cong Liu, and Wei Yang. Nicgslowdown: Evaluating the efficiency robustness of neural image caption generation models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15365–15374, 2022.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. *Advances in neural information processing systems*, 36:49250–49267, 2023.

- 540 Jianshuo Dong, Ziyuan Zhang, Qingjie Zhang, Tianwei Zhang, Hao Wang, Hewu Li, Qi Li, Chao
541 Zhang, Ke Xu, and Han Qiu. An engorgio prompt makes large language model babble on. *arXiv*
542 *preprint arXiv:2412.19394*, 2024.
- 543
- 544 Jiyuan Fu, Kaixun Jiang, Lingyi Hong, Jinglun Li, Haijing Guo, Dingkang Yang, Zhaoyu Chen, and
545 Wenqiang Zhang. Lingoloop attack: Trapping mllms via linguistic context and state entrapment
546 into endless loops. *arXiv preprint arXiv:2506.14493*, 2025.
- 547
- 548 Kuofeng Gao, Yang Bai, Jindong Gu, Shu-Tao Xia, Philip Torr, Zhifeng Li, and Wei Liu. Induc-
549 ing high energy-latency of large vision-language models with verbose images. In *International*
550 *Conference on Learning Representations*, 2024a.
- 551
- 552 Kuofeng Gao, Tianyu Pang, Chao Du, Yong Yang, Shu-Tao Xia, and Min Lin. Denial-of-service
553 poisoning attacks against large language models. *arXiv preprint arXiv:2410.10760*, 2024b.
- 554
- 555 Jonas Geiping, Alex Stein, Manli Shu, Khalid Saifullah, Yuxin Wen, and Tom Goldstein. Coercing
556 llms to do and reveal (almost) anything. *arXiv preprint arXiv:2402.14020*, 2024.
- 557
- 558 Jiaming Han, Renrui Zhang, Wenqi Shao, Peng Gao, Peng Xu, Han Xiao, Kaipeng Zhang, Chris Liu,
559 Song Wen, Ziyu Guo, et al. Imagebind-llm: Multi-modality instruction tuning. *arXiv preprint*
560 *arXiv:2309.03905*, 2023.
- 561
- 562 Mirazul Haque, Simin Chen, Wasif Haque, Cong Liu, and Wei Yang. Antinode: Evaluating effi-
563 ciency robustness of neural odes. In *Proceedings of the IEEE/CVF International Conference on*
564 *Computer Vision*, pp. 1507–1517, 2023.
- 565
- 566 Sanghyun Hong, Yiğitcan Kaya, Ionuț-Vlad Modoranu, and Tudor Dumitraș. A panda? no, it’s
567 a sloth: Slowdown attacks on adaptive multi-exit neural network inference. *arXiv preprint*
568 *arXiv:2010.02432*, 2020.
- 569
- 570 Zhang-wei Hong, Idan Shenfeld, Tsun-hsuan Wang, Yung-sung Chuang, Aldo Pareja, James Glass,
571 Akash Srivastava, and Pulkit Agrawal. Curiosity-driven red-teaming for large language models.
572 In *International Conference on Learning Representations*, 2024.
- 573
- 574 Neel Jain, Avi Schwarzschild, Yuxin Wen, Gowthami Somepalli, John Kirchenbauer, Ping-yeh Chi-
575 ang, Micah Goldblum, Aniruddha Saha, Jonas Geiping, and Tom Goldstein. Baseline defenses
576 for adversarial attacks against aligned language models. *arXiv preprint arXiv:2309.00614*, 2023.
- 577
- 578 Xiaojun Jia, Tianyu Pang, Chao Du, Yihao Huang, Jindong Gu, Yang Liu, Xiaochun Cao, and Min
579 Lin. Improved techniques for optimization-based jailbreaking on large language models. *arXiv*
580 *preprint arXiv:2405.21018*, 2024.
- 581
- 582 Ankur Joshi, Saket Kale, Satish Chandel, and D Kumar Pal. Likert scale: Explored and explained.
583 *British journal of applied science & technology*, 7(4):396, 2015.
- 584
- 585 Sarada Krithivasan, Sanchari Sen, and Anand Raghunathan. Sparsity turns adversarial: Energy
586 and latency attacks on deep neural networks. *IEEE Transactions on Computer-Aided Design of*
587 *Integrated Circuits and Systems*, 39(11):4129–4141, 2020.
- 588
- 589 Sarada Krithivasan, Sanchari Sen, Nitin Rathi, Kaushik Roy, and Anand Raghunathan. Efficiency
590 attacks on spiking neural networks. In *Proceedings of the 59th ACM/IEEE Design Automation*
591 *Conference*, pp. 373–378, 2022.
- 592
- 593 Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan
Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv preprint*
arXiv:2408.03326, 2024.
- Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Nau-
mann, Hoifung Poon, and Jianfeng Gao. Llava-med: Training a large language-and-vision as-
sistant for biomedicine in one day. *Advances in Neural Information Processing Systems*, 36:
28541–28564, 2023.

- 594 Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-
595 training for unified vision-language understanding and generation. In *International conference on*
596 *machine learning*, pp. 12888–12900. PMLR, 2022.
- 597 Zeyi Liao and Huan Sun. Amplegcg: Learning a universal and transferable generative model of
598 adversarial suffixes for jailbreaking both open and closed llms. *arXiv preprint arXiv:2404.07921*,
599 2024.
- 600 Bin Lin, Yang Ye, Bin Zhu, Jiayi Cui, Munan Ning, Peng Jin, and Li Yuan. Video-llava: Learning
601 united visual representation by alignment before projection. *arXiv preprint arXiv:2311.10122*,
602 2023.
- 603 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances*
604 *in Neural Information Processing Systems*, 36:34892–34916, 2023.
- 605 Xin Liu, Yichen Zhu, Jindong Gu, Yunshi Lan, Chao Yang, and Yu Qiao. Mm-safetybench: A
606 benchmark for safety evaluation of multimodal large language models. In *European Conference*
607 *on Computer Vision*, pp. 386–403, 2024a.
- 608 Yi Liu, Chengjun Cai, Xiaoli Zhang, Xingliang Yuan, and Cong Wang. Arondight: Red teaming
609 large vision language models with auto-generated multi-modal jailbreak prompts. In *Proceedings*
610 *of the 32nd ACM International Conference on Multimedia*, pp. 3578–3586, 2024b.
- 611 Mansi Phute, Alec Helbling, Matthew Hull, ShengYun Peng, Sebastian Szyller, Cory Cornelius, and
612 Duen Horng Chau. Llm self defense: By self examination, llms know they are being tricked.
613 *arXiv preprint arXiv:2308.07308*, 2023.
- 614 Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng
615 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei.
616 ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*
617 (*IJCV*), 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
- 618 Avishag Shapira, Alon Zolfi, Luca Demetrio, Battista Biggio, and Asaf Shabtai. Phantom sponges:
619 Exploiting non-maximum suppression to attack deep object detectors. In *Proceedings of the*
620 *IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 4571–4580, 2023.
- 621 Iliia Shumailov, Yiren Zhao, Daniel Bates, Nicolas Papernot, Robert Mullins, and Ross Anderson.
622 Sponge examples: Energy-latency attacks on neural networks. In *2021 IEEE European symposium*
623 *on security and privacy (EuroS&P)*, pp. 212–231. IEEE, 2021.
- 624 Shivalika Singh, Angelika Romanou, Cl  mentine Fourrier, David I. Adelani, Jian Gang Ngui, Daniel
625 Vila-Suero, Peerat Limkonchotiwat, Kelly Marchisio, Wei Qi Leong, Yosephine Susanto, Ray-
626 mond Ng, Shayne Longpre, Wei-Yin Ko, Madeline Smith, Antoine Bosselut, Alice Oh, Andre
627 F. T. Martins, Leshem Choshen, Daphne Ippolito, Enzo Ferrante, Marzieh Fadaee, Beyza Ermiş,
628 and Sara Hooker. Global mmlu: Understanding and addressing cultural and linguistic biases in
629 multilingual evaluation, 2024. URL <https://arxiv.org/abs/2412.03304>.
- 630 Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut,
631 Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly
632 capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- 633 Qwen Team. Qwen2.5-vl, January 2025. URL <https://qwenlm.github.io/blog/qwen2.5-vl/>.
- 634 Hengyi Wang, Haizhou Shi, Shiwei Tan, Weiyi Qin, Wenyuan Wang, Tunyu Zhang, Akshay Nambi,
635 Tanuja Ganu, and Hao Wang. Multimodal needle in a haystack: Benchmarking long-context
636 capability of multimodal large language models. *arXiv preprint arXiv:2406.11230*, 2024a.
- 637 Weihang Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang,
638 Lei Zhao, Song XiXuan, et al. Cogvlm: Visual expert for pretrained language models. *Advances*
639 *in Neural Information Processing Systems*, 37:121475–121499, 2024b.
- 640 Xiaosen Wang, Shaokang Wang, Zhijin Ge, Yuyang Luo, and Shudong Zhang. Attention! you vision
641 language model could be maliciously manipulated. *arXiv preprint arXiv:2505.19911*, 2025.

648 Peng Xia, Siwei Han, Shi Qiu, Yiyang Zhou, Zhaoyang Wang, Wenhao Zheng, Zhaorun Chen,
649 Chenhang Cui, Mingyu Ding, Linjie Li, et al. Mmie: Massive multimodal interleaved compre-
650 hension benchmark for large vision-language models. *arXiv preprint arXiv:2410.10139*, 2024.
651

652 Hanrong Ye, De-An Huang, Yao Lu, Zhiding Yu, Wei Ping, Andrew Tao, Jan Kautz, Song Han, Dan
653 Xu, Pavlo Molchanov, et al. X-vila: Cross-modality alignment for large language model. *arXiv*
654 *preprint arXiv:2405.19335*, 2024.

655 Jiaming Zhang, Rui Hu, Qing Guo, and Wei Yang Bryan Lim. Cavalry-v: A large-scale generator
656 framework for adversarial attacks on video mllms. *arXiv preprint arXiv:2507.00817*, 2025a.
657

658 Yuanhe Zhang, Zhenhong Zhou, Wei Zhang, Xinyue Wang, Xiaojun Jia, Yang Liu, and Sen Su.
659 Crabs: Consuming resource via auto-generation for llm-dos attack under black-box settings. *arXiv*
660 *preprint arXiv:2412.13879*, 2024.

661 Yuanhe Zhang, Xinyue Wang, Haoran Gao, Zhenhong Zhou, Fanyu Meng, Yuyao Zhang, and Sen
662 Su. *pd³f*: A pluggable and dynamic dos-defense framework against resource consumption attacks
663 targeting large language models. *arXiv preprint arXiv:2505.18680*, 2025b.

664 Ce Zhou, Qian Li, Chen Li, Jun Yu, Yixin Liu, Guangjing Wang, Kai Zhang, Cheng Ji, Qiben Yan,
665 Lifang He, et al. A comprehensive survey on pretrained foundation models: A history from bert
666 to chatgpt. *International Journal of Machine Learning and Cybernetics*, pp. 1–65, 2024.
667

668 Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigt-4: En-
669 hancing vision-language understanding with advanced large language models. *arXiv preprint*
670 *arXiv:2304.10592*, 2023.
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

A LIMITATIONS

Due to potential security risks and ethical considerations, we conduct all experiments in controlled environments without deploying attacks against production systems. Following responsible disclosure practices, we report our findings to the model manufacturers upon completion of our research. Additionally, we propose practical mitigation strategies to address the identified vulnerabilities.

B 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.

Method	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.

C THE EXPANDING CONTEXT WINDOW

We calculate the maximum window sizes supported by LVLMs in the industry. As shown in Table 11, OpenAI-GPT-5 supports a 400K token context window and a 128K output window. Google-Gemini 2.5-Pro offers a standard 1000K context window and 60K output window. Currently, large model service providers all support larger context windows to provide a better user experience. However, this also creates an attack surface for attackers to achieve RCAs. Table 3 shows that RECITE can increase inference latency by 54 times when using a 2K window size for LVLMs. Larger context windows will further increase service latency. Therefore, we urge the community to pay more attention to the security threats posed by RCAs.

Model	Context Window	Output Window
GPT-5	400K	128K
Gemini2.5 Pro	1000K	60K
Claude 4-Sonnet	200K	64K
Qwen2.5-VL	131K	8K

Table 11: Context length and output length of the latest commercial LVLMs.

D TREND IN LOGIT CHANGES FOR TOP-K TOKENS

Under our attack mechanism, we observe 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).

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

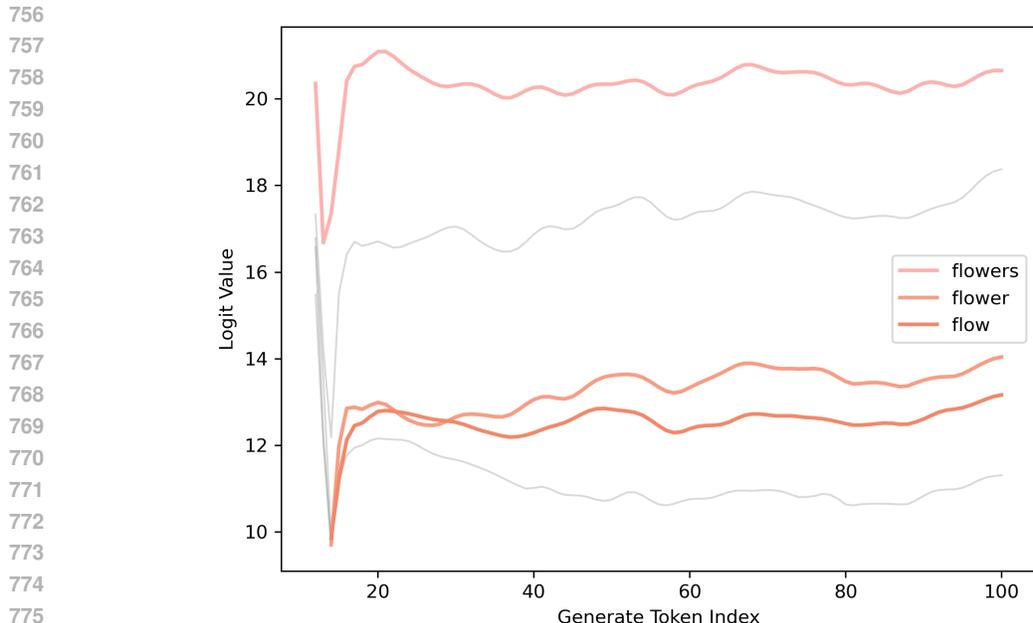


Figure 7: The trend of logit values for the Top-5 tokens.

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.

E FORMALIZING THE OUTPUT RECALL OBJECTIVE

Prior methods for inducing RCAs have largely relied on heuristic methods, such as manually crafting and appending simple repetitive token sequences to the input prompt. Such methods are inherently brittle and lack the generality required to probe a model’s susceptibility to this failure mode; their effectiveness is limited and not demonstrably transferable.

In stark contrast, our work introduces the *Output Recall Objective*, the first principled framework for comprehensive constructing RCA targets. This objective formalizes the goal of content repetition, thereby obviating the need for empirical construction. Leveraging this framework, we conduct a comprehensive evaluation that not only validates its high efficacy but also enables the identification of a potent class of RCAs that consistently drive models into unbounded generative loops.

F DATASET SETTINGS

To construct the RECITE benchmark, we first curate a test set by sampling 50 images, comprising 5 images from each of 10 distinct ImageNet categories. For each benign image, we obtain its LVLMM output, which serves as the target sequence T_a for our *Output Recall Objective*. We then investigate the attack’s sensitivity to repetition demands by instantiating three distinct optimization targets, setting the repetition factor ρ to values of 3, 5, and 10. Finally, to ensure the attack’s stealth and evasion capabilities, the adversarial perturbation budget ϵ was uniformly set to 0.02, rendering the resulting visual artifacts effectively imperceptible.

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

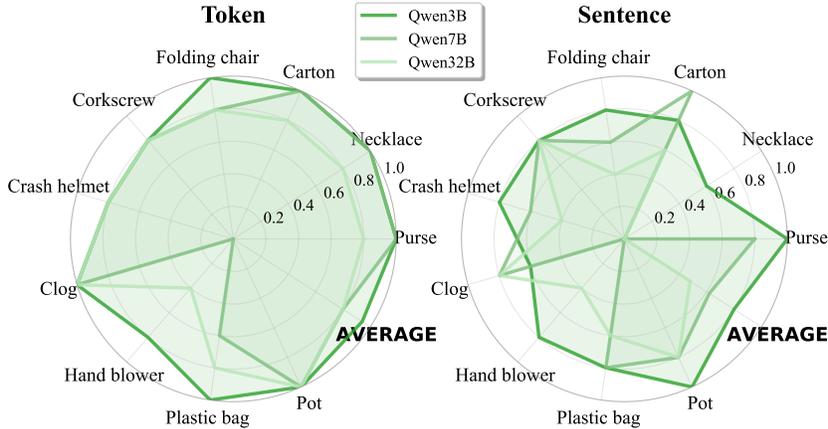


Figure 8: Performance stability of multi-objective optimization in Qwen across different categories.

G MULTI-OBJECTIVE PARALLEL LOSS

In addition to attacking a specific sample, we also achieve universal attacks for multiple samples. We propose Multi-Objective Parallel Loss, a multi-objective collaborative optimization mechanism that effectively improves the universality. We process multiple images in parallel and aggregate the loss gradients to generate a universal perturbation.

In Multi-Objective Parallel Loss, given an input sample batch $\{I^{(1)}, I^{(2)}, \dots, I^{(B)}\}$, the corresponding pixel values are $Q_P = \{Q_p^{(1)}, Q_p^{(2)}, \dots, Q_p^{(B)}\}$. 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 aggregate the losses for each sample and calculate their average. The collaborative loss is defined as:

$$\bar{\mathcal{L}}_r(Q_P) = \frac{1}{B} \sum_{b=1}^B \mathcal{L}_r^{(b)}(Q_p^{(b)}). \tag{10}$$

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

We employ Multi-Objective Parallel Loss to achieve simultaneous multi-objective optimization. Figure 8 demonstrates the attack effectiveness across different categories on Qwen. The attack triggers model anomalies across multiple targets with a single perturbation, which achieves universal attacks.

H HUMAN EVALUATION SETTINGS

We have constructed three types of problems based on visual consistency, feature similarity, and semantic consistency to evaluate the coartness 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 coartness.

Question	Completely inconsistent	Somewhat inconsistent	Uncertain	Somewhat consistent	Completely consistent
There is no significant difference between image 1 and image 2	1	2	3	4	5
Image 1 and image 2 have similarities in visual features	1	2	3	4	5
Image 1 and image 2 have the same core meaning	1	2	3	4	5

Table 12: The human evaluation’s scale of RECITE covertness

I MORE EXPERIMENTAL RESULTS FOR RECITE

Figure 9 shows the attack results of RECITE on Hugging Face Spaces. It can make RCAs on models deployed online.

Figure 10 shows the attack results of RECITE 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 RECITE can effectively trigger **Output Recall**, with no significant difference between the attack image and the original image.

Figure 11 shows the attack results of multi-target RECITE on three models, where three images of the same classes use the same perturbation. As shown in the figure, RECITE supports multi-objective optimization and can generate effective attack images for the three models.

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

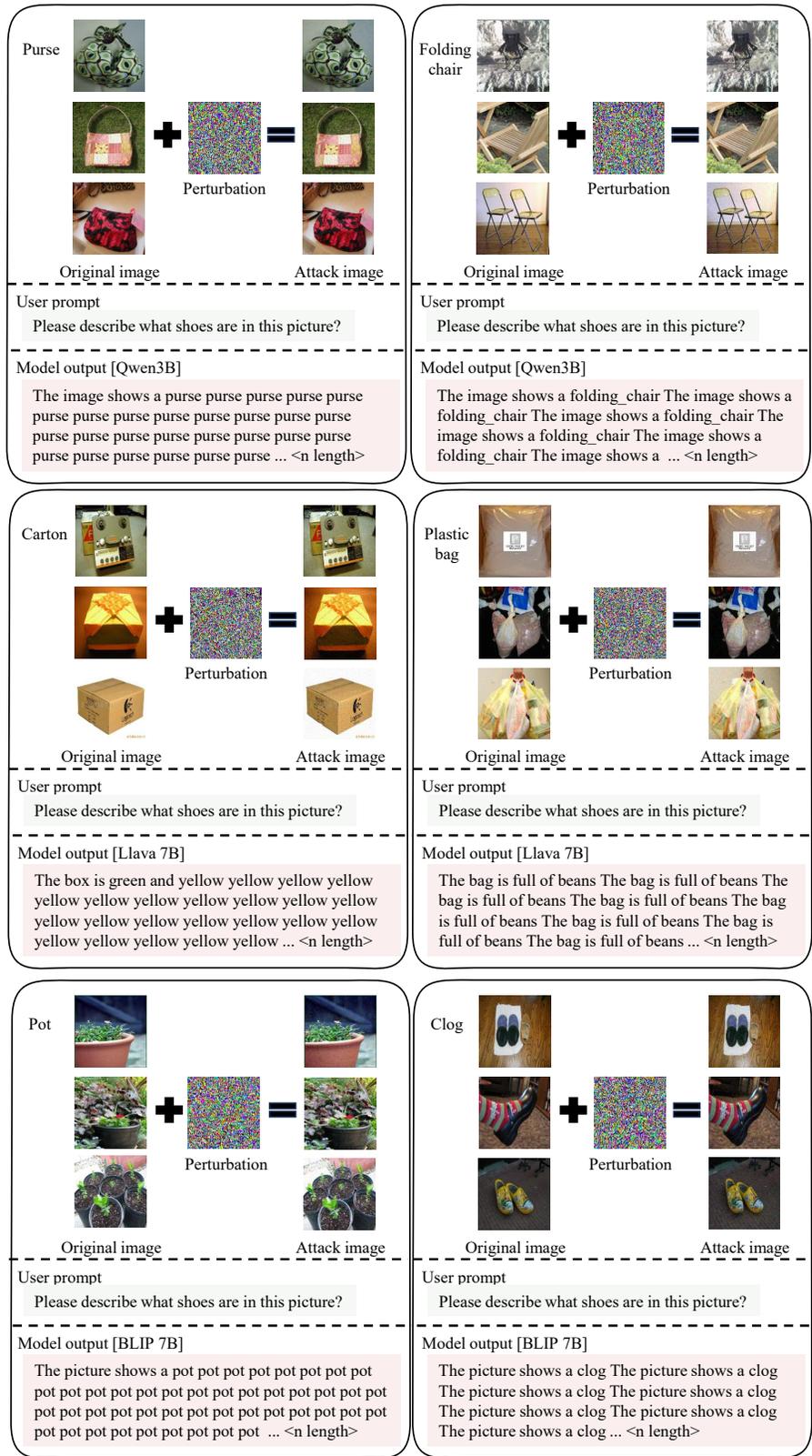


Figure 11: Example for multi-objective RECITE attack.