ReLoop: "Seeing Twice and Thinking Backwards" via Closed-loop Training to Mitigate Hallucinations in Multimodal understanding

Jianjiang Yang^{1*}, Yanshu Li^{2*}, Ziyan Huang^{3*}

¹University of Bristol, ²Brown University, ³South China University of Technology edisonyang109@gmail.com, yanshu_li1@brown.edu, bonnie.ziyan.huang@gmail.com

Abstract

While Multimodal Large Language Models (MLLMs) have achieved remarkable progress in open-ended visual question answering, they remain vulnerable to hallucinations. These are outputs that contradict or misrepresent input semantics, posing a critical challenge to the reliability and factual consistency. Existing methods often rely on external verification or posthoc correction, lacking an internal mechanism to validate outputs directly during training. To bridge this gap, we propose ReLoop, a unified closed-loop training framework that encourages multimodal consistency for cross-modal understanding in MLLMs. ReLoop adopts a ring-shaped structure that integrates three complementary consistency feedback mechanisms, obliging MLLMs to "seeing twice and thinking backwards". Specifically, ReLoop employs the frozen Consistency Feedback Plugin (CFP), comprising semantic reconstruction and visual description modules, along with an attention supervision module for attention alignment. These components collectively enforce semantic reversibility, visual consistency, and interpretable attention, enabling the model to correct its outputs during training. Extensive evaluations and analyses demonstrate the effectiveness of ReLoop in reducing hallucination rates across multiple benchmarks, establishing a robust method for hallucination mitigation in MLLMs. The code is available at: https://github.com/ZiyanHuang11/Reloophallucinations.

1 Introduction

In recent years, MLLMs (Liu et al., 2023b; OpenAI, 2023; Li et al., 2023a) have demonstrated significant progress in bridging vision and language, addressing tasks such as visual question answering (VQA), image captioning, and instruction adherence. However, a fundamental difficulty that

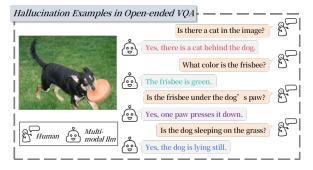


Figure 1: Illustration of four major hallucination types in open-ended VQA. Despite being visually grounded, MLLMs produce fluent but hallucinated responses across object, attribute, relation, and event dimensions.

persists is hallucination, where the generation of outputs that are inconsistent with or unsupported by visual inputs (Kalavasis et al., 2024). Hallucinations are especially common in open-ended VQA circumstances, in which unclear or underspecified questions can result in factual mistakes. These hallucinations span diverse categories, including Object, Attribute, Relation, and Event. Figure 1 illustrates that a singular image of "a dog grasping an orange frisbee" can elicit various forms of hallucination: a fictitious "cat" (object), an incorrectly identified "green" frisbee (attribute), an erroneous spatial relation "under the paw" (relation), or a temporal misrepresentation "sleeping" (event). These errors are semantically plausible yet visually unfounded, posing major challenges for trustworthiness and safety of MLLMs across critical applications, including medical decision-making (Kim et al., 2025), robotic perception (Park et al., 2023), and autonomous navigation (Alsulaimawi, 2025).

Existing works (Sun et al., 2023; Ayala and Béchard, 2024; Sun et al., 2024) often regard hallucination as an output-level anomaly that is corrected post hoc, overlooking its underlying cause. In practice, hallucinations frequently arise from misalignment between the input, visual content,

^{*}Equal contribution

and the model's latent reasoning. Without an internal supervision mechanism, models may produce fluent yet ungrounded answers. We argue that hallucination stems from the model's inability to validate its own output across modalities and recommend injecting this ability directly into training.

We subsequently derive inspiration from human cognitive processes. When answering visual questions, individuals rarely rely on a single forward guess. Instead, after answering, they may reassess the question's intent, examine the visual scene, and refine conclusions—especially in the face of ambiguity or uncertainty. However, most models operate in a unidirectional manner, mapping $(Q,I \to A)$. As a result, once the model makes a prediction, there is no structured way to assess whether it actually understood the question, if the answer aligns with the visual evidence, or whether the model attended to the right regions in the image.

To address this issue, we propose ReLoop, a cognitively inspired unified training framework that encourages multimodal consistency for crossmodal understanding in MLLMs. ReLoop implements a feedback-driven closed-loop supervision process, allowing the model to reassess its predictions and validate their consistency with the original input through multi-level supervision during training. Specifically, after MLLMs produce an answer from the image-question pair, Reloop enables the model to recapture input semantics and assess internal consistency via: a Consistency Feedback Plugin (CFP), comprising two frozen modules: (1) CFP-Lang reconstructs the question \hat{Q}^* from (A, I)to supervise semantic alignment, and (2) CFP-Vis generates a description I^* to assess factual grounding. In parallel, an attention supervision module extracts the model's token-to-image attention map \mathcal{H} and compares it with an entropy-based pseudoground truth. All signals are integrated as differentiable losses in the overall optimization objective. This design encourages the model to "see twice and think backward"—first look to answer $(Q, I \rightarrow A)$, look twice to reassess $(A, I \rightarrow \hat{Q}^*, I^*, \mathcal{H})$, and finally to correct $(\hat{Q}^*, I^*, \mathcal{H} \approx Q, I, \mathcal{H}_{\text{pseudo}})$.

ReLoop bridges the gap between perception and output. It turns the black-box understanding process into an interpretable, feedback-aware loop that continuously refines the model's internal representations. Our key contributions can be summarized clearly as follows:

• We propose **ReLoop**, a cognitively inspired

- closed-loop training framework that ensures consistency among image, question, and answer modalities, effectively mitigating hallucinations in MLLMs.
- We introduce three complementary consistency signals: semantic reconstruction, visual description, and attention alignment, to emulate the humanlike "reversible thinking" process and improve cross-modal consistency during training.
- We provide a novel use of pretrained visionlanguage models by repositioning them as frozen Consistency Feedback Plugins (CFPs) in the training loop. Rather than functioning as typical forward-only encoders, they now perform in a reflective, backward supervisory role, producing feedback signals to guide the main model's alignment with multimodal semantics.

2 Related Work

Hallucination Mitigation in MLLMs. Multimodal LLMs frequently produce hallucinations: responses conflicting with visual inputs, such as inventing entities or misaligning semantics (Li et al., 2023b). Recent mitigation efforts combine post-hoc correction and architectural refinement. Retrieval-augmented methods like (Mala et al., 2025) grounds outputs in external knowledge via hybrid retrievers, while (Ayala and Béchard, 2024) reduces hallucinations in structured outputs. Architectural solutions such as OPERA (Huang et al., 2024) penalize over-trust during decoding, and preference-aligned training like TPO (Gu and Wang, 2025) enhances vision grounding. Postgeneration verification (Woodpecker (Yin et al., 2023)) and decoding-time personalization (PAD (Chen et al., 2024)) complement training-time alignment; concurrent work improves multimodal ICL efficiency and task control (Li et al., 2025a,b) and instruction-tuned dialog grounding (Luo et al., 2024b). Beyond images, video hallucination is diagnosed via fine-grained spatio-temporal grounding (Luo et al., 2025).

Semantic Reversibility and Bidirectional Supervision. Human cognition leverages bidirectional reasoning to validate hypotheses: a principle termed "cognitive reversibility" (Johnson-Laird, 1983). Recent works explore this idea through decoding-time strategies: Self-RAG (Asai et al., 2023) integrates retrieval-augmented generation with self-reflection, enabling models to critique

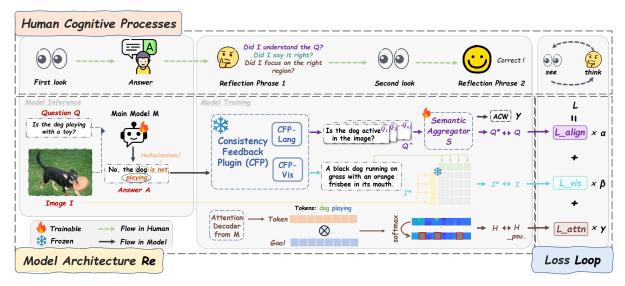


Figure 2: Seeing Twice and Thinking Backwards: ReLooping Hallucination Suppression in Multimodal Language Models. This diagram aligns human cognitive phases (left) with model modules (right) in a closed-loop process. The main model M produces an answer which is then introspected via CFP-Lang (language reconstruction), CFP-Vis (visual description), and internal cross-attention maps. Semantic aggregation, CLIP similarity, and entropy-based soft masks produce feedback losses that are summed and back-propagated to update M and the semantic aggregator S.

and refine their outputs iteratively, while DeepSeek-Math employs Group Relative Policy Optimization (GRPO) (Shao et al., 2024), enhancing mathematical reasoning by optimizing policy decisions based on group sampling strategies. Similarly, back-translation methods (Sennrich et al., 2016) enforce answer-question consistency through round-trip translation.

Cross-modal Consistency. Ensuring modal consistency is vital for mitigating hallucinations in MLLMs. Recent methods enhance visual-text alignment to reduce semantic drift. VCD (Leng et al., 2024) contrasts outputs from original and perturbed images to promote grounding and reduce unimodal bias. (Jiang et al., 2024) treats hallucinated captions as hard negatives to improve alignment. EAGLE (Villa et al., 2025) further refines visual encoders post-pretraining, yielding better grounding and fewer hallucinations. Broader benchmark efforts catalog the MLLM evaluation landscape (Li et al., 2024), including conversational aspect-based sentiment settings (Luo et al., 2024a).

3 Preliminaries

3.1 Task Formulation: Open-ended Visual Question Answering

We consider the task of open-ended VQA, where the model receives an image I and a natural language question Q, and produces a free-form answer A. Unlike multiple-choice settings, this task requires the model to produce linguistically coherent and visually grounded responses without predefined options.

In this case, hallucination refers to answers that contradict the image I, misinterpret the question Q, or introduce unsupported content.

3.2 Consistency Signals

To encourage faithful understanding, we supervise the model using three types of cross-modal consistency signals:

Linguistic Consistency. We verify whether the model's answer A implies the same question intent as the original Q, by attempting to reconstruct Q from (A, I). This tests whether the model understood the question meaningfully.

Visual Consistency. We evaluate whether the answer A is factually grounded in image I, by generating a descriptive caption I^* based on (A, I) and checking its alignment with the image. This ensures that the response reflects the actual visual content.

Attention Consistency. We examine whether the model attends to the correct regions of the image while producing A. This is assessed by comparing its internal attention map \mathcal{H} with a soft pseudoground truth \mathcal{H}_{pseudo} derived from entropy-based cues.

Together, these consistency signals serve as indirect evidence of whether the model truly grasps both the visual input and the question semantics.

4 ReLoop Framework: Reflect, Recapture, and Optimize through a Closed-Loop Process

We introduce **ReLoop**, a unified training framework aimed at reducing hallucinations in MLLMs for open-ended VQA answering. As illustrated in Figure 2, the framework incorporates three complementary consistency feedback mechanisms: **semantic reconestruction**, **visual description**, and **attention alignment** to supervise the model toward producing answers faithful to both the question and the image.

These feedback signals are instantiated through a frozen Consistency Feedback Plugin (CFP): semantic reconstruction (CFP-Lang) and visual description (CFP-Vis), and attention supervision from the model itself. The CFP module is broadly compatible with a range of encoder-decoder or decoder-only MLLMs. During inference (*First Look* \rightarrow *Answer*), the model receives a question-image pair and produces an initial answer. The training process then begins with *Reflect* \rightarrow *Second Look* \rightarrow *Correct*: the model examines its output through structured consistency feedback. Specifically, it "introspectively" asks:

- "Did I understand the Q?" (\rightarrow semantic reconstruction)
- "Did I say it right?" (→ visual description)
- "Did I focus on the right region?" (→ attention alignment)

ReLoop decomposes hallucination mitigation into two interacting components:

- "Re" emphasizes recapturing details, encouraging the model to reassess the semantic and visual cues from both question and image through CFP modules and token-level attention heatmaps.
- "Loop" denotes a feedback-driven training loop. After each forward prediction, feedback from the three consistency pathways is aggregated into the loss function ($L_{\rm align}$, $L_{\rm vis}$, $L_{\rm attn}$), driving iterative updates that refine the model's multimodal grounding and answer reliability.

4.1 A Closed-loop Training

The entire training process follows a closed-loop pattern, emulating "seeing twice and thinking backward". Each training step proceeds as follows:

- First Look: The main model M takes the image I and question Q as input to produce an initial answer A.
- 2. **Reflect:** The model introspects on A by reconstructing a proxy question \hat{Q} , generating a visual description I^* , and extracting token-level attention \mathcal{H} .
- Second Look: The reconstructions are compared against the original inputs to compute consistency losses, capturing discrepancies in semantics, visual grounding, and attention focus.
- 4. **Correct:** All feedback signals are aggregated into L_{total} to update M and the semantic aggregator S via backpropagation.

This multi-stage loop is repeated across training epochs, leading to the model M that gradually reduces hallucinations.

4.2 Re: Recapturing Details for Consistency Supervision

This stage corresponds to the training-time processes of "Reflect" and "Second Look", where the model reassesses its answers to recapture overlooked semantic and visual details. Three feedback pathways modules examine whether the model understood the question, correctly grounded its answer in the image, and attended to salient regions.

4.2.1 CFP-lang: Language Reconstruction and Adaptive Consistency Weighting

To evaluate whether the model correctly interprets the input question, we introduce a frozen language reconstruction module, CFP-lang. Given the answer-image pair (A,I), CFP-lang produces a set of candidate reverse questions $\{\hat{Q}_1,\hat{Q}_2,\ldots,\hat{Q}_k\}$ that approximate possible intents underlying the predicted answer. A lightweight semantic aggregator S, composed of a BERT encoder and a single-layer MLP, scores each candidate against the original question Q using BERTScore. The highest-ranked proxy \hat{Q}^* is selected to reflect the model's inferred intent.

However, directly enforcing alignment on all reconstructed questions may introduce noise, particularly when the produced answer is short or underspecified. To mitigate this, we introduce an Adaptive Consistency Weighting (ACW) mechanism, which adjusts the attention supervision (mentioned in section 4.2.3) strength based on the similarity

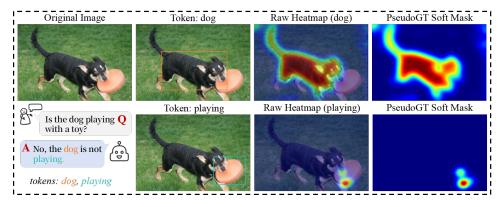


Figure 3: **Token-Level Attention Supervision.** Visualization of predicted attention \mathcal{H} and entropy-based pseudo ground truth \mathcal{H}_{pseudo} for two key answer tokens: dog (top row) and playing (bottom row).

between Q and \hat{Q}^* :

$$\gamma = \begin{cases} 1.0 & \text{if BERTScore}(Q, \hat{Q}^*) \geq 0.8 \\ 0.1 & \text{if } 0.6 \leq \text{BERTScore}(Q, \hat{Q}^*) < 0.8 \\ 0.01 & \text{if BERTScore}(Q, \hat{Q}^*) < 0.6 \end{cases} \tag{1}$$

Rather than discarding low-confidence pairs, this soft weighting ensures that stronger semantic matches contribute more prominently to the learning objective. The language consistency loss is computed as:

$$L_{\text{align}} = 1 - \text{BERTScore}(Q, \hat{Q}^*)$$
 (2)

4.2.2 CFP-visual: Visual Description and Similarity Supervision

To validate whether the produced answer A is visually grounded in the image I, we employ a frozen visual description module, CFP-visual. Given (A,I), it generates a caption I^* describing the image content implied by the answer. We then compute the cosine similarity between the CLIP-encoded vectors of I and I^* , and derive the visual consistency loss as:

$$L_{\text{vis}} = 1 - \cos(\text{CLIP}_{\text{img}}(I), \text{CLIP}_{\text{text}}(I^*))$$
 (3)

4.2.3 Attention Supervision via Heatmap Consistency

To enhance interpretability and mitigate hallucinations arising from inattentive or unstable decoding, we explicitly supervise the model's token-level cross-attention patterns. From the decoder of the main model M, we extract attention maps \mathcal{H} , which indicate the spatial focus during answer generation. We construct a soft pseudo-ground-truth heatmap $\mathcal{H}_{\text{pseudo}}$ using entropy-based masking (Detailed explanation can be found in Appendix C).

This method preserves uncertainty information and avoids brittle hard labels. As illustrated in Figure 3, well-grounded tokens (e.g., dog) yield concentrated heatmaps aligned with visual evidence, while hallucinated tokens (e.g., playing) produce offset patterns. We enforce alignment between \mathcal{H} and \mathcal{H}_{pseudo} by minimizing the KL divergence:

$$L_{\text{attn}} = \text{KL}(\mathcal{H} \parallel \mathcal{H}_{\text{pseudo}}) \tag{4}$$

4.3 Loop: Feedback Aggregation, Alignment, and Optimization

After consistency signals are computed from language, vision, and attention supervision, ReLoop aggregates them into a unified training objective. This stage corresponds to the "Correction" step in the loop, where the model updates its parameters based on multi-perspective feedback. The total loss combines standard supervision with the three consistency terms:

$$L_{\text{total}} = L_{\text{sft}} + \alpha \cdot L_{\text{align}} + \beta \cdot L_{\text{vis}} + \gamma \cdot L_{\text{attn}} + \lambda \cdot \Omega(\theta)$$
(5)

where $L_{\rm sft}$ is the token-level cross-entropy loss, and $\Omega(\theta)$ is an L2 regularization term. The consistency weights are empirically set as $\alpha=1.0$, $\beta=0.7$, $\lambda=10^{-5}$ and γ is defined in Equation 1.

Only the parameters of the main model M and the semantic aggregator S are updated during training. All feedback modules, including CFP-Lang, CFP-Vis, attention supervision, and CLIP, remain frozen.

5 Experimental Setup

Training Data. We curate 30K high-quality $\{I, Q, A\}$ from LLaVA-Instruct-150K. To simulate hallucination supervision, we generate contrastive examples by perturbing key semantics (e.g., *objects*,

| Туре | Module | Signal Type | Baseline | ReLoop | Δ Mean | Baseline Hallu. | ReLoop Hallu. | Δ Rate |
|-----------|---------------------------------|--|--|--|--------------------------|--------------------|------------------|---------------|
| Object | Visual Language Attention | $\begin{array}{c} \operatorname{CLIP}(I,I^*) \\ \operatorname{BERT}(Q,\hat{Q}) \\ \operatorname{Entropy}(\mathcal{H}) \end{array}$ | 28.02 ± 3.10 0.862 ± 0.022 1.31 ± 0.40 | 29.46 ± 3.27 0.873 ± 0.024 1.28 ± 0.45 | ↑1.44 ↑0.011 ↓0.03 | 24.5% | 10.3% | ↓14.2% |
| Attribute | Visual Language Attention | $	ext{CLIP}(I, I^*)$ $	ext{BERT}(Q, \hat{Q})$ $	ext{Entropy}(\mathcal{H})$ | 26.59 ± 3.31 0.868 ± 0.025 1.36 ± 0.46 | $26.81 \pm 3.41 \\ 0.894 \pm 0.028 \\ 1.32 \pm 0.52$ | ↑0.22 ↑0.026 ↓0.04 | 7.3% | 4.0% | ↓3.3% |
| Relation | Visual Language Attention | $\begin{array}{c} \operatorname{CLIP}(I,I^*) \\ \operatorname{BERT}(Q,\hat{Q}) \\ \operatorname{Entropy}(\mathcal{H}) \end{array}$ | 27.22 ± 3.26 0.855 ± 0.020 1.39 ± 0.43 | 28.01 ± 3.38 0.875 ± 0.023 1.34 ± 0.50 | ↑0.79 ↑0.020 ↓0.05 | 13.2% | 7.6% | ↓5.6% |
| Event | Visual Language Attention | $\begin{array}{c} \operatorname{CLIP}(I,I^*) \\ \operatorname{BERT}(Q,\hat{Q}) \\ \operatorname{Entropy}(\mathcal{H}) \end{array}$ | 26.63 ± 3.08 0.861 ± 0.024 1.33 ± 0.42 | $26.94 \pm 3.37 \\ 0.877 \pm 0.029 \\ 1.51 \pm 0.55$ | ↑0.31 ↑0.016 ↑0.18 | 10.4% | 5.2% | ↓5.2% |

Table 1: Effect of ReLoop on consistency and hallucination reduction across different hallucination types. We compare MiniGPT-4 (baseline) and ReLoop in terms of signal outputs from three frozen feedback modules: visual grounding (CLIP similarity), semantic alignment (BERTScore), and attention focus (entropy). Δ denotes the absolute change in signal quality after applying ReLoop.

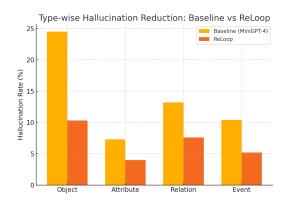


Figure 4: Type-wise hallucination rates (%) for baseline (MiniGPT-4) and ReLoop models.

attributes, relations, event), enabling fine-grained control over hallucination types. Details can be found in Appendix A.1 A.2.

Evaluation Benchmarks and Metrics. We evaluate ReLoop on a broad range of hallucination and multimodal understanding benchmarks, including POPE (Li et al., 2023b), CHAIR (Rohrbach et al., 2018), AMBER (Wang et al., 2023), MMHal-B (Sun et al., 2023), HallusionBench (Guan et al., 2024), Faith/FaithS (Jing et al., 2024), and MME (Fu et al., 2023). Full definitions can be found in Appendix A.3 A.4.

Baselines. We use MiniGPT-4 as the baseline model in Experiment 6.1 and compare against LLaVA-1.5 variants trained with LLaVA-RLHF (Sun et al., 2023), HA-DPO (Zhao et al., 2023), POVID (Zhou et al., 2024), and Visual

Contrastive Decoding (VCD) (Sicong Leng, 2023). For robustness analyses, we adopt *LLaVA-1.5* + *ReLoop* as the canonical setting and report stress tests under noisy external supervision and nonsensical answers (Table 5). Unless otherwise specified, all baselines share the same backbone, data, and training protocol for a fair comparison. Implementation details are provided in Appendix A.5.

6 Results and Analysis

6.1 Identify Internal Causes of Hallucinations: Module Signals vs. Hallucination States

We first aim to pinpoint internal representation deficiencies that drive hallucination behaviors across different hallucination types. We analyze consistency signal deviations produced by ReLoop's frozen supervision modules, with hallucinated versus non-hallucinated samples. Responding: "Did I understand the question?" (language, via BERTScore); "Did I say it right?" (visual, via CLIP similarity); "Did I focus on the right region?" (attention, via entropy).

Multimodal hallucinations stem from structured, modality-specific representation gaps. As shown in Table 1, hallucinated responses are consistently associated with lower CLIP similarity (–2.25), reduced BERTScore (–0.034), and higher attention entropy (+0.31). Figure 5 reveals distinct signal patterns associated with different hallucination types. Object hallucinations correspond to a clear leftward shift in CLIP similarity, indicat-

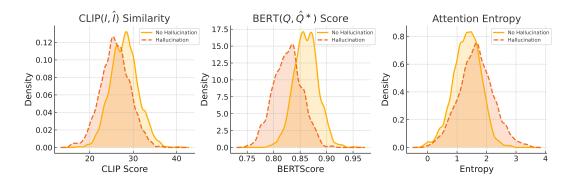


Figure 5: KDE distributions of CLIP similarity, BERTScore, and attention entropy for hallucinated and non-hallucinated samples. ReLoop's frozen modules exhibit sharp signal shifts that serve as reliable supervision sources.

| Model | Hallu | cination Supp | ression | Cross | s-modal Faith | fulness |
|-----------------------|---------------------|---------------------|---|---------------------|------------------|------------------|
| Wiodei | POPE↑ | $CHAIR_s\downarrow$ | $\overline{\operatorname{CHAIR}_i\downarrow}$ | F1↑ | Faith↑ | FaithS↑ |
| MiniGPT-4 | 82.3 | 49.0 | 22.7 | 63.2 | 86.7 | 68.5 |
| + ReLoop | 83.9 | 38.8 | 20.5 | 69.9 | 88.6 | 71.3 |
| InstructBLIP | 83.8 | 47.8 | 20.6 | 68.4 67.0 | 87.3 | 69.8 |
| + ReLoop | 85.3 | 36.9 | 17.5 | | 88.5 | 73.2 |
| LLaVA-1.5 | 85.7 | 53.5 | 24.2 | 65.8 | 89.5 | 75.8 75.3 |
| + ReLoop | 86.3 | 40.2 | 16.2 | 70.3 | 89.2 | |
| LLaVA-1.6 | 86.8 | 52.0 | 21.8 | 67.3 | 89.4 | 76.6 76.2 |
| + ReLoop | 87.9 | 38.5 | 16.1 | 71.1 | 89.2 | |
| Qwen-VL-2.5 | 89.3 | 47.3 | 20.8 | 69.1 | 89.7 | 76.5 |
| + ReLoop | 90.7 | 37.6 | 16.6 | 72.5 | 90.4 | 77.8 |
| mPLUG-owl + ReLoop | 89.1 90.9 | 62.5 42.5 | 31.0 21.8 | 58.9 66.5 | 88.3 87.9 | 72.7 71.0 |
| ShareGPT4V | 88.2 | 50.2 | 21.8 | 68.0 | 88.2 | 73.6 |
| + ReLoop | 89.7 | 44.9 | 21.5 | 69.2 | 89.3 | 74.8 |

Table 2: Performance comparison of various LVLMs with and without ReLoop. Hallucination is measured by POPE, CHAIR_s, and CHAIR_i; cross-modal faithfulness is evaluated using F1, Faith, and FaithS. \downarrow indicates lower is better; \uparrow indicates higher is better.

ing weaker visual grounding. Attribute hallucinations are marked by lower BERTScore, reflecting reduced semantic alignment. Event hallucinations show higher attention entropy, suggesting that the model distributes focus more broadly, which may help in capturing complex scenes but also increases the risk of focusing on irrelevant regions.

Signal dynamics vary by hallucination type. (1) *Object hallucinations* are primarily rooted in the visual module. They often manifest as hallucinated entities not present in the image. ReLoop yields a significant gain in CLIP similarity (\uparrow 1.44) and a decrease in attention entropy (\downarrow 0.03), suggesting enhanced image-text alignment and focused visual grounding. (2) *Attribute hallucinations* show the largest improvement in BERTScore (\uparrow 0.026) and only a slight change in CLIP similarity (\uparrow 0.22), indicating that semantic reconstruction plays a more important role than visual grounding. This aligns

with their nature: attributes often relate to textual misinterpretation (e.g., color or size), even when visual cues are present. (3) Relation hallucinations involve complex spatial or relational semantics and display moderate improvements across all three signals (CLIP \(^1\)0.79, BERT \(^1\)0.020, Entropy \(\lambda 0.05 \), suggesting that ReLoop's multi-signal supervision addresses cross-modal misalignment collaboratively. (4) Event hallucinations are primarily tied to attention misallocation. ReLoop improves CLIP ($\uparrow 0.31$) and BERT ($\uparrow 0.016$) slightly, but entropy increases ($\uparrow 0.18$), reflecting broader attention scopes. This likely helps avoid fixation on irrelevant regions, especially in dynamic or temporally inferred scenes. Figure 4 shows that ReLoop successfully mitigates hallucinations compared to MiniGPT-4 across four hallucination types.

| Ablation Version | Hallucination Suppression | | | Cross-modal Faithfulness | | |
|-----------------------------|---------------------------|---------------------|---|--------------------------|--------|---------------------|
| 11011011 (0101011 | POPE↑ | $CHAIR_s\downarrow$ | $\overline{\operatorname{CHAIR}_i\downarrow}$ | F1↑ | Faith† | FaithS [↑] |
| MiniGPT-4 | 83.0 | 49.0 | 22.7 | 60.2 | 84.3 | 64.2 |
| w/o Consistency Supervision | 84.2 | 47.4 | 21.6 | 60.7 | 86.7 | 68.5 |
| w/o Gating & Aggregator | 85.4 | 39.8 | 19.7 | 60.4 | 88.1 | 71.6 |
| w/o Attention Supervision | 83.6 | 40.2 | 20.1 | 61.9 | 86.3 | 67.5 |
| Full ReLoop | 84.9 | 38.3 | 18.9 | 63.1 | 88.6 | 72.8 |

Table 3: Performance comparison of ReLoop under different ablation configurations on MiniGPT-4. Removing consistency supervision results in the worst faithfulness and hallucination rate, while full ReLoop delivers the best overall performance. Although gating removal slightly improves POPE, it hurts precision (F1) and consistency.

| Method | Hallu | Hallucination Suppression | | | Cross-modal Faithfulness | | |
|--------------|-------------|---------------------------|---|-------------------|--------------------------|---------------------|--|
| Memou | POPE↑ | $CHAIR_s\downarrow$ | $\overline{\operatorname{CHAIR}_i\downarrow}$ | F1↑ | Faith↑ | FaithS [↑] | |
| LLaVA-1.5 | 83.5 | 53.9 | 23.5 | 63.2 | 86.9 | 70.5 | |
| + LLaVA-RLHF | 88.2 | 44.5 | 20.1 | 67.0 | 89.0 | 74.4 | |
| + HA-DPO | 86.7 | 52.3 | $\frac{20.1}{21.6}$ | $\overline{65.4}$ | $\frac{89.0}{88.4}$ | $\frac{74.4}{73.5}$ | |
| + POVID | 84.3 | 53.2 | 24.2 | 64.7 | 87.3 | 71.8 | |
| + VCD | 86.8 | <u>43.1</u> | 20.2 | 66.9 | 88.8 | 73.6 | |
| + ReLoop | <u>87.9</u> | 42.0 | 19.5 | 67.4 | 89.5 | 75.1 | |

Table 4: Performance comparison of ReLoop with alignment-enhancing baselines for LLaVA-1.5 on hallucination suppression and cross-modal faithfulness. Best scores are in **bold** and the second best are <u>underlined</u>.

6.2 Effects of Structured Feedback in ReLoop

Motivated by earlier findings, we evaluate how effectively ReLoop's structured feedback enhances semantic grounding across five representative LVLMs (Table 2). The observed improvements span models with diverse architectures and training paradigms, showing that ReLoop is broadly compatible and easily integrable into various LVLMs.

Hallucination Suppression. ReLoop significantly reduces references to non-existent entities. InstructBLIP shows 22.8%/15.0% reductions. LLaVA-1.5 improves by 24.9%/33.1%, and strong backbones exhibit the same trend: LLaVA-1.6 achieves $\sim\!26\%/\sim\!26\%$ drops, while Qwen-VL-2.5 yields $\sim\!20.5\%/\sim\!20.2\%$. Similar effects hold for mPLUG-owl and ShareGPT4V. These reductions confirm that ReLoop enhances visual grounding and spatial precision across backbones.

Cross-modal Faithfulness. ReLoop also enhances cross-modal faithfulness. F1 increases on MiniGPT-4 (+10.6%), LLaVA-1.5 (+6.8%), LLaVA-1.6 (~5.6%), and Qwen-VL-2.5(~4.9%); InstructBLIP maintains comparable F1 while gaining on faith metrics. FaithS improves for MiniGPT-4 (+2.8), InstructBLIP (+3.4), Qwen-VL-2.5 (+1.3), and ShareGPT4V (+1.2), and remains near-parity on LLaVA-1.6. These gains suggest that the model not only grounds responses more accurately in the image but also maintains semantic alignment with the question intent.

6.3 Robustness Under Noisy Supervision

We next ask whether closed-loop training remains stable when external supervision is imperfect or when the initial answer signal is degenerate. We stress-test ReLoop under (i) noisy teacher feedback on the visual channel and (ii) nonsensical answers that could mislead the loop.

Closed-loop supervision is resilient to noisy teachers. We corrupt 15% of visual descriptions fed to the CFP-Vis branch with pseudo-random text. As shown in Table 5, the average CLIP similarity decreases (Δ CLIP = -0.11), yet core hallucination metrics remain stable (POPE 87.9 \rightarrow 86.8, CHAIR-s 42.0 \rightarrow 42.6, CHAIR-i 19.5 \rightarrow 20.2). This suggests that multi-signal aggregation dilutes spurious teacher cues, preserving cross-modal consistency even when the visual supervisor is noisy.

ACW suppresses nonsensical answers without destabilizing training. We replace 15% of answers with meaningless strings while keeping ACW active. Table 5 shows a predictable reallocation of per-sample weights: the mass at high confidence shrinks (γ =1.0: 52% \rightarrow 35%), medium/low weights grow (0.1/0.01: 41/9% \rightarrow 43/22%), and semantic alignment only mildly drops (Δ BERTScore= -0.06), while POPE/CHAIR remain essentially unchanged. This confirms that ACW down-weights misleading answer signals before they influence learning.

| Setting | POPE ↑ | CHAIR-s ↓ | CHAIR-i↓ | $\Delta \text{CLIP-sim} \ (\uparrow)$ | γ dist. (1.0/0.1/0.01) | $\Delta \mathbf{BERTScore}~(\uparrow)$ |
|----------------------------|---------------|------------------|----------|---------------------------------------|-------------------------------|--|
| LLaVA-1.5 + ReLoop (clean) | 87.9 | 42.0 | 19.5 | _ | 52/41/9 | <u> </u> |
| + Teacher noise (15%) | 86.8 | 42.6 | 20.2 | -0.11 | _ | _ |
| + Answer noise (15%) | 87.2 | 43.1 | 20.4 | _ | 35/43/22 | -0.06 |

Table 5: Robustness under noisy supervision. Teacher-side visual description corruption (15%) and answer-side nonsense injection (15%) have limited impact on core hallucination metrics. CLIP similarity drops under teacher noise, whereas ACW re-allocates per-sample weights under answer noise (smaller high-confidence mass).

Natural noise robustness via fourfold filtering.

ReLoop's robustness emerges from a fourfold filter: three orthogonal supervision signals (language via BERTScore, visual via CLIP, attention via entropy-aware \mathcal{H}_{pseudo}) plus ACW's discrete gating $\gamma \in \{1,0.1,0.01\}$ (Sec. 4.2). Noisy teacher feedback is first cross-validated across modalities, then attenuated by ACW, and finally diluted in the multiloss objective (Sec. 4.3), so that biased cues have limited influence on the update direction. Empirically, the stability of POPE/CHAIR under teacherside corruption and the expected shift in the γ distribution under answer nonsense together indicate that the closed loop self-regularizes and converges stably despite imperfect teachers.

6.4 Ablation Study

To assess the contribution of each component in ReLoop, we perform a coarse-grained ablation study over four configurations (Table 3). Removing consistency supervision leads to the highest hallucination rates (CHAIR_s: 47.4) and lowest semantic faithfulness (FaithS: 68.5), highlighting its central role. Attention supervision also proves important, as its removal moderately reduces FaithS. While removing gating slightly improves POPE, it harms F1 and hallucination suppression. Full ReLoop achieves the best overall results, reducing CHAIR_s by 10.7 and increasing FaithS by 8.6 over the baseline. These findings underscore the complementary roles of all modules and the importance of structured feedback for robust alignment.

6.5 Unified Comparison with Alignment Strategies

We compare ReLoop with representative alignment methods, LLaVA-RLHF, HA-DPO, and POVID on both fine-grained hallucination metrics and broader benchmark evaluations. As shown in Table 4, ReLoop consistently outperforms alternatives on POPE, CHAIR, F1, and faithfulness metrics, indicating stronger hallucination suppression and cross-modal faithfulness. On benchmark-level evalua-

| Method | AMBER↑ | MME↑ | MMHal-B↑ | Hallu-B↑ |
|--------------|--------|------|----------|----------|
| LLaVA-1.5 | 73.9 | 1513 | 65.4 | 48.6 |
| + LLaVA-RLHF | 73.8 | 1231 | 64.3 | 43.2 |
| + HA-DPO | 77.2 | 1374 | 65.6 | 49.9 |
| + POVID | 75.8 | 1421 | 65.9 | 51.4 |
| + ReLoop | 80.3 | 1505 | 68.9 | 52.3 |

Table 6: Benchmark-level comparison of ReLoop with alignment strategies across four evaluation baselines.

tions (Table 6), ReLoop leads on AMBER, MMHal-B, and HallusionBench, while remaining competitive on MME. The slight MME drop may reflect a common trade-off between alignment supervision and low-level perception, also observed in other alignment-based methods like LLaVA-RLHF. These findings underscore ReLoop's effectiveness across both targeted and comprehensive settings.

6.6 Additional Analyses

To further substantiate the effectiveness and practicality of ReLoop, we provide supplementary analyses in the appendices. Appendix A.6 reports a Training Cost Breakdown: cross-method cost, CFP overhead attribution, and convergence/epoch-level timing. Appendix A.7 ablates contrastive augmentation to isolate closed-loop gains. Appendix B provides a Case Study over four hallucination types and a nonsensical-answer failure mode, illustrating early rejection, ACW γ -downweighting, and entropy-aware masking.

7 Conclusion

We present **ReLoop**, a closed-loop training framework that mitigates hallucinations in MLLMs by enforcing semantic and visual consistency through bidirectional feedback. By incorporating language reconstruction, visual description, and attention alignment, ReLoop allows models to verify and refine predictions during training. Experiments show consistent gains in hallucination suppression and interpretability, establishing ReLoop as a general foundation for building more reliable MLLMs.

Potential Limitations

Performance Variability Across Hallucination

Types. While ReLoop substantially improves hallucination suppression in object and attribute categories, its effectiveness on relation and event hallucinations remains relatively modest. These hallucination types often involve higher-order reasoning and temporal or spatial understanding, which are less easily corrected through current consistency signals. Future extensions may incorporate specialized supervision tailored to relational semantics or causal cues to address this gap.

Supervision Dependency and Domain Adaptability. ReLoop relies on access to paired image—question—answer data to compute consistency signals. This requirement poses challenges in domains with limited high-quality supervision, such as medical or scientific imaging. Moreover, the training framework assumes reasonably clean and grounded reference answers, which may not hold in low-resource or noisy environments. Reducing ReLoop's dependence on strongly supervised inputs and exploring semi-supervised or synthetic feedback generation remain promising directions for broader applicability.

Scope of evaluation. Our study focuses on standard VQA-style image benchmarks and general-domain LVLMs. We have not evaluated text-heavy or long-tail domains (e.g., *dense OCR*, *charts*) or temporal reasoning tasks, where attention allocation and signal reliability may differ. Extending evaluation to these regimes is left for future work.

Ethics Statement

All datasets utilized in this work are either publicly released or ethically sourced, ensuring full compliance with associated data usage policies. For evaluation purposes, we additionally include AI-generated content produced under controlled prompting conditions. These samples are clearly labeled and subjected to careful human verification to ensure factual accuracy and annotation quality. We acknowledge the broader implications of hallucination mitigation in AI systems and advocate for responsible model development that prioritizes reliability, fairness, and interpretability.

References

Zahir Alsulaimawi. 2025. Feedback-enhanced hallucination-resistant vision-language model for

- real-time scene understanding. arXiv preprint arXiv:2504.04772.
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2023. Self-rag: Learning to retrieve, generate, and critique through self-reflection. *arXiv preprint arXiv:2310.11511*.
- Orlando Marquez Ayala and Patrice Béchard. 2024. Reducing hallucination in structured outputs via retrieval-augmented generation. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 6: Industry Track)*, pages 228–238. Association for Computational Linguistics.
- Ruizhe Chen, Xiaotian Zhang, Meng Luo, Wenhao Chai, and Zuozhu Liu. 2024. Pad: Personalized alignment at decoding-time. *arXiv e-prints*, pages arXiv–2410.
- Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, Yunsheng Wu, and Rongrong Ji. 2023. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv* preprint arXiv:2306.13394.
- Jihao Gu and Yingyao Wang. 2025. Token preference optimization with self-calibrated visual-anchored rewards for hallucination mitigation. *arXiv* preprint *arXiv*:2412.14487.
- Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, Dinesh Manocha, and Tianyi Zhou. 2024. Hallusionbench: An advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14375–14385. IEEE.
- Qidong Huang, Xiaoyi Dong, Pan Zhang, Bin Wang, Conghui He, Jiaqi Wang, Dahua Lin, Weiming Zhang, and Nenghai Yu. 2024. Opera: Alleviating hallucination in multi-modal large language models via over-trust penalty and retrospection-allocation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. Highlight Paper.
- Chaoya Jiang, Haiyang Xu, Mengfan Dong, Jiaxing Chen, Wei Ye, Ming Yan, Qinghao Ye, Ji Zhang, Fei Huang, and Shikun Zhang. 2024. Hallucination augmented contrastive learning for multimodal large language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 27036–27045. Code available at https://github.com/X-PLUG/mPLUG-HalOwl/tree/main/hacl.
- Liqiang Jing, Ruosen Li, Yunmo Chen, and Xinya Du. 2024. Faithscore: Fine-grained evaluations of hallucinations in large vision-language models. In *Findings of the Association for Computational Linguis*

- *tics: EMNLP 2024*, pages 5042–5063. Association for Computational Linguistics.
- Philip N. Johnson-Laird. 1983. *Mental Models: Towards a Cognitive Science of Language, Inference, and Consciousness*. Harvard University Press, Cambridge, MA.
- Alkis Kalavasis, Anay Mehrotra, and Grigoris Velegkas. 2024. On the limits of language generation: Tradeoffs between hallucination and mode collapse. *arXiv* preprint arXiv:2411.09642.
- Yubin Kim, Hyewon Jeong, Shan Chen, Shuyue Stella Li, Mingyu Lu, Kumail Alhamoud, Jimin Mun, Cristina Grau, Minseok Jung, Rodrigo Gameiro, Lizhou Fan, Eugene Park, Tristan Lin, Joonsik Yoon, Wonjin Yoon, Maarten Sap, Yulia Tsvetkov, Paul Liang, Xuhai Xu, and 6 others. 2025. Medical hallucinations in foundation models and their impact on healthcare. arXiv preprint arXiv:2503.05777.
- Sicong Leng, Hang Zhang, Guanzheng Chen, Xin Li, Shijian Lu, Chunyan Miao, and Lidong Bing. 2024. Mitigating object hallucinations in large vision-language models through visual contrastive decoding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1316–1325. Poster Highlight.
- Jian Li, Weiheng Lu, Hao Fei, Meng Luo, Ming Dai, Min Xia, Yizhang Jin, Zhenye Gan, Ding Qi, Chaoyou Fu, and 1 others. 2024. A survey on benchmarks of multimodal large language models. *arXiv* preprint arXiv:2408.08632.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023a. Blip-2: Bootstrapping language-image pretraining with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*.
- Yanshu Li, Yi Cao, Hongyang He, Qisen Cheng, Xiang Fu, Xi Xiao, Tianyang Wang, and Ruixiang Tang. 2025a. M²IV: Towards efficient and fine-grained multimodal in-context learning via representation engineering. In *Second Conference on Language Modeling*.
- Yanshu Li, Tian Yun, Jianjiang Yang, Pinyuan Feng, Jinfa Huang, and Ruixiang Tang. 2025b. Taco: Enhancing multimodal in-context learning via task mapping-guided sequence configuration. *arXiv* preprint arXiv:2505.17098.
- Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. 2023b. Evaluating object hallucination in large vision-language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023a. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023b. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*.

- Meng Luo, Hao Fei, Bobo Li, Shengqiong Wu, Qian Liu, Soujanya Poria, Erik Cambria, Mong-Li Lee, and Wynne Hsu. 2024a. Panosent: A panoptic sextuple extraction benchmark for multimodal conversational aspect-based sentiment analysis. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 7667–7676.
- Meng Luo, Shengqiong Wu, Liqiang Jing, Tianjie Ju, Li Zheng, Jinxiang Lai, Tianlong Wu, Xinya Du, Jian Li, Siyuan Yan, Jiebo Luo, William Yang Wang, Hao Fei, Mong-Li Lee, and Wynne Hsu. 2025. Dr.v: A hierarchical perception-temporal-cognition framework to diagnose video hallucination by fine-grained spatial-temporal grounding. *Preprint*, arXiv:2509.11866.
- Meng Luo, Han Zhang, Shengqiong Wu, Bobo Li, Hong Han, and Hao Fei. 2024b. Nus-emo at semeval-2024 task 3: Instruction-tuning llm for multimodal emotion-cause analysis in conversations. *arXiv* preprint arXiv:2501.17261.
- Chandana Sree Mala, Gizem Gezici, and Fosca Giannotti. 2025. Hybrid retrieval for hallucination mitigation in large language models: A comparative analysis. *arXiv preprint arXiv:2504.05324*.
- OpenAI. 2023. Gpt-4 technical report. arXiv preprint arXiv:2303.08774.
- Jin-Soo Park, Xuesu Xiao, Garrett Warnell, Harel Yedidsion, and Peter Stone. 2023. Learning perceptual hallucination for multi-robot navigation in narrow hallways. In *Proceedings of the 2023 IEEE International Conference on Robotics and Automation (ICRA)*, London, England.
- Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. 2018. Object hallucination in image captioning. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4035–4045. Association for Computational Linguistics.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Edinburgh neural machine translation systems for wmt 16. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*, pages 371–376, Berlin, Germany. Association for Computational Linguistics.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*.
- Guanzheng Chen Xin Li Shijian Lu Chunyan Miao Lidong Bing Sicong Leng, Hang Zhang. 2023. Mitigating object hallucinations in large vision-language models through visual contrastive decoding. *arXiv* preprint arXiv:2311.16922.

Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan, Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, Kurt Keutzer, and Trevor Darrell. 2023. Aligning large multimodal models with factually augmented rlhf. *arXiv preprint arXiv:2309.14525*.

Zhongxiang Sun, Xiaoxue Zang, Kai Zheng, Yang Song, Jun Xu, Xiao Zhang, Weijie Yu, Yang Song, and Han Li. 2024. Redeep: Detecting hallucination in retrieval-augmented generation via mechanistic interpretability. *arXiv preprint arXiv:2410.11414*.

Andrés Villa, Juan León Alcázar, Motasem Alfarra, Vladimir Araujo, Alvaro Soto, and Bernard Ghanem. 2025. Eagle: Enhanced visual grounding minimizes hallucinations in instructional multimodal models. *arXiv preprint arXiv:2501.02699*.

Junyang Wang, Yuhang Wang, Guohai Xu, Jing Zhang, Yukai Gu, Haitao Jia, Jiaqi Wang, Haiyang Xu, Ming Yan, Ji Zhang, and Jitao Sang. 2023. An Ilm-free multi-dimensional benchmark for mllms hallucination evaluation. *arXiv preprint arXiv:2311.07397*.

Shukang Yin, Chaoyou Fu, Sirui Zhao, Tong Xu, Hao Wang, Dianbo Sui, Yunhang Shen, Ke Li, Xing Sun, and Enhong Chen. 2023. Woodpecker: Hallucination correction for multimodal large language models. arXiv preprint arXiv:2310.16045.

Zhiyuan Zhao, Bin Wang, Linke Ouyang, Xiaoyi Dong, Jiaqi Wang, and Conghui He. 2023. Beyond hallucinations: Enhancing lvlms through hallucinationaware direct preference optimization. *arXiv* preprint *arXiv*:2311.16839.

Yiyang Zhou, Chenhang Cui, Rafael Rafailov, Chelsea Finn, and Huaxiu Yao. 2024. Aligning modalities in vision large language models via preference finetuning. *arXiv* preprint arXiv:2402.11411.

A Additional Experimental Details

A.1 Implementation Details

Backbone and Setup. We apply ReLoop to five representative LVLMs with diverse architectures: MiniGPT-4, InstructBLIP, LLaVA-1.5, mPLUG-owl, and ShareGPT4V. To assess generalizability on stronger backbones, we further evaluate on LLaVA-1.6 and Qwen-VL-2.5 (Table 2). Importantly, we do not alter the internal structures of these models. ReLoop is introduced as a lightweight, external consistency-supervision framework during training. All backbones are initialized with their public checkpoints and keep their visual encoders (e.g., ViT, CLIP) frozen.

ReLoop Components. ReLoop introduces three frozen feedback modules: (1) CFP-Lang (MiniGPT-4-based reverse question reconstructor); (2) CFP-Vis (BLIP-2-based visual describer); (3) Attention Supervision that aligns decoder attention maps with entropy-based soft pseudo-labels. A frozen BERT encoder plus an MLP scorer serves as a lightweight semantic aggregator. All feedback modules remain frozen; only the backbone and the aggregator are updated.

Training Details. Experiments are performed on $8\times A100$ GPUs (80GB) using mixed-precision training (fp16) for 3 epochs. We adopt the AdamW optimizer with parameters $\beta_1=0.9$, $\beta_2=0.98$, and a weight decay of 0.05. The effective batch size is 128, with a gradient accumulation step of 8. The initial learning rate is set to 5×10^{-5} , along with 1,000 warm-up steps and cosine learning rate decay scheduling. Unless otherwise stated, all accuracy and resource measurements follow this main setup. For efficiency-only timing and method-level comparability, we additionally report a controlled regime: $4\times A100$ GPUs, batch size 12/GPU, and a fixed 2k-step schedule.

Loss Function. The overall objective is

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{sft}} + \alpha \, \mathcal{L}_{\text{align}} + \beta \, \mathcal{L}_{\text{vis}} + \gamma \, \mathcal{L}_{\text{attn}} + \lambda \, \Omega(\theta)$$
(6)

We set the hyper-parameters as $\alpha=1.0$, $\beta=0.7$, and $\lambda=10^{-5}$. The weight γ is dynamically adjusted by the Adaptive Consistency Weighting (ACW) mechanism, which modulates γ based on the BERTScore between the original and reconstructed questions (see Section 4.2.1).

A.2 Training Dataset Construction

We curated approximately 30K high-quality QAimage triplets from the LLaVA-Instruct-150K corpus (Liu et al., 2023a), each containing an image, an open-ended question, and a human-annotated answer. To simulate hallucination supervision, we generated semantically contradictory answers by modifying key elements (e.g., objects, attributes, or relations) in the references. These hallucinated samples were automatically constructed and manually verified for quality and type diversity. In Experiment 6.1, we selected 500 representative QA-image pairs from the filtered validation set based on POPE and MMHalBench, equally split between hallucinated and non-hallucinated cases. In Experiment 6.2, we evaluated five LVLMs on this curated set to assess the impact of ReLoop. Models with open alignment architectures (e.g., MiniGPT-4, InstructBLIP) showed the greatest improvement, while high-performing black-box models (e.g., ShareGPT4V) saw minimal gains, suggesting ReLoop's effectiveness hinges on alignment signal compatibility.

A.3 Evaluation Metrics

To comprehensively evaluate the effectiveness of ReLoop in mitigating hallucinations and enhancing visual grounding, we adopt a structured set of metrics covering both hallucination suppression and cross-modal consistency. In particular, shown in Table 4, we group the metrics into two key categories: *Hallucination Suppression*, which quantifies the presence of non-existent or spurious content, and *Cross-modal Faithfulness*, which assesses the semantic and perceptual alignment between generated text and visual input.

A.3.1 Metrics on Hallucination Suppression

For hallucination evaluation, we incorporate CHAIR (Rohrbach et al., 2018) to measure hallucination frequencies at instance levels and include POPE (Li et al., 2023b), a probing-based diagnostic benchmark to evaluate object hallucinations through direct VQA-style interactions. Together, these metrics allow us to holistically assess ReLoop's ability to suppress hallucinated content while preserving descriptive quality.

CHAIR (Rohrbach et al., 2018) (Caption Hallucination Assessment with Image Relevance) quantifies hallucinations by detecting whether the model-generated captions mention objects

that do not exist in the image. It provides two variants:

$$CHAIR_{I} = \frac{|\{\text{hallucinated objects}\}|}{|\{\text{all objects}\}|}$$
 (7)

$$CHAIR_S = \frac{|\{\text{hallucinated responses}\}|}{|\{\text{all responses}\}|} \quad (8)$$

where $CHAIR_I$ measures instance-level hallucination (object granularity) and $CHAIR_S$ measures sentence-level hallucination (response granularity).

- POPE (Li et al., 2023b) (Polling-based Object Probing Evaluation) automates hallucination detection via instance-level object probing. It:
 - Segments objects in the image;
 - Asks the model about object existence and introduces distractor queries;
 - Computes metrics such as F1 score to measure detection precision.

POPE offers direct insights into a model's visual grounding capability through objective visual questioning.

A.3.2 Metrics on Cross-modal Faithfulness

On the side of Cross-modal Faithfulness, we adopt Faith and Faith $_S$ (Jing et al., 2024), which evaluate how well the generated text is grounded in the visual input. Faith focuses on overall alignment, while Faith $_S$ specifically checks whether statements are supported by the visual evidence in a token-level or segment-wise manner. In addition, we report the F1 score, a standard metric that captures the harmonic mean of precision and recall between the predicted and reference entities. In our context, it reflects how well the model identifies relevant visual content without fabricating or omitting essential elements, thus serving as a practical indicator of the model's grounding precision and completeness.

• **F1 Score** reflects the harmonic mean of precision and recall in detecting whether queried objects exist. High F1 indicates accurate recognition and rejection of hallucinated entities:

$$F1 = 2 \cdot \frac{Precision \cdot Recall}{Precision + Recall}$$
 (9)

• **Faith** (Jing et al., 2024) measures the overall semantic alignment between image and response.

It uses automated matching or human verification to assess whether the content is factually grounded in the image:

$$Faith = \frac{|Aligned Statements|}{|Total Statements|}$$
 (10)

• **Faith**_S (Jing et al., 2024) extends Faith to a finer granularity by evaluating the support of specific sentence segments or tokens using cross-modal supervision or saliency alignment:

$$Faith_S = \frac{|Grounded Segments or Tokens|}{|Total Segments or Tokens|}$$
(11)

A.4 Evaluation Benchmark

Besides, to provide a fine-grained and multiperspective assessment of ReLoop's effectiveness in suppressing hallucinations and enhancing crossmodal faithfulness, we adopt four complementary benchmarks. AMBER (Wang et al., 2023) targets object-level hallucinations, while MMHal-B (Sun et al., 2023) and HallusionBench (Guan et al., 2024) assess errors in attributes, spatial relations, and perceptual consistency. MME (Fu et al., 2023) covers general multimodal capabilities such as OCR and counting. These benchmarks collectively evaluate generative and discriminative capabilities, entity grounding, perceptual consistency, and multimodal reasoning:

- AMBER (Wang et al., 2023): An LLM-free multi-dimensional benchmark that diagnoses hallucinations in both generative and discriminative tasks. It explicitly tests object *existence*, *attributes*, and *relations*, allowing us to assess ReLoop's object-level grounding fidelity, attribute correctness, and relational accuracy. This supports the evaluation of semantic precision in visual grounding.
- MMHal-B (Sun et al., 2023): A benchmark built upon fact-augmented reinforcement learning (RLHF) that penalizes hallucinated attributes and spatial configurations. MMHal-B offers targeted diagnostics for hallucination suppression in factual and compositional dimensions, particularly assessing whether ReLoop can resist overgeneralization and maintain factual grounding under complex prompts.
- HallusionBench (Guan et al., 2024): A benchmark that probes visual-linguistic robustness under ambiguous image-text settings. It emphasizes contextual grounding, requiring models to

handle subtle visual cues and nuanced linguistic traps. HallusionBench evaluates ReLoop's ability to maintain perceptual consistency and reject misleading contextual cues that typically trigger hallucinations.

 MME (Fu et al., 2023): A broad-spectrum benchmark measuring multimodal perception and cognition across 14 sub-tasks, including OCR, object counting, spatial reasoning, and commonsense grounding. MME validates whether ReLoop's structured supervision translates into generalized improvements in visual understanding and multimodal reasoning, beyond hallucination mitigation.

Together, these benchmarks offer layered supervision signals from fine-grained object hallucination detection to holistic multimodal cognition, providing strong empirical evidence of ReLoop's reliability across diverse real-world tasks.

A.5 Baseline Implementation

To evaluate ReLoop's generalizability and additive benefit, we compare it with three representative alignment-based hallucination mitigation strategies: LLaVA-RLHF (Sun et al., 2023), HA-DPO (Zhao et al., 2023), and POVID (Zhou et al., 2024). These baselines span a diverse range of supervision paradigms, from reinforcement learning to contrastive grounding. Importantly, all methods are applied on top of the same backbone (LLaVA-1.5) with consistent training configurations, ensuring fair comparison.

- LLaVA-RLHF (Sun et al., 2023) aligns responses to human preferences through reinforcement learning from human feedback. While effective for improving general fluency and tone, it does not explicitly penalize visual or factual inconsistencies.
- HA-DPO (Zhao et al., 2023) adopts hallucination-aware preference optimization by contrasting faithful versus hallucinated generations. This method introduces targeted loss signals during fine-tuning, encouraging the model to avoid semantically spurious content.
- **POVID** (Zhou et al., 2024) enhances visual grounding via perturbed image inputs, injecting contrastive visual signals to reduce reliance on textual priors and promote visual fidelity.

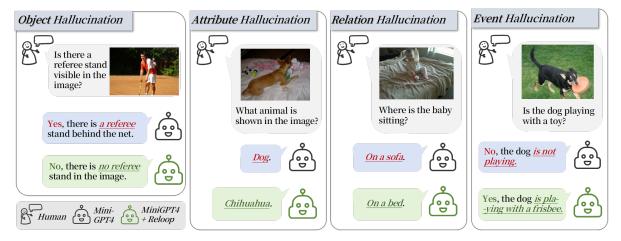


Figure 6: **Case Study:** Comparison between MiniGPT-4 and ReLoop across four types of hallucination in openended VQA: *Object*, *Attribute*, *Relation*, and *Event*. ReLoop produces more accurate and grounded responses by aligning its outputs with both the visual evidence and the question semantics.

• VCD (Sicong Leng, 2023) is a training-free, inference-time method that contrasts output distributions from the original image and a visually distorted counterpart, down-weighting tokens supported mainly by unimodal/language priors and thereby mitigating object hallucinations without further tuning.

Results from both fine-grained hallucination metrics (Table 4) and benchmark-level evaluations (Table 6) demonstrate that ReLoop consistently outperforms all competing methods. These results validate ReLoop as a robust and generalizable framework capable of enhancing multimodal model performance beyond what is achievable by current alignment-based techniques alone.

A.6 Training Cost Breakdown

We quantify ReLoop's computational footprint and efficiency trade-offs from four angles: (i) crossmethod training cost, (ii) overhead attribution via CFP-component ablations, (iii) convergence versus training steps, and (iv) per-epoch cost under a controlled regime. Unless otherwise stated, CFP modules are frozen (no gradient), and only the main model M and the Semantic Aggregator S are updated.

A.6.1 Cross-Method Training Cost

We compare ReLoop against representative alignment methods. As summarized in Table 7, ReLoop updates only M and S with frozen CFPs, substantially reducing optimizer/memory footprint compared to end-to-end baselines. Despite extra forward passes for CFPs, ReLoop avoids optimizer

states for large frozen modules, yielding substantially lower GPU hours and memory than end-to-end RLHF/DPO. This design preserves training affordability while enabling closed-loop supervision.

A.6.2 Overhead Attribution via CFP Components

We attribute runtime and memory overheads to individual CFP branches on LLaVA-1.5, keeping data/backbone constant. Table 8 reports GPU hours, memory, latency, and POPE. As seen in Table 8, both language and visual CFPs contribute positively to hallucination suppression, with additive gains. The per-branch overhead is bounded and predictable. Since CFPs are optional at inference, the training-time penalty does not translate to deployment latency. Both language and visual CFPs contribute positively to hallucination suppression, with additive gains. The overhead is bounded and predictable per component. Since CFPs are optional at inference, the training-time penalty does not translate to deployment latency.

A.6.3 Convergence and Epoch-Level Cost

We study returns versus training steps and report per-epoch cost under the controlled regime (4×A100, batch 12/GPU, fixed 2k steps) for method-level comparability. Step-wise gains are summarized in Table 10; the end-to-end throughput and slowdown are reported in Table 9. Table 10 shows fast early gains with diminishing returns beyond 2k steps. Under a fixed budget, Table 9 in-

¹Scores here follow the controlled 2k-step regime and are not directly comparable to full-training results elsewhere.

| Method | Trainable Params | Feedback Modules | GPU Hours | GPU Type | Peak VRAM | Notes |
|-----------------------|----------------------------|---------------------|--------------|-------------|--------------|-----------------------------|
| ReLoop (MiniGPT-4) | $M + S (\sim 1.2B)$ | Frozen | 3 | A100-40G | ∼26GB | Fastest; ablation |
| ReLoop (LLaVA-1.5) | $M+S~(\sim 13 \mathrm{B})$ | Frozen | 6 | A100-80G | ∼48GB | Core experiment. (Sec. 5.2) |
| RLHF finetuning | Whole model $(\sim 13B)$ | End-to-end | 20 | A100-80G | ~70GB | _ |
| DPO | Whole model | End-to-end | 8-16 | A100-80G | >50GB | _ |
| Contrastive alignment | Whole model | Encoder fusion | ~10 | A100-80G | ~32GB | - |

Table 7: Cross-method training cost. ReLoop updates only M and S with frozen CFPs, reducing optimizer/memory footprint relative to end-to-end baselines.

| Variant | CFP Modules | GPU Hours | Peak Mem | Train Latency | POPE ↑ |
|---------------|-------------|-----------|----------|---------------|--------|
| ReLoop (Full) | Lang + Vis | 6.0 | 48GB | 1.8× | 82.2 |
| w/o Lang CFP | Vis only | 5.2 | 40GB | 1.4× | 81.5 |
| w/o Vis CFP | Lang only | 4.7 | 38GB | 1.3× | 81.1 |
| No CFP | None | 3.0 | 26GB | $1.0 \times$ | 77.2 |

Table 8: CFP-component ablations (LLaVA-1.5). Each CFP adds \sim 1 GPU hour and \sim 10–12GB VRAM, yielding +3–5 POPE. CFPs are frozen yet incur forward attention/encoding; batch size may reduce (e.g., $16 \rightarrow 12$ on 100).

dicates that a moderate training slow-down trades for sizable hallucination reduction, while inference cost remains unchanged by dropping CFPs at test time.

A.7 Contrastive Augmentation Ablation

To evaluate whether RELOOP depends on the manually perturbed semantic negatives (Appendix A.2), we ablate this contrastive augmentation and train on the standard LLaVA-Instruct data only, keeping all other settings identical to the main experiments. Even without contrastive augmentation, RELOOP significantly improves over the LLaVA-1.5 baseline on all metrics (Table 11), demonstrating that the core closed-loop alignment mechanism is effective when trained purely on standard supervision data. The contrastive examples provide further refinement, most notably on AMBER and MMHal-B, but they are not essential for RELOOP to outperform existing alignment strategies.

B Case Study

We present a qualitative case study to analyze how ReLoop mitigates hallucination across four representative types:

• **Object Hallucination**: The baseline model incorrectly asserts the presence of a "referee stand" which is not in the image. ReLoop corrects this by recognizing the absence of such an entity.

- Attribute Hallucination: An animal is mislabeled as "dog" instead of "chihuahua." ReLoop identifies the finer-grained attribute correctly.
- **Relation Hallucination**: The spatial relationship "on a sofa" is incorrectly predicted; ReLoop grounds the child's location more accurately.
- Event Hallucination: The action "not playing" contradicts visual evidence; ReLoop revises the answer to match the depicted motion.

As shown in Figure 6, baseline models such as MiniGPT-4 frequently produce fluent yet inaccurate answers that are not grounded in the image. ReLoop corrects these errors by leveraging consistency feedback to align its answers with both the question intent and visual content. The examples highlight ReLoop's capacity to suppress diverse hallucination patterns and improve factual reliability in open-ended VQA.

B.1 Handling Nonsensical or Unrelated Initial Answers

This section complements the above qualitative cases by focusing on an orthogonal failure mode raised by reviewers: how ReLoop handles instances where the initial answer A is nonsensical or unrelated to the question. ReLoop employs a set of structured safeguards: early rejection, γ -based downweighting via ACW, and entropy-aware masking of attention. The set of structured safeguards

| Model | Time / Epoch | Total GPU Hrs (2k) | Peak Memory | Throughput (img/s/GPU) | Slowdown vs. Base | POPE ↑ | CFP @ Inference |
|-----------------------|-----------------|-----------------------|----------------|------------------------|----------------------|---------------|--------------------|
| LLaVA-1.5 (base) | ∼58 min | 3.0 | 26GB | 11.6 | 1.0× | 77.2 | N/A |
| LLaVA-1.5 + ReLoop | ~103 min | 6.0 | 48GB | 6.4 | 1.77× | 82.2 | Optional (off) |

Table 9: Per-epoch cost under the controlled regime. ReLoop increases epoch time by \sim 77% yet converges within comparable steps. CFPs are disabled at inference, so deployment latency remains unchanged.

| Steps | GPU hours | POPE ↑ | ΔΡΟΡΕ | Gain / hour |
|--------------------|--------------|---------------|-------|----------------|
| 0 (no Training) | 0.0 | 77.2 | - | _ |
| 1k | 3.6 | 80.1 | +2.9 | 0.81 |
| 2k | 7.2 | 81.6 | +1.5 | 0.21 |
| 3k | 10.9 | 82.2 | +0.6 | 0.08 |

Table 10: Diminishing returns with steps (LLaVA-1.5). >90% of gains appear within 2k steps (\sim 7 GPU hours).

| Method | AMBER ↑ | MME↑ | MMHal-B↑ | Hallu-B↑ |
|-----------------------------------|----------------|------|----------|----------|
| LLaVA-1.5 (baseline) ReLoop | 73.9 | 1513 | 65.4 | 48.6 |
| (w/o aug- mentation) | 78.5 | 1452 | 67.1 | 50.1 |
| Full ReLoop | 80.3 | 1505 | 68.9 | 52.3 |

Table 11: Ablation on contrastive augmentation (LLaVA-1.5 backbone). Even without augmentation, RELOOP improves over the base model across all benchmarks; augmentation yields additional gains.

can prevent misleading updates when the feedback is deemed unreliable. The representative examples in Table 12 illustrate how these filters operate at the case level before gradients are applied. In addition, the quantitative stress tests in the main text (see Table 5). demonstrate stable performance under injected noise, corroborating the robustness of these safeguards.

C Entropy-based Pseudo Ground-Truth Attention \mathcal{H}_{pseudo}

This section details the construction of the entropybased pseudo ground-truth attention used in Sec. 4.2.3. We describe tensor shapes, multilayer/head aggregation, token filtering, smoothing/normalization, and edge cases to facilitate reproduction.

Notation and Shapes. Let the decoder cross-attention at decoding step t be $\{A_t^{(\ell,h)} \in \mathbb{R}^S\}_{\ell=1..L,\,h=1..H}$ over S visual patches (keys), for

L layers and H heads. We first aggregate heads and layers to obtain a single distribution over patches for token t:

$$\bar{a}_t = \sum_{\ell=1}^L \sum_{h=1}^H w^{(\ell)} u^{(h)} A_t^{(\ell,h)}$$
 (12)

$$\sum_{\ell} w^{(\ell)} = \sum_{h} u^{(h)} = 1 \tag{13}$$

where $w^{(\ell)}$ and $u^{(h)}$ are fixed convex weights. In a default setting, we set uniform across heads $(u^{(h)}{=}1/H)$ and a back-loaded layer prior $(w^{(\ell)} \propto \exp(\kappa\,\ell/L))$ with $\kappa{=}1.5)$ to emphasize later layers. We row-normalize to obtain a probability over spatial patches:

$$p_{t,s} = \frac{\bar{a}_t[s]}{\sum_{s'} \bar{a}_t[s']} \in [0,1], \qquad \sum_{s} p_{t,s} = 1.$$
(14)

Per-token Entropy and Confident Set. For each generated token t, we excludes BOS/EOS/padding/special tokens and computes entropy:

$$\mathcal{E}_{t} = -\sum_{s=1}^{S} p_{t,s} \log p_{t,s} \in [0, \log S] \quad (15)$$

and form a confident token set:

$$\mathcal{T}_{conf} = \{ t \mid \mathcal{E}_t \le \tau \}$$
 (16)

In a default setting, we set τ =2.0 (nats). In low-entropy regimes ($\mathcal{E}_t \approx 0$) the attention is highly focused; high entropy indicates diffuse/unstable focus. In an optional variant (ablation only) setting, we reweight votes by $w_t = \max \left(0, 1 - \frac{\mathcal{E}_t}{\log S}\right)$; we do not apply this by default to keep the scheme simple and robust.

Voting and Temperature-normalized Map. Confident tokens vote for patch importance by summation:

$$\tilde{h}[s] = \sum_{t \in \mathcal{T}_{conf}} p_{t,s} \tag{17}$$

| Failure Type | Q | A | Action Taken |
|---------------------------|-------------------------------|--------------------------------|--|
| Empty Output | "What is on the table?" | (empty) | Early rejection: sample skipped; no feedback loss computed. |
| Overly Generic | "What sport is being played?" | "I'm not sure." | γ -downweighting: low-confidence ACW weight; skipped if below length threshold. |
| Nonsensical Repetition | "What is the man holding?" | help!" | Syntactic abnormality detection: reject via regex/#token heuristics; no update. |
| Unrelated Semantic | "What color is the bus?" | "Apples grow in the summer." | Semantic mismatch: $\gamma \approx 0$ from BERTScore \Rightarrow loss suppressed. |
| Hallucinated Words | "What are the people doing?" | "Grockling spinners do fleeb!" | Entropy + token validation: flat attention masked; spurious tokens filtered before loss. |

Table 12: Case-level handling of nonsensical or unrelated initial answers in ReLoop. Structured safeguards (early rejection, ACW γ -based downweighting, and entropy-aware masking) prevent misleading gradients from invalid feedback.

We convert \tilde{h} to a soft target by temperature soft-max (Default: T_a =0.7):

$$\mathcal{H}_{\text{pseudo}}[s] = \frac{\exp(\tilde{h}[s]/T_a)}{\sum_{s'} \exp(\tilde{h}[s']/T_a)}$$
(18)

Spatial Reshaping and Smoothing. Let the S patches correspond to a $(H_p \times W_p)$ grid from the vision encoder (e.g., ViT patch tokens). We reshape $\mathcal{H}_{pseudo} \in \mathbb{R}^S$ to $\mathbb{R}^{H_p \times W_p}$, apply light Gaussian smoothing (3×3 kernel, σ =0.8; reflect padding), then flatten back to length S. A final ℓ_1 normalization ensures $\sum_s \mathcal{H}_{pseudo}[s]$ =1.

Special Tokens, Padding, and Masking. We exclude special tokens (BOS/EOS, padding) and punctuation-only tokens from \mathcal{T}_{conf} . For subword tokenization, all subpieces are treated uniformly; no external POS/saliency tools are used to preserve the unsupervised nature.

Empty-set and Degenerate Cases. If $\mathcal{T}_{conf} = \emptyset$ (rare; e.g., extremely diffuse attention), we fall back to a min-entropy top-k strategy: pick the $k = \max(1, \lceil 0.01T \rceil)$ lowest-entropy tokens to form \mathcal{T}_{conf} and proceed with Eq. (17). If \tilde{h} is flat (numerically), we use a near-uniform prior slightly peaked at the global min-entropy token's argmax.

Objective and Gradient Flow. The attention supervision minimizes a KL divergence between the model's cross-attention map \mathcal{H} (from the current forward pass) and \mathcal{H}_{pseudo} :

$$\mathcal{L}_{\text{attn}} = \text{KL}(\mathcal{H} \parallel \mathcal{H}_{\text{pseudo}})$$

$$= \sum_{s} \mathcal{H}[s] \Big(\log \mathcal{H}[s] - \log \mathcal{H}_{\text{pseudo}}[s] \Big)$$
(19)

Algorithm 1 Constructing Entropy-based Pseudo Attention \mathcal{H}_{pseudo}

```
Require: Cross-attn \{A_t^{(\ell,h)} \in \mathbb{R}^S\}, weights \{w^{(\ell)}\},
         \{u^{(h)}\}, threshold 	au, temperature T_a
Ensure: \mathcal{H}_{pseudo} \in \mathbb{R}^{S} (stop-gradient)
   1: \mathcal{T}_{\mathrm{conf}} \leftarrow \emptyset
  2: for each token t in generated tokens (exclude spe-
                 \begin{array}{l} \bar{a} \leftarrow \sum_{\ell=1}^{L} \sum_{h=1}^{H} w^{(\ell)} u^{(h)} A_t^{(\ell,h)} \\ p \leftarrow \bar{a} / \sum_{s} \bar{a}[s] \\ \mathcal{E}_t \leftarrow - \sum_{s} p[s] \log p[s] \\ \text{if } \mathcal{E}_t \leq \tau \text{ then} \end{array} 
                                                                                                   ⊳ Eq. 13⊳ Eq. 14⊳ Eq. 15
  3:
                          \mathcal{T}_{\text{conf}} \leftarrow \mathcal{T}_{\text{conf}} \cup \{(t,p)\}
  7:
  8:
  9: end for
10: if \mathcal{T}_{conf} = \emptyset then
                  \mathcal{T}_{\text{conf}} \leftarrow \text{top-}k \text{ lowest-entropy tokens}
11:
13: h \leftarrow \sum_{(t,p) \in \mathcal{T}_{conf}} p
                                                                                                      ⊳ Eq. 17
14: \mathcal{H}_{pseudo} \leftarrow \operatorname{softmax}(\tilde{h}/T_a)
                                                                                                      ⊳ Eq. 18
15: \mathcal{H}_{pseudo} \leftarrow smooth(reshape\_grid(\mathcal{H}_{pseudo}))
 16: \mathcal{H}_{pseudo} \leftarrow \mathcal{H}_{pseudo} / \|\mathcal{H}_{pseudo}\|_1
                                                                                                  \triangleright \ell_1 norm
17: return stopgrad(\mathcal{H}_{pseudo})
```

We stop gradients through \mathcal{H}_{pseudo} ; only the main model's attention is updated. The loss is weighted by γ in the total objective (Sec. 4.2.1).

Pseudocode The procedure (Algorithm 1) is unsupervised and self-adaptive, relying solely on model-internal attentional confidence and requiring no external saliency or human annotation. Its design (entropy gating, temperature control, and mild spatial smoothing) yields stable targets for the KL alignment used in Sec. 4.2.3.