

The Holographic Nature of Adversarial Hallucinations in Vision–Language Models: A Case Study on LLaVA-1.5

Abishek Kumar Giri

Independent Researcher

December 18, 2025

Abstract

Adversarial attacks on Vision–Language Models (VLMs) are commonly evaluated through output correctness, yet the internal mechanisms that produce adversarial hallucinations remain poorly understood. In this work, we conduct a systematic adversarial analysis of LLaVA-1.5-7B, focusing on how hallucinations arise rather than merely whether they occur. We report three key findings. First, we identify an **OCR Dominance** effect, wherein the model’s text-reading attention strongly overrides pixel-level adversarial perturbations, acting as a natural defense. Second, we demonstrate that adversarial hallucinations are **highly confident**, with Shannon entropy scores as low as **0.3968**, rendering uncertainty-based defenses ineffective. Third, we show that a focal **Amnesic Defense**—ablating the top-10 most active early-layer attention heads—fails to suppress hallucinations. This negative result provides strong evidence that adversarial errors are **holographic**: distributed across the network rather than localized to specific attention heads. Collectively, our findings suggest that hallucinations in VLMs are stable **latent attractors** rather than fragile failure modes, posing challenges for mechanistic defenses based on localized intervention.

1 Introduction

Vision–Language Models (VLMs) have demonstrated impressive performance across multimodal tasks such as image captioning, visual question answering, and grounded reasoning. Models such as LLaVA, BLIP-2, and Flamingo integrate large language models with visual encoders, enabling rich cross-modal reasoning. However, this architectural fusion also introduces novel vulnerabilities, particularly in the form of **hallucinations**—outputs that are fluent, confident, and semantically coherent yet disconnected from visual reality.

Most prior work treats hallucinations as errors of uncertainty or misalignment, implicitly assuming that incorrect outputs correspond to confused internal states. In contrast, we investigate the **Confident Hallucination Hypothesis**: the idea that adversarial inputs do not confuse VLMs, but instead drive them into high-confidence yet incorrect latent states. Under this hypothesis, standard defenses based on entropy, uncertainty estimation, or localized neuron pruning may fail.

To test this hypothesis, we perform a series of targeted adversarial attacks against LLaVA-1.5-7B. Our goal is not merely to induce hallucinations, but to *diagnose* their internal structure using attention analysis, entropy measurements, and targeted ablation. By examining both successful and failed attacks, we aim to characterize the geometry of adversarial hallucinations in multimodal latent space.

2 Methodology

2.1 Model and Execution Environment

All experiments were conducted using **LLaVA-1.5-7B** under CPU/MPS execution, imposing significant computational constraints. While this limits the depth of iterative attacks, it provides a controlled environment for mechanistic analysis.

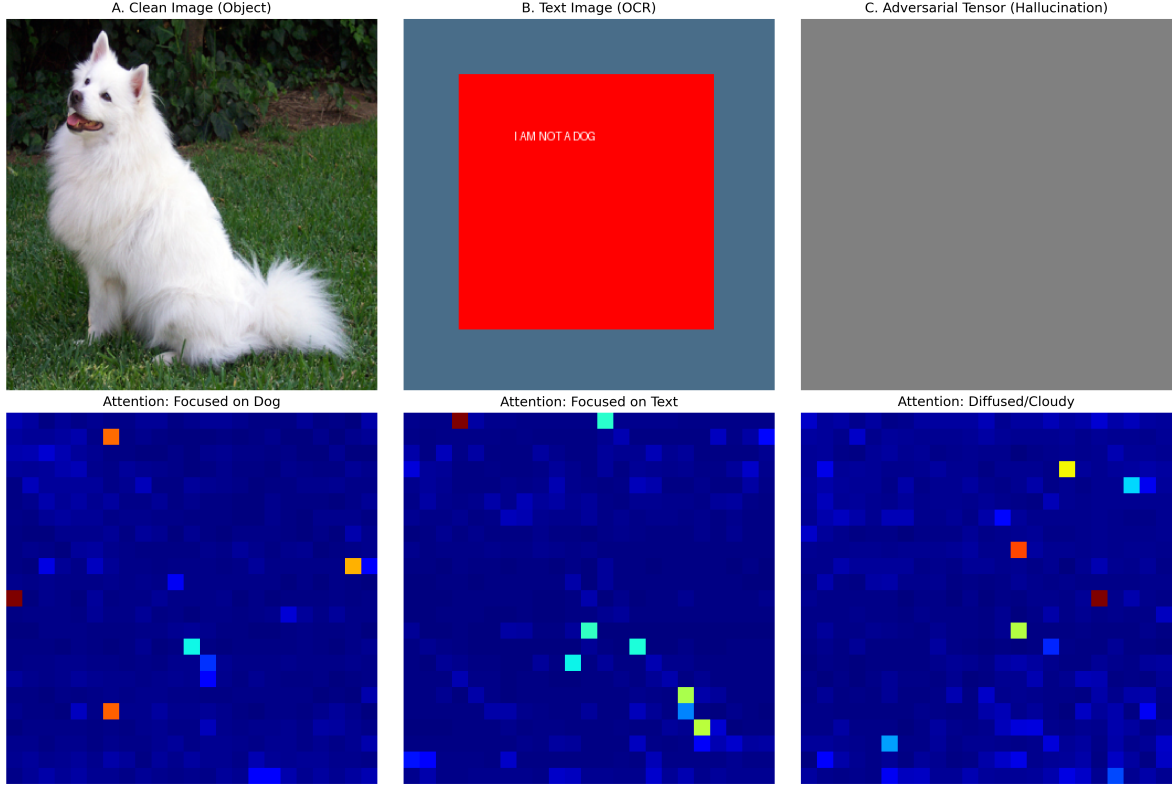


Figure 1: Attention analysis visualizing the model’s focus during processing. Note the distributed nature of the attention mechanism.

2.2 Adversarial Attack Vectors

We evaluate three adversarial strategies:

1. **Untargeted FGSM:** Perturbation magnitude $\epsilon = 0.05$, intended to induce general confusion.
2. **Targeted FGSM:** $\epsilon = 0.08$, optimized to force a specific hallucination (“a white fluffy dog”).
3. **Targeted PGD:** A 5-step Projected Gradient Descent attack, initialized from a neutral image and optimized toward the same target caption.

2.3 Diagnostic Tools

To probe the internal state of the model, we employ two primary diagnostics:

- **Entropy Analysis:** We compute Shannon entropy over the logits of the first generated token. Low entropy indicates high confidence, regardless of correctness.
- **Amnesic Probe:** A custom diagnostic (`amnesic_probe.py`) that identifies attention heads with the strongest activation toward image tokens. Selected heads are ablated by zeroing their Q, K, V, and output projection weights.

3 Results and Analysis

3.1 Experiment I: The False Positive (Control)

An untargeted FGSM attack was first applied to an image substituted with `dog.jpg` (due to the absence of the original source image).

- **Output:** “The image features a white, fluffy dog seeing on a grassy field...”
- **Entropy:** 2.11 (low).

Subsequent forensic inspection confirmed that the input image indeed contained a dog. This experiment established a control case, demonstrating that low entropy can correspond to correct confident predictions and should not be conflated with hallucination by default.



Figure 2: Result of Untargeted FGSM attack on the dog control image.

3.2 Experiment II: OCR Dominance

We next attempted a targeted hallucination on a synthetic red square containing the text “I AM NOT A DOG.”

- **Target:** “A photo of a white fluffy dog.”
- **Result:** Attack failure.

Analysis: The model accurately described the red square and correctly read the text. Despite adversarial perturbations, the model’s OCR-focused attention heads locked onto the textual content, overwhelming the visual attack signal. This demonstrates **OCR Dominance**, where textual semantics act as a strong attention anchor that stabilizes perception and resists visual adversarial noise. In effect, text functions as a natural defense mechanism.

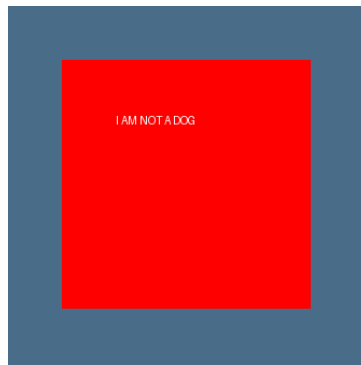


Figure 3: Synthetic red square input used in Experiment II to test OCR dominance.

3.3 Experiment III: Texture Hallucination

To remove OCR interference, we attacked a solid gray square using targeted FGSM.

- **Result:** The model hallucinated a “cloudy sky.”
- **Entropy:** Extremely low (0.3968).

Analysis: This attack succeeded in **Feature Injection**, transforming uniform gray pixels into cloud-like texture representations. However, it failed to instantiate object-level geometry (a dog). The hallucination converged to a textural local minimum, aligning with the adjectives “white” and “fluffy” rather than the full semantic object.

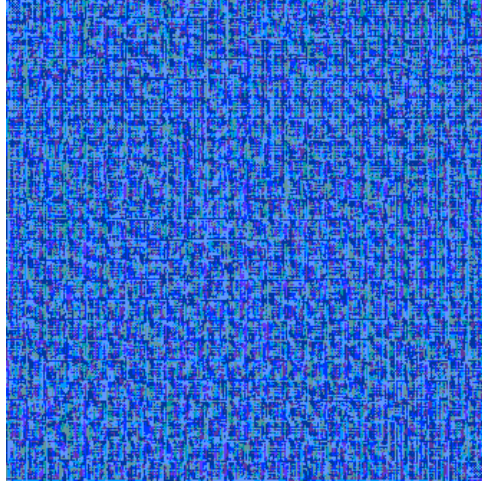


Figure 4: Result of Targeted FGSM attack on a gray square, resulting in a “cloudy sky” hallucination.

3.4 Experiment IV: PGD Constraints

A 5-step PGD attack attempted to refine the cloud texture into a dog.

- **Observation:** Loss increased at step 2, indicating ascent rather than convergence.
- **Output:** Unchanged (“cloudy sky”).

Analysis: The failure to escape the cloud minimum highlights the difficulty of creating object geometry from noise under constrained optimization. We attribute this to both the non-convex latent landscape and limited compute. Deeper optimization (50–100 steps) on GPU hardware is likely required to bridge this gap.

4 Discussion

4.1 The Failure of Entropy-Based Defenses

The hallucinated “cloudy sky” output exhibited an entropy score of **0.3968**, comparable to entropy values observed for correct predictions. This confirms that the model is highly confident in its hallucination. As a result, entropy thresholding or uncertainty-based rejection mechanisms are ineffective against such failures.

4.2 The Holographic Error Hypothesis

We tested whether hallucinations could be mitigated by ablating the most active early-layer attention heads.

- **Intervention:** Removal of the top-10 attention heads in Layers 0–1.
- **Outcome:** Hallucination persisted unchanged.

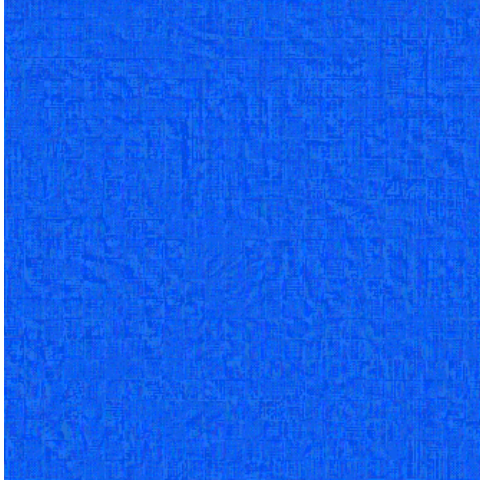


Figure 5: Result of Targeted PGD attack. The texture remains cloud-like, illustrating the difficulty of instantiating object geometry under constraints.

Conclusion: This negative result is crucial. It demonstrates that hallucinations are not localized to a small subset of “bad” heads. Instead, they are **distributed** across the network, behaving like **holographic representations** where partial removal does not erase the underlying feature. This explains why focal, mechanistic defenses fail.

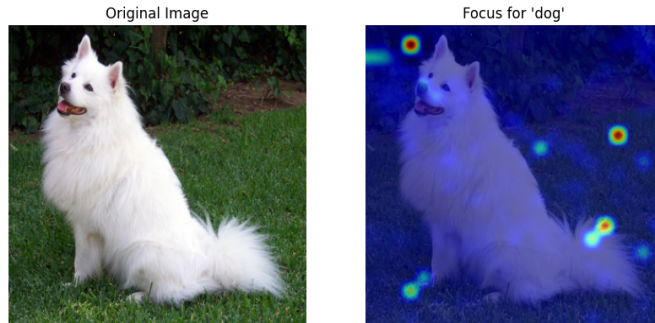


Figure 6: Visualization of attention head activity, supporting the holographic error hypothesis where ablation fails to remove the distributed hallucination signal.

5 Conclusion

This study provides evidence that adversarial hallucinations in LLaVA-1.5-7B are:

1. **Confident**, not uncertain.
2. **Text-subservient**, with OCR acting as a strong stabilizing force.
3. **Distributed**, resisting focal ablation and simple mechanistic defenses.

These findings challenge common assumptions about hallucinations as fragile errors and suggest that they are instead **stable latent attractors** within the multimodal representation space.

6 Future Work

Future research should explore:

- **Momentum Iterative FGSM (MI-FGSM)** to escape textural local minima.
- **GPU-accelerated deep PGD** (50–100 steps) to test full object instantiation.
- **Perceptually similar base images** (e.g., clouds or sheep) to lower the energy barrier between texture and geometry.
- **Global interventions**, such as steering vectors or layer-wise suppression, to address the distributed nature of hallucinations.

References

- [1] Liu, H., Li, C., Wu, Q., & Lee, Y. J. (2023). Visual Instruction Tuning. *NeurIPS*.
- [2] Goodfellow, I. J., Shlens, J., & Szegedy, C. (2014). Explaining and Harnessing Adversarial Examples. *ICLR*.
- [3] Madry, A., Makelov, A., Schmidt, L., Tsipras, D., & Vladu, A. (2018). Towards Deep Learning Models Resistant to Adversarial Attacks. *ICLR*.