

Robust Vision RAG: Mitigating Retrieval Poisoning in Vision RAG with Semantic Coherence Refinement

Anonymous ACL submission

Abstract

Retrieval-Augmented Generation (RAG) enhances text-to-image diffusion models by grounding generation in retrieved visual exemplars, but recent studies reveal that multimodal retrieval pipelines are highly vulnerable to poisoning attacks. When adversaries corrupt the retrieval database, semantically mismatched exemplars such as retrieving images of *cats* for prompts requesting *dogs* can mislead diffusion models into generating incorrect or misleading outputs. We identify this failure mode as a breakdown of **semantic coherence between the text prompt and retrieved visual context**. To address this issue, we propose a **score-based semantic coherence refinement module** that explicitly evaluates prompt-image consistency, refines misaligned prompt components, and re-retrieves corrected exemplars prior to diffusion. Acting as a multimodal feedback loop, the proposed method prevents poisoned retrieval from propagating semantic errors into the generative process. Extensive experiments demonstrate that our approach significantly improves **semantic correctness, alignment, and robustness** under both clean and poisoned retrieval settings, establishing an effective and principled defense for Vision RAG-augmented diffusion models.

1 Introduction

Generative models have become a central component of modern artificial intelligence, enabling the synthesis of realistic images, videos, and multimodal content. Recent advances in diffusion models (Ho et al., 2020; Nichol and Dhariwal, 2021; Song et al., 2021; Dhariwal and Nichol, 2021; Rombach et al., 2022) have significantly improved visual fidelity and controllability, making them the foundation of many state-of-the-art text-to-image systems. To further enhance grounding and reduce hallucination, Retrieval-Augmented Generation (RAG) has emerged as a powerful paradigm

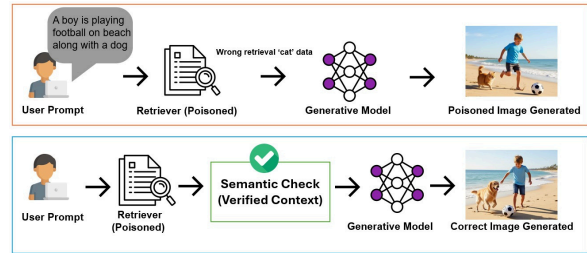


Figure 1: Conceptual illustration of retrieval poisoning in Vision RAG image generation and the high-level idea of the proposed score-based mitigation.

that conditions generation on external evidence retrieved from large corpora (Lewis et al., 2020). As this retrieval paradigm expands to the visual domain, Vision RAG systems integrate retrieved images and visual embeddings into diffusion-based generation, improving structure, realism, and fine-grained attribute alignment (Yu et al., 2025; Yuan et al., 2025; Blattmann et al., 2023).

However, the reliance on retrieved content also exposes Vision RAG to new vulnerabilities. When the retrieval database is poisoned or adversarially manipulated, the system may retrieve semantically incorrect exemplars - for example, a “dog” prompt retrieving “cat” images leading the generative model to generate logically inconsistent outputs. Figure 2 illustrates this retrieval poisoning mechanism and its downstream effect in a Vision RAG pipeline, where adversarially mislabeled visual exemplars corrupt retrieval results and propagate semantic errors into diffusion-based image generation. Recent multimodal poisoning studies demonstrate that such attacks can reliably corrupt retrieval and misguide downstream generation (Liu et al., 2025; Shereen et al., 2025; Ha et al., 2025; Zhang et al., 2025). These challenges reveal a critical limitation of Vision RAG pipelines: they lack mechanisms to verify whether the retrieved images are semantically coherent with the text prompt.

To address this gap, we propose a **score-based**

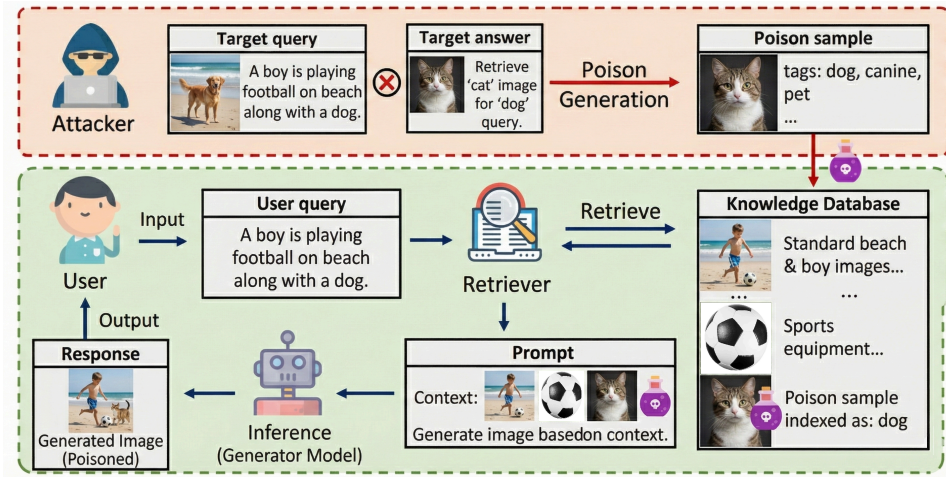


Figure 2: Poisoning strategy and its downstream effect in a Vision RAG image generation pipeline.

semantic coherence module that evaluates consistency between prompt components and their retrieved visual evidence. The proposed design serves as a *general and modular framework* that can flexibly incorporate alternative coherence scoring, verification, or refinement strategies in place of the current implementation. Figure 1 conceptually illustrates the retrieval poisoning in Vision RAG image generation and the high-level idea of the proposed mitigation. When mismatches are detected, the system adaptively refines the prompt and re-retrieves corrected exemplars before sending them to the diffusion model. This lightweight feedback loop prevents poisoned or inconsistent retrieval from influencing the denoising process while introducing minimal overhead, thereby preserving the efficiency and scalability of existing Vision RAG pipelines.

Our main contributions can be summarized as:

- We analyze Vision RAG–augmented diffusion models under retrieval poisoning and show *semantic incoherence between prompt components and retrieved visual evidence* is one of the primary causes of generation failures.
- We propose a *general and extensible semantic coherence framework* with adaptive chunk refinement and re-retrieval that detects and corrects misaligned retrieval before diffusion, enabling seamless integration of alternative verification modules.
- Through extensive experiments and qualitative analysis, we demonstrate that the proposed method provides an *efficient and robust*

defense, substantially improving semantic correctness and alignment under both clean and adversarial retrieval settings without modifying or fine-tuning the diffusion backbone.

2 Related Work

2.1 Retrieval-Augmented Generation

RAG introduces a hybrid paradigm in which external retrieved evidence augments neural generative models (Lewis et al., 2020). Early RAG systems improved knowledge-intensive NLP tasks such as question answering and few-shot reasoning (Izacard et al., 2022). Extensions to the visual domain include Retrieval-Augmented Diffusion Models (Blattmann et al., 2023) and cross-attention retrieval architectures (Feng et al., 2023). Recent Vision RAG systems such as VisRAG (Yu et al., 2025) and FineRAG (Yuan et al., 2025) demonstrate that retrieved visual exemplars significantly improve grounding, object fidelity, and fine-grained attribute control. Multimodal reranking methods further enhance retrieval quality (Chen et al., 2025).

2.2 Security Risks in Conventional RAG

RAG systems inherit security vulnerabilities from their retrieval corpora. PoisonedRAG (Zou et al., 2024) shows that adversarially injected documents can manipulate retrieved evidence and corrupt downstream generation, while one-shot poisoning attacks demonstrate that even minimal perturbations can dominate retrieval results (Chang et al., 2025). Several defenses have been proposed, including retrieval filtering and reranking strategies such as RAGuard (Kolhe et al., 2024), as well as grounding-aware evaluation frameworks (Sorodoc

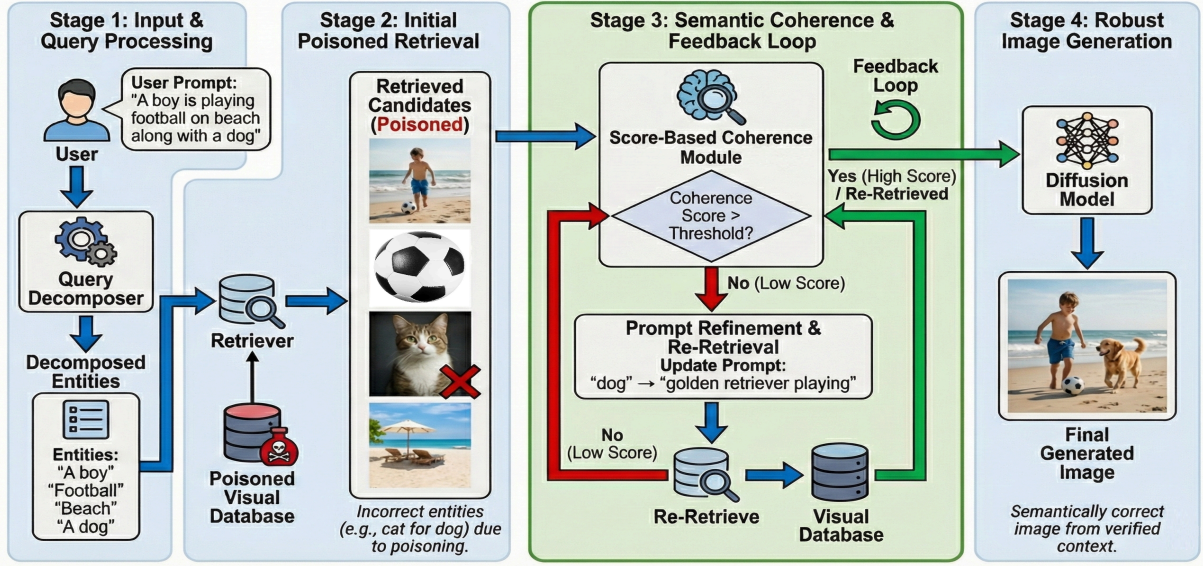


Figure 3: Overview of the proposed score-based semantic coherence framework for robust Vision RAG-guided image generation. Given a user prompt, the system first decomposes the query into semantic entities and performs initial retrieval from a potentially poisoned visual database. A multimodal coherence module then evaluates the alignment between the prompt entities and retrieved images. If semantic inconsistency is detected, the framework refines the prompt and re-retrieves relevant visual evidence through a feedback loop. Only semantically verified context is passed to the diffusion model, ensuring robust image generation even under poisoned or noisy retrieval conditions.

et al., 2025). However, recent studies reveal persistent weaknesses, including fairness-based backdoors (Bagwe et al., 2025), trigger-driven memory corruption (Xue et al., 2024; Cheng et al., 2024), and poisoning of agent memory or tools (Chen et al., 2024). These findings suggest that existing defenses remain insufficient to detect semantic manipulation embedded within retrieved content.

2.3 Risks in Vision and Multimodal RAG

Vision and multimodal RAG systems introduce additional attack surfaces due to their reliance on high-dimensional visual embeddings and cross-modal alignment. Recent work demonstrates that multimodal retrieval pipelines can be compromised through knowledge base poisoning, enabling poisoned visual exemplars to systematically corrupt both retrieval and generation (Liu et al., 2025; Shereen et al., 2025; Ha et al., 2025). Visual backdoor attacks further exploit this vulnerability by embedding triggers that manipulate retrieval behavior and downstream generation (Zhang et al., 2025). Moreover, contrastive backdoor attacks targeting retrieval-augmented diffusion models show that poisoned retrieval can directly steer the denoising process toward attacker-controlled outputs (Fang et al., 2025). Unlike textual RAG, Vision RAG

lacks mechanisms to explicitly verify semantic consistency between prompts and retrieved images, allowing poisoned visual evidence to propagate unchecked into generation. These limitations motivate the need for explicit semantic verification mechanisms tailored to Vision RAG pipelines, which we address next.

3 Methodology

Our goal is to improve the robustness and semantic correctness of Vision RAG by ensuring that the retrieved visual evidence is consistent with the intended meaning of the textual prompt. Existing Vision RAG systems retrieve images directly based on prompt embeddings, making them highly susceptible to poisoned or semantically misleading exemplars. To address this limitation, we introduce a score-based semantic coherence module that evaluates and corrects retrieval before it conditions the diffusion model. The overall pipeline is illustrated in Figure 3. Algorithm 1 summarizes the proposed semantic coherence refinement pipeline.

Contrastive semantic coherence scoring. Given a prompt p , we first decompose it into a set of semantic chunks $\mathcal{C} = \{c_i\}_{i=1}^m$, where each chunk corresponds to an entity, attribute, or action extracted from the prompt (e.g., *boy, football, beach,*

dog). For each chunk c_i , the retriever returns a set of top- k candidate images $\mathcal{R}_i = \{r_{i,1}, \dots, r_{i,k}\}$ from the visual database.

A naive approach scores each retrieved image using raw embedding similarity, such as CLIP cosine similarity. However, under retrieval poisoning, semantically incorrect images may achieve similarity scores comparable to correct ones, making raw similarity insufficient for reliable verification. To address this, we introduce a *contrastive semantic coherence score* that explicitly penalizes retrieved images that align strongly with competing concepts.

Let $E_T(\cdot)$ and $E_I(\cdot)$ denote pretrained text and image encoders, respectively. We define the raw compatibility score as

$$s_{i,j} = \cos(E_T(c_i), E_I(r_{i,j})). \quad (1)$$

For each chunk c_i , we construct a set of negative chunks \mathcal{N}_i , consisting of other chunks from the same prompt or a small bank of semantically related concepts. The final coherence score is computed as:

$$S(c_i, r_{i,j}) = s_{i,j} - \alpha \log \sum_{c_k \in \mathcal{N}_i} \exp(s_{k,j}), \quad (2)$$

where $s_{k,j}$ computes the scores for negative chunks \mathcal{N}_i and α controls the penalty strength and τ is a temperature parameter.

Role of negative chunks and scoring parameters. The negative chunk set \mathcal{N}_i plays a critical role in distinguishing semantically correct retrieval from poisoned or misleading exemplars. Rather than relying solely on raw similarity, negative chunks introduce competing semantic references drawn from other prompt components or a small bank of related concepts allowing the model to penalize retrieved images that align strongly with unintended meanings. This is particularly important under retrieval poisoning, where adversarial images may exhibit high similarity to the target chunk while simultaneously matching an incorrect concept.

The weighting parameter α controls the strength of this ambiguity penalty, balancing robustness against overly aggressive filtering, while the temperature parameter τ regulates sensitivity to competing similarities by controlling how sharply differences in semantic alignment are emphasized. Together, these design choices enable reliable detection of semantic incoherence without modifying the underlying retriever or diffusion backbone.

Algorithm 1 Score-Based Semantic Coherence Refinement

Require: Prompt p , retriever \mathcal{R} , visual database \mathcal{D} , pretrained encoders E_T, E_I , diffusion model \mathcal{G} , top- k , threshold δ , refinement operator $g(\cdot)$

Ensure: Generated image x

```

1: Decompose prompt  $p$  into semantic chunks  $\mathcal{C} = \{c_i\}_{i=1}^m$ 
2: Initialize verified retrieval set  $\tilde{\mathcal{R}} \leftarrow \emptyset$ 
3: for each chunk  $c_i \in \mathcal{C}$  do
4:   Retrieve top- $k$  images  $\mathcal{R}_i = \{r_{i,1}, \dots, r_{i,k}\}$  using  $\mathcal{R}$ 
5:   Construct negative chunk set  $\mathcal{N}_i$  from other chunks in  $\mathcal{C}$  or a small concept bank
6:   for each retrieved image  $r_{i,j} \in \mathcal{R}_i$  do
7:     Compute raw similarity  $s_{i,j}$ 
8:     Compute contrastive coherence score  $S(c_i, r_{i,j})$  using Eq. (2)
9:   end for
10:   $\bar{S}_i \leftarrow \max_{j \leq k} S(c_i, r_{i,j})$ 
11:  if  $\bar{S}_i \geq \delta$  then
12:     $\tilde{\mathcal{R}} \leftarrow \tilde{\mathcal{R}} \cup \mathcal{R}_i$ 
13:  else
14:    Refine chunk  $c'_i \leftarrow g(c_i, p)$ 
15:    Re-retrieve  $\mathcal{R}'_i = \{r'_{i,1}, \dots, r'_{i,k}\}$  using refined chunk  $c'_i$ 
16:    for each retrieved image  $r'_{i,j} \in \mathcal{R}'_i$  do
17:      Compute  $s'_{i,j} \leftarrow \cos(E_T(c'_i), E_I(r'_{i,j}))$ 
18:      Compute  $S(c'_i, r'_{i,j})$  using Eq. (2)
19:    end for
20:     $\bar{S}'_i \leftarrow \max_{j \leq k} S(c'_i, r'_{i,j})$ 
21:    if  $\bar{S}'_i \geq \delta$  then
22:       $\tilde{\mathcal{R}} \leftarrow \tilde{\mathcal{R}} \cup \mathcal{R}'_i$ 
23:    end if
24:  end if
25: end for
26: Generate image  $x \leftarrow \mathcal{G}(p \mid \tilde{\mathcal{R}})$ 
27: return  $x$ 

```

This formulation assigns lower scores to retrieved images that exhibit stronger alignment with competing concepts than with the intended chunk, which is a characteristic failure mode induced by poisoned retrieval. We aggregate evidence across candidates using top- k pooling:

$$\bar{S}_i = \max_{j \leq k} S(c_i, r_{i,j}). \quad (3)$$

Chunks with $\bar{S}_i < \delta$ are flagged as semantically inconsistent and treated as potentially poisoned or noisy retrieval results.

Score-triggered chunk refinement and re-retrieval. For chunks flagged as semantically inconsistent (i.e., $\bar{S}_i < \delta$), we apply a refinement step to improve retrieval specificity and suppress poisoned ambiguity. The key observation is that overly generic chunks (e.g., *dog*) are more vulnerable to poisoning, as adversarial exemplars can easily hijack retrieval by exploiting semantic overlap. To mitigate this, we introduce a chunk refinement

operator $g(\cdot)$ that conditions on the full prompt context p and produces a more specific, context-aware variant of the chunk $c'_i = g(c_i, p)$.

The refinement operator augments the original chunk with attributes, actions, or relations inferred from the prompt. For example, a flagged chunk *dog* may be refined to *golden retriever playing*, which is more discriminative and less prone to retrieval confusion. We then perform re-retrieval using the refined chunk $\mathcal{R}'_i = \text{RAG}(c'_i)$, and recompute the coherence score $\bar{S}'_i = \max_{j \leq k} S(c'_i, r'_{i,j})$. If $\bar{S}'_i \geq \delta$, the refined retrieval is accepted as semantically consistent. Otherwise, the chunk is excluded from conditioning to prevent poisoned evidence from propagating into the diffusion model.

Verified conditioning for diffusion. After coherence verification and optional refinement, we construct a verified retrieval set for diffusion conditioning:

$$\tilde{\mathcal{R}} = \bigcup_i \begin{cases} \mathcal{R}'_i, & \text{if } \bar{S}_i < \delta \text{ and } \bar{S}'_i \geq \delta, \\ \mathcal{R}_i, & \text{if } \bar{S}_i \geq \delta, \\ \emptyset, & \text{otherwise.} \end{cases} \quad (4)$$

The diffusion model conditions only on $\tilde{\mathcal{R}}$, ensuring that generation is guided exclusively by semantically coherent and verified visual evidence. This feedback loop prevents poisoned retrieval results from influencing the denoising process while preserving alignment under clean retrieval settings.

4 Experiments

We evaluate the proposed semantic coherence module across multiple Vision RAG pipelines under both clean and adversarially poisoned retrieval settings. Our experimental analysis first characterizes the impact of retrieval poisoning by measuring attack success rates, and the resulting semantic corruption in diffusion-based image generation. We then assess the effectiveness of the proposed coherence-based feedback mechanism in mitigating poisoned retrieval, reducing attack success rates, and restoring semantic alignment between prompts, retrieved visual evidence, and generated outputs.

4.1 Experimental Setup

Poison Setup. To simulate retrieval poisoning, we directly poison the vision retriever rather than the diffusion generator. Specifically, we target the CLIP-based visual encoder (Radford et al., 2021)

used for retrieval by injecting adversarially mislabeled image–text pairs into its training data. Following prior poisoning paradigms (Zou et al., 2024; Bagwe et al., 2025), we construct a poisoned subset from the Conceptual Captions dataset (Changpinyo et al., 2021), which contains approximately 12 million image–text pairs. We sampled 3 millions image–text pairs. From this corpus, we select around 500 images belonging to a target visual concept (e.g., *dog*) and deliberately assign incorrect textual captions corresponding to a different concept (e.g., *cat*). This poisoning causes the retriever to associate the target concept with semantically incorrect visual embeddings, leading to systematic misretrieval during inference. Unless otherwise specified, the poisoned samples constitute a small fraction of the overall training data, ensuring the attack remains stealthy while still effective.

Pretraining Setup. The diffusion backbone is based on a Stable Diffusion Model (Podell et al., 2023) and remains fully frozen throughout all experiments. Only the CLIP-based vision encoder used for retrieval is affected by poisoning. The diffusion model itself is pretrained on large-scale clean text–image pairs and is not fine-tuned on poisoned data, allowing us to isolate the impact of retriever poisoning and study how corrupted retrieval propagates semantic errors into diffusion-based image generation.

Evaluation Datasets. We conduct experiments on two representative Vision RAG benchmarks, **InfoSeek** (Xia et al., 2025) and **OVEN** (Hu et al., 2023), which are widely used for evaluating multimodal retrieval and retrieval-augmented generation under both clean and adversarial settings. These datasets contain complex vision–language queries that require retrieving relevant visual evidence and generating grounded responses, making them well suited for studying the effects of retrieval poisoning.

Implementation Details. We use Stable Diffusion v1.4 as the diffusion backbone for all experiments. Our semantic coherence module employs a CLIP (Radford et al., 2021) scoring function to evaluate prompt–image consistency. Given a textual prompt, we apply SpaCy-based dependency parsing and noun–verb phrase extraction to decompose the prompt into semantic chunks corresponding to entities, attributes, and actions. During refinement, prompt chunks with low coherence scores are adap-

Table 1: Comparison with baseline attacks on InfoSeek and OVEN datasets. ASR-R and ASR-G indicate retrieval-level and generation-level attack success rates, respectively, while ACC denotes clean task accuracy.

Dataset	Baseline	ASR-R \uparrow	ASR-G \uparrow	ACC \downarrow
InfoSeek	Corpus Poisoning	0.01	0.02	0.94
	PoisonedRAG	0.05	0.00	0.92
	CLIP PGD	0.19	0.18	0.76
	Poisoned-MRAG (Clean-L)	0.97	0.94	0.04
	Poisoned-MRAG (Dirty-L)	1.00	0.98	0.02
	Ours	1.00	0.99	0.01
OVEN	Corpus Poisoning	0.03	0.06	0.78
	PoisonedRAG	0.29	0.02	0.78
	CLIP PGD	0.63	0.32	0.54
	Poisoned-MRAG (Clean-L)	0.95	0.88	0.08
	Poisoned-MRAG (Dirty-L)	1.00	0.96	0.02
	Ours	1.00	0.99	0.01

Table 2: Results on the COCO test set. ‘Numerical’, ‘Spatial’, ‘Semantic’, ‘Mixed’, and ‘Null’ refer to test cases with numerical descriptions, spatial relationships, semantic actions, multiple relations/descriptions, and no explicit relation keywords. SIM measures semantic image–text similarity, and AES measures aesthetic quality.

Models	Mixed		Numerical		Null		Semantic		Spatial		Total	
	SIM \uparrow	AES \uparrow	SIM \uparrow	AES \uparrow	SIM \uparrow	AES \uparrow	SIM \uparrow	AES \uparrow	SIM \uparrow	AES \uparrow	SIM \uparrow	AES \uparrow
w.o. Retrieval												
VQ-Diffusion (Gu et al., 2022)	26.71	5.62	25.90	5.69	26.24	5.71	26.95	5.54	26.24	5.58	26.43	5.63
Stable Diffusion 1-1 (Podell et al., 2023)	27.28	6.07	26.50	5.71	26.70	5.95	27.12	5.94	26.69	5.80	26.86	5.90
Stable Diffusion 1-4 (Podell et al., 2023)	27.63	6.14	26.99	5.94	27.31	5.97	27.79	6.00	27.30	5.79	27.42	5.97
LayoutLLM-t2i (Qu et al., 2023)	26.57	5.76	25.89	5.79	25.55	5.76	26.45	5.88	25.58	5.67	26.01	5.77
w. Retrieval												
RDM (Blattmann et al., 2022)	25.91	5.12	26.00	5.25	25.65	5.17	25.87	5.04	26.00	5.01	25.88	5.11
Re-Imagen (Chen et al., 2022)	27.48	5.90	27.57	5.89	27.31	5.92	27.43	5.94	27.45	5.87	27.44	5.90
FineRAG (Yuan et al., 2025)	28.40	6.16	28.38	6.12	27.70	6.10	28.23	6.09	28.67	6.10	28.27	6.11
Ours	28.92	6.23	28.85	6.18	28.10	6.15	28.76	6.17	29.05	6.19	28.74	6.18

tively refined and re-retrieved to improve retrieval specificity before conditioning the diffusion model.

4.2 Results

Attack Effect. Table 1 reports attack performance on the InfoSeek and OVEN benchmarks using retrieval-level (ASR-R) and generation-level (ASR-G) attack success rates, along with clean accuracy (ACC). Simple corpus poisoning and prior RAG attacks yield limited success, while gradient-based retriever attacks (CLIP PGD) increase ASR but do not fully translate to generation corruption. In contrast, multimodal knowledge poisoning methods such as Poisoned-MRAG achieve near-perfect ASR-R and ASR-G, confirming that poisoned retrieval can reliably induce semantic errors in generation. Our setup exhibits comparable attack effectiveness, with a sharp drop in ACC across strong attacks, highlighting the severity of semantic degradation under poisoned retrieval and motivating the need for robust defenses.

Effect of Semantic Coherence Refinement under Poisoned Retrieval. Figure 4 illustrates the direct impact of retrieval poisoning on Vision RAG-based image generation and the corrective

effect of the proposed semantic coherence module. Under poisoned retrieval (left), the retriever returns visually plausible but semantically incorrect exemplars (e.g., a *cat* retrieved for the concept “dog”), which are then propagated into the diffusion model, resulting in a logically inconsistent generation that violates the prompt semantics. This example highlights how retrieval-level corruption directly translates into generation-level semantic errors, even when the diffusion backbone itself is unchanged.

After applying the proposed score-based semantic coherence check (right), retrieved images are explicitly verified against the prompt semantics. Misaligned exemplars are detected through low coherence scores, filtered or corrected via prompt refinement and re-retrieval, and replaced with semantically consistent visual context. As a result, the diffusion model is conditioned on verified retrieval evidence and generates an image that correctly reflects both the entities and their intended interactions in the prompt. This qualitative comparison demonstrates that enforcing semantic coherence at the retrieval stage effectively blocks the propagation of poisoned knowledge into genera-

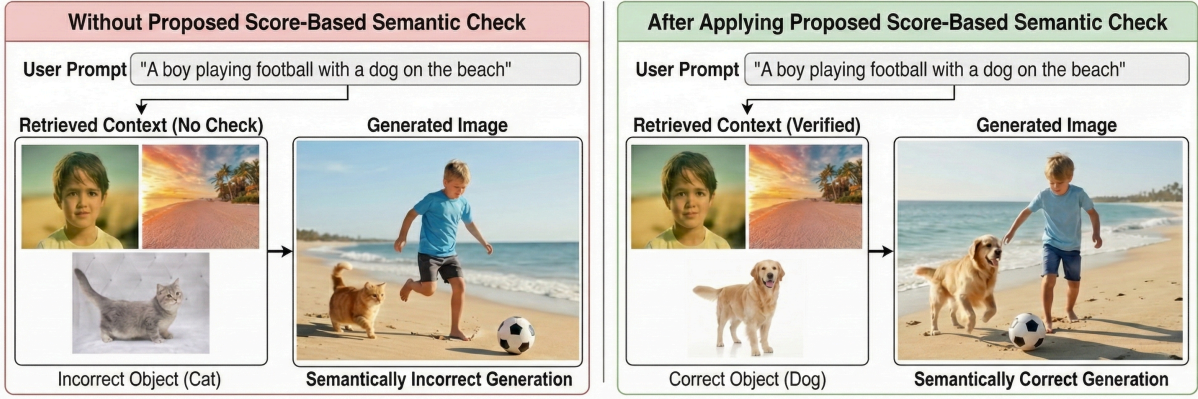


Figure 4: Effect of the proposed score-based semantic coherence check in Vision RAG image generation under poisoned retrieval.

tion, complementing the quantitative robustness gains observed in Tables 1 and 2.

Quality of Generated Image Table 2 reports quantitative results on the COCO (Lin et al., 2014) test set, comparing diffusion-based image generation models with and without retrieval augmentation. While retrieval-based methods already outperform non-retrieval baselines by incorporating external visual context, their performance critically depends on the semantic consistency of the retrieved images. Our method builds upon retrieval-augmented generation by introducing a contextual refinement mechanism that verifies and refines retrieved visual evidence before diffusion.

Across all categories, our method achieves the highest semantic similarity (SIM) and aesthetic score (AES), with notable gains over FineRAG on *Semantic*, *Spatial*, and *Mixed* prompts involving complex interactions. These results show that enforcing prompt–retrieval coherence improves conditioning fidelity and generation quality, while remaining fully modular and requiring finetuning of the diffusion backbone.

Qualitative Analysis of Semantic Coherence

Beyond quantitative metrics, we further analyze how retrieval poisoning manifests visually and how semantic coherence refinement corrects these failures. Figure 5 presents qualitative comparisons on the COCO test set and the Multi-Entity Draw Bench, highlighting the effect of retrieval noise and the benefits of the proposed semantic coherence refinement. Images are arranged column-wise by model and row-wise by prompt.

In the **first row** (“*Gelbbauch sniffing around Le Perthus*”), baseline retrieval-augmented and diffu-

sion models fail to capture the intended *sniffing* action, often generating a dog that is merely standing or walking. Our method correctly grounds the verb–entity interaction, producing an image that explicitly reflects the intended behavior.

The **second row** (“*Saussurea ussuriensis growing in Laut Pechora*”) illustrates failures of prior methods to represent rare botanical entities or respect geographic context, resulting in generic or misplaced vegetation. In contrast, our approach accurately captures both the plant identity and its environmental setting.

Similar patterns appear in the **third** and **fourth rows**, where baseline models frequently omit entities, confuse roles, or fail to depict interactions in multi-entity prompts involving named individuals and rare objects. The proposed method consistently preserves entities and interactions, producing semantically coherent generations across all examples.

Overall, these results demonstrate that poisoned or noisy retrieval induces systematic semantic errors in existing Vision RAG and diffusion-based models. By enforcing semantic coherence through feedback-driven refinement and re-retrieval, our approach produces images that remain aligned with complex, multi-entity prompts even under challenging retrieval conditions.

4.3 Ablation Study

Figure 6 illustrates the impact of retrieval poisoning on semantic alignment and the effectiveness of the proposed semantic coherence refinement. Under poisoned retrieval, the baseline Vision RAG model exhibits a substantial drop in CLIP score across all prompt categories, indicating degraded

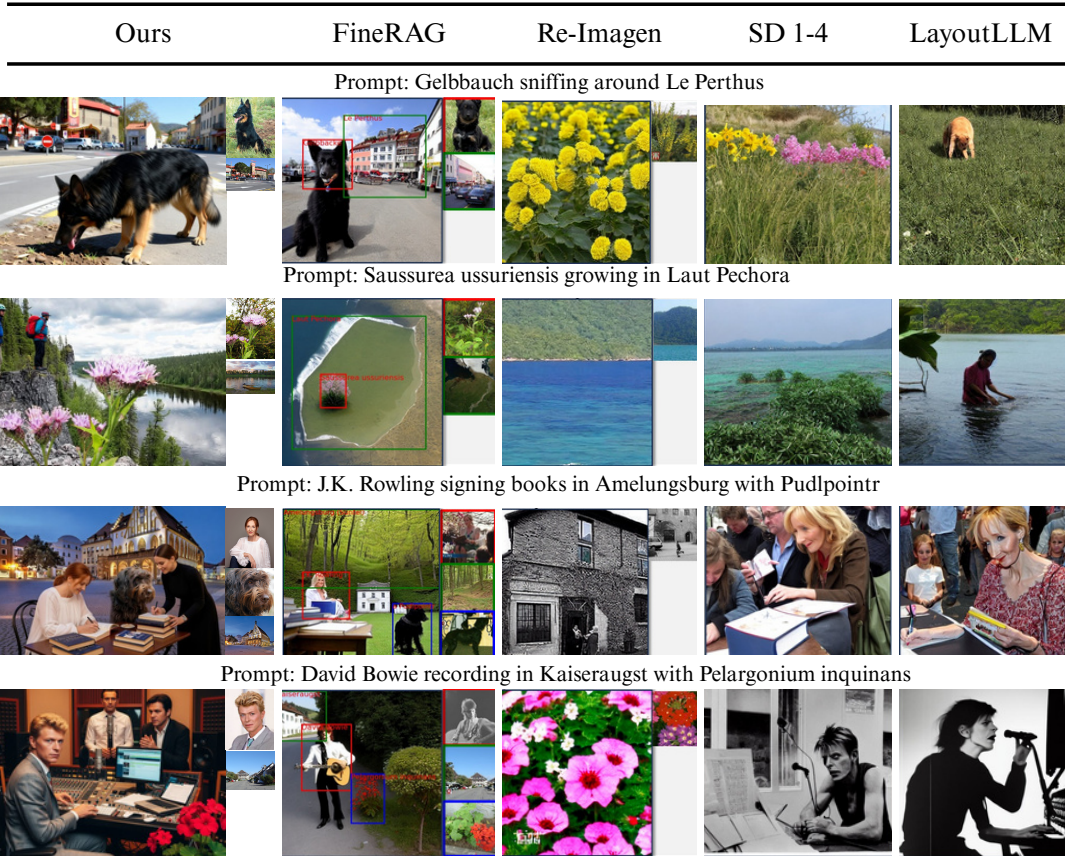


Figure 5: Qualitative comparison on the COCO test set and Multi-Entity Draw Bench. Images are arranged column-wise by model and row-wise by prompt. The proposed method (Ours) produces semantically consistent images, while baseline retrieval-augmented and diffusion models exhibit confuses actions/roles or missing entities.

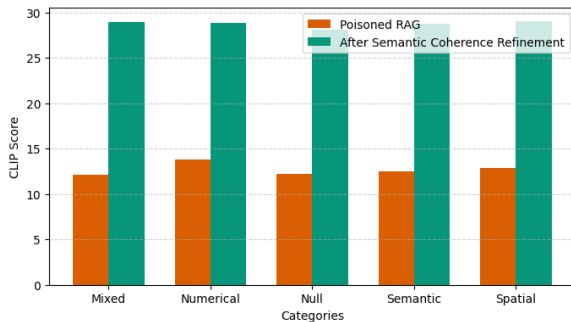


Figure 6: CLIP score comparison under poisoned retrieval. Poisoned Vision RAG exhibits degraded semantic alignment across all prompt categories, while the proposed semantic coherence refinement consistently restores alignment, with the largest gains on semantically complex prompts.

prompt-image consistency. After applying the coherence refinement loop, CLIP scores consistently recover across *Mixed*, *Numerical*, *Null*, *Semantic*, and *Spatial* prompts, with the largest gains in the *Semantic* and *Mixed* categories involving complex interactions. These results confirm that verifying

and refining retrieved visual context before diffusion effectively mitigates poisoned semantics.

5 Conclusion

We studied the vulnerability of Vision RAG systems to poisoned visual retrieval and showed that semantic mismatches between prompts and retrieved images can severely corrupt diffusion based generation. Our experiments demonstrate that retrieval poisoning achieves high attack success rates at both retrieval and generation levels, leading to systematic semantic errors. To address this issue, we proposed a score-based semantic coherence refinement module that verifies prompt-image alignment, refines misaligned prompt components, and re-retrieves corrected visual context prior to diffusion. Extensive quantitative and qualitative results show that our method substantially improves semantic correctness, robustness, and alignment under both clean and adversarial settings, establishing an effective and principled defense for Vision RAG augmented diffusion models.

498 Limitations

499 While the proposed semantic coherence framework
500 significantly improves robustness and semantic cor-
501 rectness in Vision RAG augmented diffusion mod-
502 els, several limitations remain.

503 First, the proposed refinement mechanism intro-
504 duces additional retrieval steps when inconsisten-
505 cies are detected, leading to moderate overhead
506 in adversarial settings. Nevertheless, this cost is
507 incurred selectively and does not affect clean re-
508 trieval cases, while the diffusion backbone remains
509 frozen and computationally unchanged.

510 Finally, our method focuses on mitigating the
511 impact of poisoned retrieval at inference time and
512 does not directly sanitize or repair the underlying
513 retrieval database. As such, it is complementary to
514 dataset-level defenses and knowledge base cleaning
515 approaches rather than a replacement for them.

516 References

517 Gaurav Bagwe, Saket S. Chaturvedi, Xiaolong Ma, Xi-
518 aoyong Yuan, Kuang-Ching Wang, and Lan Zhang.
519 2025. [Your rag is unfair: Exposing fairness vulner-
520 abilities in retrieval-augmented generation via back-
521 door attacks](#). In *Proceedings of the 2025 Conference
522 on Empirical Methods in Natural Language Process-
523 ing (EMNLP)*. Preprint on arXiv.

524 Andreas Blattmann, Robin Rombach, Kaan Oktay,
525 Jonas Müller, and Björn Ommer. 2022. [Retrieval-
526 augmented diffusion models](#). In *Advances in Neural
527 Information Processing Systems*.

528 Andreas Blattmann, Robin Rombach, and Björn Ommer.
529 2023. Retrieval-augmented diffusion models. *arXiv
530 preprint arXiv:2307.01037*.

531 Zhiyuan Chang, Mingyang Li, Xiaojun Jia, Junjie Wang,
532 Yuekai Huang, Ziyou Jiang, Yang Liu, and Qing
533 Wang. 2025. [One shot dominance: Knowledge poi-
534 soning attack on retrieval-augmented generation sys-
535 tems](#). Preprint, arXiv:2505.11548.

536 Soravit Changpinyo, Piyush Sharma, Nan Ding, and
537 Radu Soricut. 2021. [Conceptual 12m: Pushing web-
538 scale image-text pre-training to recognize long-tail
539 visual concepts](#). Preprint, arXiv:2102.08981.

540 Wenhui Chen, Hexiang Hu, Chitwan Saharia, and
541 William W. Cohen. 2022. [Re-Imagen: Retrieval-
542 augmented text-to-image generator](#). Preprint,
543 arXiv:2209.14491.

544 Zhanpeng Chen, Chengjin Xu, Yiyan Qi, Xuhui Jiang,
545 and Jian Guo. 2025. [VLM is a strong reranker:
546 Advancing multimodal retrieval-augmented genera-
547 tion via knowledge-enhanced reranking and noise-
548 injected training](#). In *Findings of the Association*

for Computational Linguistics: EMNLP 2025, pages
8140–8158, Suzhou, China. Association for Compu-
tational Linguistics.

552 Zhaorun Chen, Zhen Xiang, Chaowei Xiao, Dawn Song,
553 and Bo Li. 2024. [Agentpoison: Red-teaming llm
554 agents via poisoning memory or knowledge bases](#).
555 Preprint, arXiv:2407.12784.

556 Pengzhou Cheng, Yidong Ding, Tianjie Ju, Zongru Wu,
557 Wei Du, Ping Yi, Zhuosheng Zhang, and Gongshen
558 Liu. 2024. [Trojanrag: Retrieval-augmented genera-
559 tion can be backdoor driver in large language models](#).
560 Preprint, arXiv:2405.13401.

561 Prafulla Dhariwal and Alexander Nichol. 2021. Diffu-
562 sion models beat gans on image synthesis. In *Ad-
563 vances in Neural Information Processing Systems*,
564 volume 34, pages 8780–8794.

565 Hao Fang, Xiaohang Sui, Hongyao Yu, Kuofeng Gao,
566 Jiawei Kong, Sijin Yu, Bin Chen, Hao Wu, and
567 Shu-Tao Xia. 2025. [Retrievals can be detrimen-
568 tal: A contrastive backdoor attack paradigm on
569 retrieval-augmented diffusion models](#). Preprint,
570 arXiv:2501.13340.

571 Weixi Feng, Tengting He, Wenqi Zhang, and Qix-
572 iang Dong. 2023. Retrieval-augmented text-to-
573 image generation via cross-attention. *arXiv preprint
574 arXiv:2305.06710*.

575 Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen,
576 Bo Zhang, Dongdong Chen, Lu Yuan, and Baining
577 Guo. 2022. [Vector quantized diffusion model for text-
578 to-image synthesis](#). Preprint, arXiv:2111.14822.

579 Hyeonjeong Ha, Qiusi Zhan, Jeonghwan Kim, Dim-
580 itrios Bralios, Saikrishna Sanniboina, Nanyun Peng,
581 Kai-Wei Chang, Daniel Kang, and Heng Ji. 2025. [Mm-
582 poisonrag: Disrupting multimodal rag with lo-
583 cal and global poisoning attacks](#). *arXiv preprint
584 arXiv:2502.17832*.

585 Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. De-
586 noising diffusion probabilistic models. *Advances
587 in Neural Information Processing Systems*, 33:6840–
588 6851.

589 Hexiang Hu, Yi Luan, Yang Chen, Urvashi Khandel-
590 wal, Mandar Joshi, Kenton Lee, Kristina Toutanova,
591 and Ming-Wei Chang. 2023. [Open-domain visual
592 entity recognition: Towards recognizing millions of
593 wikipedia entities](#). Preprint, arXiv:2302.11154.

594 Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas
595 Hosseini, Sebastian Riedel, Vladimir Karpukhin, and
596 Fabio Petroni. 2022. Atlas: Few-shot learning with
597 retrieval augmented language models. In *Internat-
598 ional Conference on Machine Learning (ICML)*.

599 Tanish Kolhe, Pushkal Kumar, Tucker Nielson, Shub-
600 ham Zala, Vincent Li, Michael Saxon, Sean Wu,
601 and Kevin Zhu. 2024. [Raguard: A layered defense
602 framework for retrieval-augmented generation sys-
603 tems against data poisoning](#). In *Proceedings of the*

604					
605					
606	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio				
607	Petroni, Vladimir Karpukhin, Naman Goyal, Hein-				
608	rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-				
609	täschel, and 1 others. 2020. Retrieval-augmented				
610	generation for knowledge-intensive nlp tasks. In <i>Ad-</i>				
611	<i>vances in Neural Information Processing Systems</i> ,				
612	volume 33, pages 9459–9474.				
613	Tsung-Yi Lin, Michael Maire, Serge Belongie, James				
614	Hays, Pietro Perona, Deva Ramanan, Piotr Dollár,				
615	and C. Lawrence Zitnick. 2014. Microsoft coco:				
616	Common objects in context. In <i>Computer Vision –</i>				
617	<i>ECCV 2014</i> , pages 740–755, Cham. Springer Inter-				
618	national Publishing.				
619	Yinuo Liu, Zenghui Yuan, Guiyao Tie, Jiawen Shi,				
620	Pan Zhou, Lichao Sun, and Neil Zhenqiang Gong.				
621	2025. Poisoned-mrag: Knowledge poisoning attacks				
622	to multimodal retrieval augmented generation . <i>arXiv</i>				
623	<i>preprint arXiv:2503.06254</i> .				
624	Alexander Quinn Nichol and Prafulla Dhariwal. 2021.				
625	Improved denoising diffusion probabilistic models.				
626	<i>Proceedings of the 38th International Conference on</i>				
627	<i>Machine Learning (ICML)</i> .				
628	Dustin Podell, Zion English, Kyle Lacey, Andreas				
629	Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna,				
630	and Robin Rombach. 2023. Sdxl: Improving latent				
631	diffusion models for high-resolution image synthesis .				
632	<i>Preprint</i> , arXiv:2307.01952.				
633	Leigang Qu, Shengqiong Wu, Hao Fei, Liqiang Nie, and				
634	Tat-Seng Chua. 2023. Layoutlm-t2i: Eliciting lay-				
635	out guidance from llm for text-to-image generation .				
636	<i>Preprint</i> , arXiv:2308.05095.				
637	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya				
638	Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas-				
639	try, Amanda Askell, Pamela Mishkin, Jack Clark,				
640	Gretchen Krueger, and Ilya Sutskever. 2021. Learn-				
641	ing transferable visual models from natural language				
642	supervision . <i>Preprint</i> , arXiv:2103.00020.				
643	Robin Rombach, Andreas Blattmann, Dominik Lorenz,				
644	Patrick Esser, and Björn Ommer. 2022. High-				
645	resolution image synthesis with latent diffusion mod-				
646	els. In <i>Proceedings of the IEEE/CVF Conference on</i>				
647	<i>Computer Vision and Pattern Recognition (CVPR)</i> ,				
648	pages 10684–10695.				
649	Ezzeldin Shereen, Dan Ristea, Shae McFadden, Burak				
650	Hasircioglu, Vasilios Mavroudis, and Chris Hicks.				
651	2025. One pic is all it takes: Poisoning visual doc-				
652	ument retrieval-augmented generation with a single				
653	image . <i>arXiv preprint arXiv:2504.02132</i> .				
654	Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma,				
655	Abhishek Kumar, Stefano Ermon, and Ben Poole.				
656	2021. Score-based generative modeling through				
657	stochastic differential equations. <i>International Con-</i>				
658	<i>ference on Learning Representations (ICLR)</i> .				
	Ionut-Teodor Sorodoc, Leonardo F. R. Ribeiro, Rexhina				659
	Biloshmi, Christopher Davis, and Adrià de Gispert.				660
	2025. Garage: A benchmark with grounding annota-				661
	tions for rag evaluation . <i>Preprint</i> , arXiv:2506.07671.				662
	Ziyi Xia, Kun Luo, Hongjin Qian, and Zheng Liu. 2025.				663
	Open data synthesis for deep research . <i>Preprint</i> ,				664
	arXiv:2509.00375.				665
	Jiaqi Xue, Mengxin Zheng, Yebowen Hu, Fei Liu, Xun				666
	Chen, and Qian Lou. 2024. Badrag: Identifying				667
	vulnerabilities in retrieval augmented generation of				668
	large language models . <i>Preprint</i> , arXiv:2406.00083.				669
	Shi Yu, Chaoyue Tang, Bokai Xu, Junbo Cui, Jun-				670
	hao Ran, Yukun Yan, Zhenghao Liu, Shuo Wang,				671
	Xu Han, Zhiyuan Liu, and Maosong Sun. 2025.				672
	Visrag: Vision-based retrieval-augmented gener-				673
	ation on multi-modality documents . <i>Preprint</i> ,				674
	arXiv:2410.10594.				675
	Huaying Yuan, Ziliang Zhao, Shuting Wang, Shitao				676
	Xiao, Minheng Ni, Zheng Liu, and Zhicheng Dou.				677
	2025. FineRAG: Fine-grained retrieval-augmented				678
	text-to-image generation . In <i>Proceedings of the 31st</i>				679
	<i>International Conference on Computational Linguis-</i>				680
	<i>tics</i> , pages 11196–11205, Abu Dhabi, UAE. Associa-				681
	tion for Computational Linguistics.				682
	Chenyang Zhang, Xiaoyu Zhang, Jian Lou, Kai Wu, Zi-				683
	long Wang, and Xiaofeng Chen. 2025. Poisonedeye:				684
	Knowledge poisoning attack on retrieval-augmented				685
	generation based large vision-language models . In <i>Inter-</i>				686
	<i>national Conference on Learning Representations</i>				687
	<i>(ICLR)</i> .				688
	Wei Zou, Runpeng Geng, Binghui Wang, and Jinyuan				689
	Jia. 2024. Poisonedrag: Knowledge corruption at-				690
	tacks to retrieval-augmented generation of large lan-				691
	guage models . <i>Preprint</i> , arXiv:2402.07867.				692

693	A Implementation Details	
694	Prompt Decomposition. We decompose each input prompt into semantic chunks corresponding to entities, attributes, actions, and locations using SpaCy-based dependency parsing. Specifically, we extract noun phrases, verb phrases, and associated modifiers using the <code>en_core_web_lg</code> model. Stop-words and purely functional tokens are removed, and duplicate chunks are merged to form a compact chunk set \mathcal{C} . This decomposition is performed once per prompt and incurs negligible overhead.	
695		
696		
697		
698		
699		
700		
701		
702		
703		
704	Retrieval Backbone. The retrieval module is based on a CLIP-style dual encoder, where text chunks are encoded using a pretrained text encoder E_T and images are encoded using a pretrained vision encoder E_I . For each chunk, we retrieve the top- k nearest images, with k set to 5 in all experiments unless otherwise specified.	
705		
706		
707		
708		
709		
710		
711	Diffusion Model Conditioning. We use a pre-trained Stable Diffusion Model as the image generator. The diffusion backbone remains frozen in all experiments. Verified retrieved images are injected as conditioning signals via the same retrieval-conditioning interface used by the baseline Vision RAG pipeline. No additional fine-tuning or parameter updates are performed on the diffusion model.	
712		
713		
714		
715		
716		
717		
718		
719	Poisoning Setup. To simulate retrieval poisoning, we corrupt the retriever by injecting mislabeled image–caption pairs into the retrieval corpus. We use a subset of the Conceptual Captions dataset containing approximately 3 million samples. Specifically, we inject 500 poisoned samples where images of one concept (e.g., <i>dog</i>) are paired with captions of a conflicting concept (e.g., <i>cat</i>). The diffusion model is kept frozen to isolate the effect of retrieval poisoning. We use Stable Diffusion v1.4 as the diffusion backbone.	
720		
721		
722		
723		
724		
725		
726		
727		
728		
729		
730	Hyperparameters and Runtime. All experiments are conducted on NVIDIA A100 GPUs. We train the retriever model from scratch using three NVIDIA A100 GPUs. Training is conducted for 64 epochs with a batch size of 128. We use the AdamW optimizer with an initial learning rate of 5×10^{-4} , followed by cosine learning rate scheduling and 10,000 warm-up steps.	
731		
732		
733		
734		
735		
736		
737		
738	Unless otherwise stated, we use $k = 5$, $\alpha = 1.0$	
739	for semantic coherence scoring.	
	B Computational Complexity Analysis	740
	We analyze the computational overhead of the proposed score-based semantic coherence refinement in comparison to the standard Vision RAG.	741
		742
		743
	Let m denote the number of semantic chunks extracted from a prompt, k the number of retrieved images per chunk, d the embedding dimension, and $ \mathcal{D} $ the size of the retrieval database. In a standard Vision RAG pipeline, retrieval and similarity computation incur a cost of	744
		745
		746
		747
		748
		749
		750
		751
		752
	assuming approximate nearest neighbor search for retrieval.	753
	Our method introduces an additional contrastive semantic coherence scoring step, which evaluates each retrieved image against the intended chunk and a small set of competing chunks. This adds a cost of	754
		755
		756
		757
		758
		759
	where $ \mathcal{N}_i $ denotes the number of negative chunks, bounded by the number of prompt components and typically small in practice.	760
		761
	For chunks identified as semantically inconsistent, a lightweight refinement and re-retrieval step is applied. In the worst case, all chunks are refined, resulting in an additional retrieval cost of	762
		763
		764
		765
		766
		767
		768
	However, in practice, refinement is triggered only for a small subset of chunks.	769
	Overall, the proposed method preserves the asymptotic complexity of Vision RAG:	770
		771
		772
		773
		774
		775
		776
		777
		778
		779
		780
		781
		782
		783
		784
	Discussion. The diffusion backbone remains unchanged, and no additional denoising steps are introduced. Since coherence verification operates only on top- k retrieved candidates and re-retrieval is applied conditionally, the overall computational cost remains dominated by retrieval and diffusion. This ensures that the proposed framework scales efficiently to large-scale vision retrieval databases while significantly improving robustness and semantic correctness.	