

# MagicBench: Diagnosing Visual Agency Loss and Semantic Dependency in Multimodal LLMs

Tang Da Huang, Weidong Tang, Wen Qi Xu, Xianpeng Guo\*

Xidian University

Xi'an, China

24209100251@stu.xidian.edu.cn, wdtang0705@gmail.com,

24209100083@stu.xidian.edu.cn, guoxianpeng@xidian.edu.cn

## Abstract

Multimodal Large Language Models typically assume linguistic context invariably enhances visual understanding. We study this assumption in semantic adversarial scenarios, specifically magic tricks, where narration deliberately diverges from physical reality. We introduce MagicBench, a diagnostic benchmark of 402 videos for evaluating MLLMs under hierarchical linguistic interference, together with a Physical Constraint Set (PCS) protocol for assessing adherence to physical laws. Evaluation uncovers a Semantic Dependency Paradox: (1) **Semantic anchoring**: Entity nouns act as anchors aiding localization, paradoxically boosting performance despite false predicates. (2) **Visual Agency Loss**: In semantic vacuums, multimodal performance collapses 12.4% ( $p < 0.01$ ) below the vision-only *capability probe*. This gap persists under symmetric prompting, suggesting a form of functional perception suppression in which autonomous visual search may be under-utilized in multimodal settings without linguistic triggers. Causal interventions via spatial prompting and signal magnification provide evidence that internal reasoning remains functional, supporting the interpretation of a perceptual access bottleneck. Our findings suggest MLLMs function as *language-guided passive observers*, advocating for perceptually-independent architectures that decouple sensory agency from linguistic dominance. Code and dataset are available at <https://github.com/Ink-Dawn/MagicBench>

## 1 Introduction

The evolution of Multimodal Large Language Models (MLLMs) has ushered in early-fusion architectures, where visual and textual tokens are processed jointly in a shared latent space (Alayrac et al., 2022; Gemini Team et al., 2024). Driven

by instruction tuning (Liu et al., 2023), the community largely operates under a *context-is-king* paradigm, assuming that richer linguistic context *monotonically* enhances visual grounding (Wang et al., 2023). This creates an inductive bias toward linguistic over-reliance, presuming a cooperative relationship between modalities while overlooking the risk of asymmetric modality dominance, where linguistic priors may prematurely gate or suppress visual evidence.

However, this assumption fractures in adversarial scenarios where language is deceptive or absent. While prior research quantified **hallucination** in static contexts (Yin et al., 2024; Li et al., 2023), dynamic events require models to maintain **Autonomous Physical Entailment (APE)**, defined as the capacity to prioritize sensory evidence and track causal state changes (e.g., object permanence) against a contradictory narrative. Current benchmarks (Yu et al., 2024; Maaz et al., 2024) often lack this causal complexity, failing to disentangle whether a model is truly perceiving the physics or merely defaulting to linguistic shortcuts (Leng et al., 2024; Golovanevsky et al., 2025).

To probe these boundaries, we propose Magic as a uniquely rigorous testbed for APE. A magician’s *patter* is designed to decouple symbolic narratives from physical reality, exploiting cognitive blind spots to mislead the observer. We use magic as a controlled testbed for causal state tracking and cross-modal conflict resolution under adversarial narration. We introduce MagicBench, a diagnostic benchmark of 402 high-fidelity videos organized by linguistic interference type. Adjudicated via a Physical Constraint Set (PCS) grounded in First-Order Logic (FOL), our evaluation uncovers a counter-intuitive Semantic Dependency Paradox.

This paradox manifests in two distinct effects (illustrated in Figure 1). First, semantic anchoring acts as a *spotlight* in deceptive scenarios, where entity nouns (e.g., “coin”) facilitate localization

\*Corresponding author.

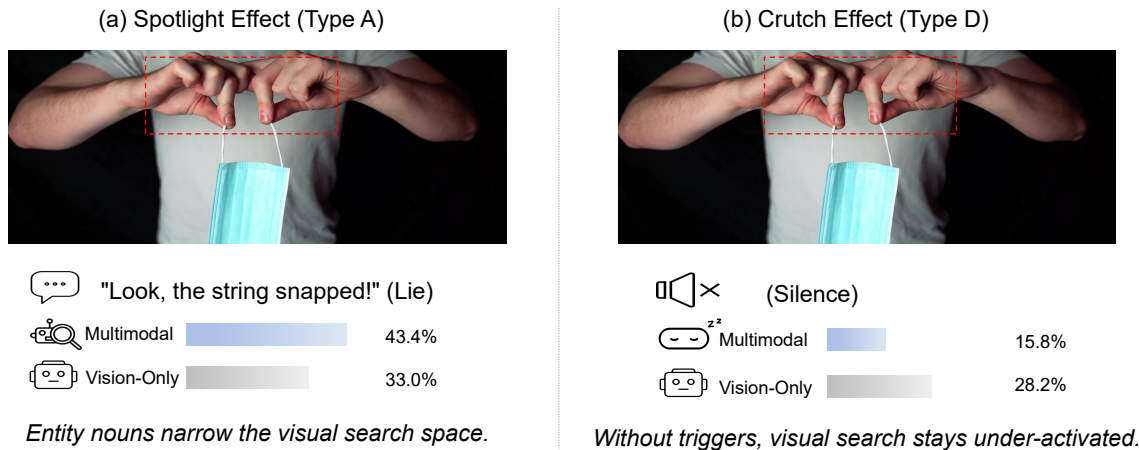


Figure 1: Illustration of the Semantic Dependency Paradox. **Left:** In deceptive scenarios, entity nouns such as “string” act as semantic anchors. Although the predicate is false (“snapped”), the noun still provides a spatial prior that helps object localization. **Right:** In semantic vacuums (Silence), multimodal performance drops sharply. Without a linguistic trigger, the model enters a form of perceptual idling, whereas the *Vision-Only* capability probe can still attend to the finger-loop maneuver.

by providing implicit spatial priors; consequently, the gain from object grounding consistently outweighs the interference introduced by false predicates (lies). Second, and more alarmingly, we identify a **Visual Agency Loss** (acting as a linguistic *crutch*) in semantic vacuums (Silence). Here, multimodal performance collapses **12.4%** ( $p < 0.01$ ) below the *Vision-Only* capability probe, revealing a substantial under-utilization of available visual evidence. This gap persists under symmetric intensity ablation (SIA), where prompts are identical, suggesting that the tested MLLMs often behave like “language-guided passive observers” that enter a state of perceptual idling without linguistic triggers, rather than actively initiating autonomous visual search processes.

Beyond observation, our interventional diagnosis using exogenous spatial prompting (e.g., red circular markers) and signal enhancement tests provides evidence that the observed failures are more consistent with a perceptual access bottleneck than with a pure reasoning deficit. Our macro-behavioral findings strongly align with recent mechanistic evidence suggesting functional suppression of visual tokens within transformer blocks (Liu et al., 2025; Nikankin et al., 2025). Our contributions are fourfold: (1) **MagicBench**, a 402-video diagnostic benchmark for semantic adversarial robustness; (2) an analysis of the spotlight and crutch effects in multimodal perception; (3) causal evidence consistent with an access bottleneck in the tested early-fusion MLLMs; and (4) empirical mo-

tivation for architectures with stronger perceptual independence.

## 2 Related Work

### 2.1 Physical Entailment and Prompt-Induced Hallucination

Hallucination research has transitioned from static object existence (Li et al., 2023) to temporal consistency in video domains (Maaz et al., 2024; Li et al., 2024). However, most benchmarks focus on summarization-level events. MagicBench introduces a more granular challenge: sub-second state tracking under adversarial narratives. We build on the framework of prompt-induced hallucination (Rudman et al., 2026), where deceptive linguistic context can divert the model from physical evidence. Unlike synthetic benchmarks like CLEVRER (Yi et al., 2020), we test **Physical Entailment** in real-world scenarios, requiring models to maintain causal consistency against narratives that actively contradict visual kinematics.

### 2.2 Modality Conflict and Competitive Circuits

The competition between pixels and priors is a focal point of recent mechanistic studies. While visual override, that is, seeing overriding knowing, is possible in simple tasks (Ortu et al., 2025; Golovanevsky et al., 2025), our results suggest that it is substantially more fragile in dynamic adversarial contexts. This aligns with recent circuit-level evidence: Nikankin et al. (2025) and Hua et al. (2025)

reveal that MLLMs employ distinct, competing circuits for different modalities, often leading to the **active suppression** of visual signals when linguistic priors are strong. Our observation of *Visual Agency Loss* is consistent with the “Bad Eyes” phenomenon (Aravindan et al., 2025), in which critical visual information may be lost or misrouted during cross-modal fusion (Liu et al., 2025).

### 2.3 Active Perception and Visual Steering

Standard benchmarks often allow models to succeed via passive recognition. In contrast, the theory of Active Perception (Bajcsy et al., 2016) posits that vision is an exploratory activity. MagicBench operationalizes this by identifying a form of “perceptual idling” in current architectures. To mitigate such deficits, recent works explore Visual Task Prompting (VTPrompt) (Xiao et al., 2025) and Image-of-Thought (IoT) pipelines (An et al., 2025) to steer attention via spatial markers or sequential auditing. Systematic frameworks like VP-Bench (Xu et al., 2025) evaluate these steering capabilities primarily in cooperative settings. MagicBench instead studies whether such steering becomes necessary because models fail to autonomously initiate visual search under semantic adversarial conditions. In this setting, attention-steering (Shi et al., 2023) can improve performance, but the underlying modality dependency remains only partially mitigated.

## 3 MagicBench Construction

### 3.1 Data Curation and APE Capabilities

We curate MagicBench, a diagnostic benchmark comprising 402 high-fidelity magic videos. We define the core capability operationalized by MagicBench as **Autonomous Physical Entailment (APE)**, defined as the capacity to maintain a consistent causal world model when linguistic narratives actively contradict visual kinematics. Magic serves as a uniquely rigorous testbed for APE because the *patter* is designed to decouple symbolic narratives from physical reality.

**Taxonomy and Annotation Rigor.** The dataset is categorized by hierarchical linguistic interference: (1) **Type A (Direct Lie,  $N = 146$ )**: Explicit falsification of the visual state. (2) **Type B (Misdirection,  $N = 91$ )**: Attentional diversion to non-critical regions. (3) **Type C (Patter,  $N = 83$ )**: High-entropy irrelevant conversational filler.

(4) **Type D (Silence,  $N = 82$ )**: A critical control group representing a Semantic Vacuum.

Each video is grounded in a Physical Constraint Set (PCS), which is a formalization of objective facts revealed in tutorial segments. To ensure logical fidelity, three experts independently drafted physical steps, resolved via a Consensus Protocol to maximize the recall of sub-second maneuvers.

To quantify objectivity, a reliability audit on 30 samples yielded a **BERTScore F1 of 0.87** (Zhang et al., 2020), confirming that experts captured the same underlying logic despite phrasing variations. On average, each PCS instance comprises  $5.7 \pm 1.4$  **atomic physical constraints**, with a mean length of  $64.3 \pm 8.5$  tokens.

### 3.2 Sampling and Sensory Sufficiency

Magic tricks hinge on transient maneuvers. We employ an Adaptive Hierarchical Sampling strategy (4fps for videos  $< 20$ s) to capture these cues within MLLM context limits.

**Recall and Existence Proof.** A human audit on 60 videos confirms that critical manipulation frames were captured in **98.3%** of cases (95% CI: [91.1%, 99.7%]). To assess whether temporal resolution (4fps) could fully account for the observed gap, we utilize the Vision-Only capability probe as an internal existence proof. If 4fps were insufficient to capture micro-sleights, the vision-only model would fail entirely. Instead, models like Gemini-2.5-Pro achieve a robust 50.7% CPA when forced to rely solely on vision. Since the Multimodal setting receives the exact same 4fps input yet collapses to significantly lower performance (e.g., GPT-4o dropping to 15.8% in Silence), we take this as evidence that the sampled visual input retains sufficient kinematic information in many cases. These results suggest that missing frames are unlikely to be the sole explanation for the observed gap. Instead, the evidence is more consistent with an access-related failure in utilizing available visual evidence under multimodal conditions.

### 3.3 Evaluation Protocol: Physical Entailment

Standard benchmarks rely on textual similarity, which fails to capture causal logic. We propose an **LLM-as-a-Judge** protocol grounded in **Physical Entailment**, assessing reasoning against three First-Order Logic (FOL) invariants: Object Permanence, Kinematic Continuity, and the Visual Supremacy Axiom (Appendix D).

<b>Step 1: The Zero-Tolerance Veto (Pre-requisite)</b>			
Score $\rightarrow 0$ if: (1) Prop Hallucination; (2) State Contradiction; (3) Magical Thinking.			
<b>Step 2: Scoring Rubric (Scale 0 – 10)</b>			
Metric	High (7–10)	Mid (4–6)	Low (0–3)
<b>VTG</b>	<b>Precise.</b> Correctly identifies Body Part, Specific Action, and Temporal Ordering aligned with PCS.	<b>Partial.</b> Identifies Body Part but is vague on Action or Timing; or correct moment but wrong location.	<b>Fail.</b> Wrong location, hallucinates visual events, or provides completely vague descriptions.
<b>CPA</b>	<b>Entailed.</b> Explanation is physically entailed by PCS (Terminology Agnostic). Describes correct mechanics.	<b>Plausible.</b> Identifies general category (e.g., Sleight of Hand) but misses key mechanical steps.	<b>Fail.</b> Triggered Hard Failure Veto OR explains a completely wrong mechanism.
<b>CFR</b>	<b>Resilient.</b> Explicitly refutes the audio lie OR provides a physical explanation contradicting the lie.	<b>Skeptical.</b> Expresses doubt (e.g., "unlikely") but fails to provide a definite physical alternative.	<b>Foiled.</b> Accepts the illusion or lie as fact (e.g., "it melted") or invents props to justify the lie.

Table 1: **The PCS-based Judgment Protocol.** Models must first pass the *Zero-Tolerance Veto*. If passed, they are graded on a 0–10 scale across three dimensions based on the degree of alignment with the Physical Constraint Set (PCS).

**Blind-Logic Adjudication.** To prevent circular reasoning, our protocol is **strictly decoupled**: the judge receives only the model’s textual output and symbolic PCS facts, but never the raw video. This transforms the judge from a subjective observer into a deterministic logic checker. Grading follows a two-stage rubric: a Zero-Tolerance Veto for physical impossibilities, followed by a Scalar Rubric (0–10) for completeness (see Table 1).

**Reliability and Bias Mitigation.** To ensure statistical rigor and address *self-preference bias*, we expanded human validation to a stratified sample of  $N = 80$  (20% of the dataset). As shown in Table 6, the protocol shows strong reliability: (1) **Cross-Architecture Objectivity**: An independent judge (*Claude-3.5-Sonnet*) achieved high correlation with human experts (Avg.  $r = 0.8398$ ), mirroring GPT-4o’s performance ( $r = 0.878, p < 0.001$ ). This suggests that adjudication is driven primarily by symbolic PCS axioms rather than model identity. (2) **Targeted Veto Audit**: We evaluated 20 cases where model explanations were fundamentally unphysical (vetted as CPA=0 by Ground Truth). Without PCS grounding, a vanilla GPT-4o judge exhibited an 85.0% False Positive Rate (FPR), accepting physical impossibilities as plausible (Table 2). Our PCS-grounded judge maintained a 0.0% FPR. (3) **Noise and Persona Robustness**: On  $N = 50$  samples, introducing 10% spatial jitter into PCS coordinates yielded negligible variance (MAPE=5.4%,  $p = 0.38$ ). Across the entire dataset ( $N = 402$ ), scores from distinct judge personas achieved a Pearson  $r = 0.7949$  ( $p < 0.001$ ), with 100% identical veto decisions.

Judge Setting	Avg. Score ( $\uparrow$ )	FPR ( $\downarrow$ )
Vanilla GPT-4o (No PCS)	5.3/10	85.0%
<b>PCS-Grounded (Ours)</b>	<b>0.0/10</b>	<b>0.0%</b>

Table 2: Ablation of PCS grounding ( $N = 20$ ). FPR denotes the rate of non-zero scores assigned to physically impossible explanations. The PCS protocol successfully mitigates the judge’s inherent bias toward accepting illusory narratives.

### 3.4 Metrics

We report performance across three dimensions: (1) **Visual-Temporal Grounding (VTG)**: Spatial and temporal localization precision. (2) **Causal-Physical Accuracy (CPA)**: Adherence to objective physical laws (the primary APE metric). (3) **Counterfactual Resilience (CFR)**: Robustness against the compounded effect of visual trickery and linguistic deception.

## 4 Experiments

### 4.1 Experimental Setup

**Models.** We evaluate three strong MLLMs representing diverse architectures: (1) **GPT-4o**, a commercial multimodal model; (2) **Gemini-2.5-Pro**, a native multimodal model; and (3) **Qwen2.5-VL-72B**, a strong open-source multimodal baseline. All inferences are conducted at a temperature of 0.6.

**Input Specification.** To isolate semantic reasoning from low-level auditory perception errors, we utilize **ground-truth textual transcripts** as the proxy for the audio modality. In Type D (Silence),

Model	Setting	Type A (Lie)	Type B (Misdir.)	Type C (Patter)	Type D (Silence)	Avg.
<b>GPT-4o</b>	Vision-Only (Probe)	33.0 $\pm$ 1.1	38.1 $\pm$ 1.5	34.5 $\pm$ 1.2	<b>28.2<math>\pm</math>1.0</b>	33.7
	Multi-Forensic (SIA)	<b>43.4<math>\pm</math>1.3</b>	<b>43.3<math>\pm</math>1.6</b>	33.6 $\pm$ 1.1	15.8 $\pm$ 0.9 $\dagger$	<b>38.1</b>
<b>Gemini-2.5-Pro</b>	Vision-Only (Probe)	47.3 $\pm$ 1.4	46.3 $\pm$ 1.2	46.3 $\pm$ 1.5	<b>50.7<math>\pm</math>1.3</b>	<b>47.6</b>
	Multi-Forensic (SIA)	42.1 $\pm$ 1.4	45.8 $\pm$ 1.2	38.4 $\pm$ 1.1	37.8 $\pm$ 1.1 $\dagger$	41.0
<b>Qwen2.5-VL-72B</b>	Vision-Only (Probe)	32.4 $\pm$ 1.0	34.5 $\pm$ 1.3	<b>35.7<math>\pm</math>1.2</b>	<b>30.6<math>\pm</math>1.1</b>	33.3
	Multi-Forensic (SIA)	<b>48.7<math>\pm</math>1.4</b>	<b>40.3<math>\pm</math>1.6</b>	32.3 $\pm$ 1.0	22.7 $\pm$ 1.2 $\dagger$	<b>36.0</b>
<i>Upper Bound (GPT-4o)</i>	Truth Context	86.7	87.0	87.0	86.3	86.8

Table 3: Main CPA results across models and linguistic conditions. (a) In deceptive scenarios (Type B), linguistic entities can act as semantic anchors that facilitate visual grounding, leading to the *Spotlight Effect*. (b) In semantic vacuums (Type D), the absence of linguistic triggers is associated with a sharp drop in performance relative to each model’s vision-only capacity, consistent with *Visual Agency Loss* (the *Crutch Effect*). Error bars denote standard error estimated from 1,000 bootstrap iterations.  $\dagger$  denotes a statistically significant degradation compared to the Vision-Only probe ( $p < 0.01$ ).

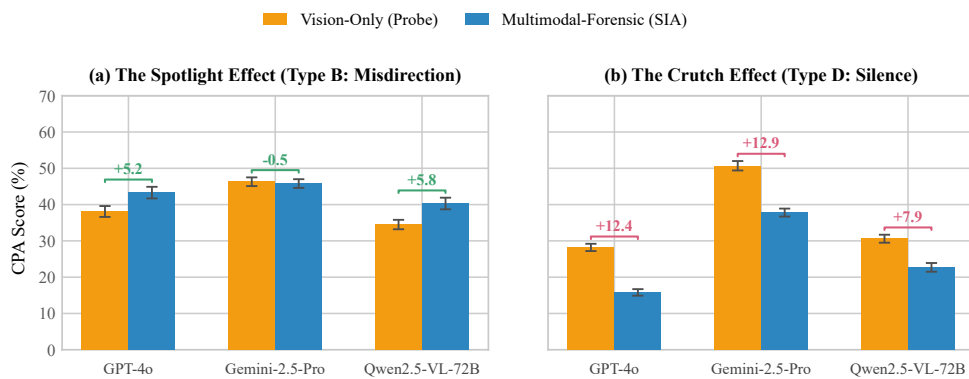


Figure 2: Visualization of the Semantic Dependency Paradox. (a) **Spotlight Effect (Type B)**: Linguistic entities can act as semantic anchors that facilitate visual grounding, allowing multimodal models to outperform the vision-only baseline. (b) **Crutch Effect (Type D)**: In semantic vacuums, all models show a substantial drop relative to their vision-only capacity, consistent with *Visual Agency Loss*. Red arrows highlight the residual performance gap under multimodal prompting.

the audio field is an **empty string**, while the system prompt continues to mandate multimodal cross-referencing.

**Symmetric Intensity Ablation (SIA)**. To ensure a fair modality-level comparison and address the confound of instruction design, we design five settings:

- **Vision-Only (Capability Probe)**: The model receives only video frames with a high-intensity forensic prompt, instructing it to ignore implied audio and rely strictly on pixels.
- **Multimodal (Neutral)**: Default behavior using generic descriptive instructions.
- **Multimodal (Skeptic)**: Content-focused intervention explicitly warning the model of deception.
- **Multimodal (Forensic) [SIA]**: Our primary experimental group uses a prompt designed to

be functionally symmetric to the Vision-Only probe, reducing prompt-design confounds while keeping instruction intensity approximately constant.<sup>1</sup>

- **Truth Context**: Deceptive transcripts are replaced with ground-truth physical principles to establish the theoretical reasoning ceiling.

## 4.2 Main Results: The Semantic Dependency Paradox

Table 3 summarizes performance across models and interference types. We highlight three observations based on Causal-Physical Accuracy (CPA):

### 1. Semantic Anchoring (The Spotlight Effect).

In deceptive scenarios (Type A/B), the **Multimodal**

<sup>1</sup>After removing the opening role-definition sentence, we assess prompt similarity in two ways: (1) BERTScore F1 over the remaining instruction text; and (2) a set-level Instruction Overlap F1 computed over manually extracted high-level logical and stylistic constraint items. The latter yields 0.9333.

**(Forensic)** setting consistently surpasses the *Vision-Only* baseline across all models. For instance, GPT-4o improves from 33.0% to 43.4% in Type A, and Qwen2.5-VL-72B improves from 32.4% to 48.7%. As shown by the component ablation in Sec. 5.1, this gain is consistent with entity nouns acting as semantic anchors that narrow the visual search space, allowing models to localize objects despite false predicates.

**2. Performance Inversion in Silence (The Crutch Effect).** In Type D (Silence), we observe a sharp drop in multimodal performance. Crucially, this occurs under our **Symmetric Intensity Ablation (SIA)** protocol: even when the Multimodal setting is explicitly instructed to act as a “skeptical forensic auditor” using a prompt functionally identical to the unimodal baseline (Instruction Overlap F1 = 0.93), all models exhibit a statistically significant performance gap compared to the Vision-Only probe (e.g., **12.4% drop for GPT-4o**,  $p < 0.01$ ). This suggests that even an inactive audio channel is associated with **Visual Agency Loss**: despite matched instructions, the model tends to enter a state of perceptual idling when linguistic support is absent.

**3. Cross-Model Stability of the Paradox.** As visualized in Figure 2, the paradox appears consistently across the tested early-fusion MLLMs. While the magnitude of the Spotlight gain varies by model, the **Crutch Effect** remains consistent across all three models. Even the high-performing Gemini-2.5-Pro suffers a 12.9% drop in Type D relative to its own intrinsic capacity.

### 4.3 Causal Evidence: Spatial Specificity

To quantify the *access bottleneck* hypothesis, we performed the spatial intervention (Figure 3) on  $N = 60$  failure cases. By highlighting ground-truth coordinates defined by the PCS, reasoning accuracy (CPA) for Qwen2.5-VL-72B recovered from **16.8% to 29.8%**. This recovery is accompanied by an increase in **Relative Spatial Density (RSD)** from 0.47 to 2.66 within the critical region. Together, these results are more consistent with a perceptual access bottleneck than with a lack of reasoning competence.

## 5 Diagnosis: The Spotlight–Crutch Framework

We analyze the behavioral patterns underlying the *Semantic Dependency Paradox*. Across the tested

models, visual reasoning appears disproportionately dependent on linguistic activation.

### 5.1 Dissecting the Spotlight: Anchoring vs. Alerting

In deceptive scenarios, multimodal performance improves over the unimodal baseline. For instance, in Type B (Misdirection), the *Multimodal (Skeptical)* setting achieves **44.0%**, surpassing the *Vision-Only* probe of **38.1%**. We attribute this to semantic anchoring: while linguistic *predicates* are false (e.g., “melted”), the *entities* (e.g., “coin”) act as spatial priors that narrow the visual search space.

#### Mechanism Isolation: Component Ablation.

To verify whether this gain stems from specific semantic guidance or generic attentional arousal, we conducted a controlled ablation on  $N = 50$  Type A samples (Table 4). The results distinguish two patterns: (1) **Alerting**: Generic cues (“Look here”) restore performance only to the visual baseline ( $33.6\% \approx 33.0\%$ ), suggesting that simple audio presence may overcome silence-induced inertia but does not by itself guide the visual encoder effectively. (2) **Anchoring**: The strongest gains occur in conditions containing specific nouns (44.2%). This pattern is consistent with the use of entity cues for localization while discounting false predicates. By contrast, high-entropy noise (Type C) reduces performance.

Text Condition	Function	CPA (%)	vs. Baseline
Vision-Only (Baseline)	–	33.0	–
Generic Cues (“Look here”)	<i>Alerting</i>	33.6	+0.6 (n.s.)
Type C (Irrelevant Patter)	<i>Interference</i>	25.6	↓ <b>7.4</b>
<b>Entity Keywords (“Coin”)</b>	<i>Anchoring</i>	<b>44.2</b>	↑ <b>11.2</b>
Full Deceptive (Original)	<i>Conflict</i>	43.4	↑ 10.4

Table 4: Component ablation in Type A scenarios ( $N = 50$ ). The null result for Generic Cues supports a noun-based anchoring account over generic attentional arousal, though matched-information alternatives are not fully ruled out.

### 5.2 Visual Agency Loss (The Crutch Effect)

The collapse in Type D (Silence) provides strong evidence for **Visual Agency Loss**, that is, modality-driven perceptual inertia. Comparing *Capacity* (Vision-Only) with *Behavior* (Multimodal) reveals a substantial gap: a score of 28.2% in the probe setting indicates that the models retain the capacity to decode the trick, yet this capacity remains largely latent (4.4%) under standard interaction. This pat-

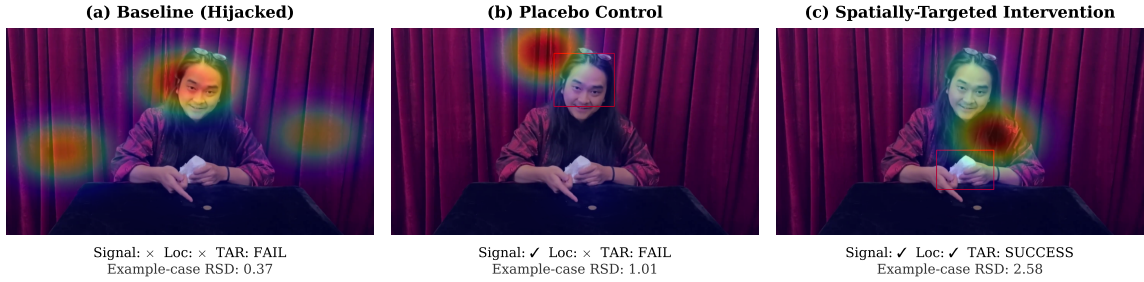


Figure 3: Causal diagnosis of the access bottleneck via spatial intervention. We analyze Visual Agency Loss using a triplet control ( $N = 60$  for Qwen2.5-VL-72B). **(a) Baseline (Hijacked):** Deceptive audio diverts attention from critical regions (mean RSD  $\approx 0.47$ ), causing reasoning failure. **(b) Placebo Control:** Highlighting non-critical regions fails to restore CPA, ruling out generic attention-arousal artifacts. **(c) Targeted Intervention:** Precise spatial prompting, analogous to an oracle-style visual prompting intervention in the spirit of VTPrompt (Xiao et al., 2025), restores CPA (16.8%  $\rightarrow$  29.8%) and increases RSD (to 2.66). These results are more consistent with a perceptual access bottleneck than with a pure reasoning deficit. Detailed analysis is provided in Sec. 5.3.

tern is consistent with the **Language-Guided Observer Hypothesis**: due to instruction-tuning biases, models may default to perceptual idling when linguistic triggers are absent. It is also compatible with the “Bad Eyes” phenomenon (Aravindan et al., 2025), in which visual information is hypothesized to be misrouted during cross-modal fusion (Liu et al., 2025; Nikankin et al., 2025).

### 5.3 Causal Diagnosis via Perceptual Access Restoration

To distinguish between reasoning deficits and an access bottleneck, we conduct two complementary interventions: exogenous spatial prompting (top-down) and signal enhancement (bottom-up).

**Restoring Access.** Highlighting critical coordinates (Red Circle) for  $N = 60$  failure cases caused Qwen2.5-VL-72B’s CPA to recover from **16.8% to 29.8%**, coinciding with a surge in **Relative Spatial Density (RSD)** from **0.47 to 2.66**. We also performed a targeted stress test on  $N = 20$  high-discrepancy samples (Multimodal CPA  $< 3.0$  vs. Vision-Only CPA  $> 7.0$ ). Magnifying critical regions increased CPA from 11.3% to 63.2% ( $p < 0.001$ ). This suggests an important possibility: **Visual Override** (Ortu et al., 2025) can occur, but may require a disproportionately high perceptual threshold in adversarial contexts. In standard settings, the observed behavior is consistent with a form of **Functional Perception Suppression** in which visual features are under-weighted or insufficiently accessed during multimodal integration.

### 5.4 Disentangling Confounding Factors

We address concerns regarding token budgets and prompt artifacts through two logical controls: (1) **Resource Control:** In Type D, the audio channel is an empty string, substantially reducing concerns about token congestion. The fact that performance is *lowest* here suggests that compute budget alone is unlikely to explain the observed failure. (2) **Symmetric Intensity Ablation (SIA):** A common critique is that performance drops are artifacts of asymmetric prompt design. By comparing the *Vision-Only* probe against the *Multimodal (Forensic)* setting using identical, high-intensity skeletal instructions, we hold the prompt variable constant. The persistence of a **12.4% Residual Gap** under these symmetric conditions suggests that Visual Agency Loss is more likely rooted in cross-modal fusion than in a transient prompting artifact.

### 5.5 Limits of Inference-Time Mitigation (IoT)

Finally, we investigate whether advanced inference strategies can mitigate this limitation. We implement an **Image-of-Thought (IoT)** pipeline (An et al., 2025), mandating a visual inventory before reasoning.

Strategy	Type A	Type B	Type C	Type D	Avg.
Standard (Neutral)	35.6	34.5	25.6	4.4	25.0
Forensic (SIA)	43.4	43.3	33.6	15.8	34.0
<b>IoT (Inference)</b>	<b>46.2</b>	<b>45.8</b>	<b>36.4</b>	21.2 $\dagger$	<b>37.9</b>
<b>Vision-Only</b>	33.0	38.1	34.5	<b>28.2</b>	33.5

Table 5: **Inference Strategy Comparison** ( $N = 402$ ). While IoT improves performance, it fails to close the 7.0% gap with the Vision-Only probe in Type D ( $p < 0.01$ , denoted by  $\dagger$ ).

As shown in Table 5, while IoT achieves peak multimodal performance (**46.2%** in Type A) by synergizing with semantic anchors, it **does not eliminate the gap** in semantic vacuums (**21.2% vs. 28.2%**). This pattern suggests that modality dependency may arise, at least in part, from the under-utilization of visual evidence during multimodal integration, and is not fully alleviated by inference-time heuristics.

## 6 Discussion

The *Semantic Dependency Paradox* highlights a limitation of the common assumption that richer context always improves multimodal reasoning. Our findings suggest that the tested MLLMs often behave more like **language-guided observers** than fully integrated reasoners, tending to prioritize linguistic priors over sensory evidence.

### 6.1 Implications for Multimodal Integration

Our behavioral diagnosis provides macro-scale evidence consistent with the **competing circuits** identified by Nikankin et al. (2025) and the **functional suppression** described by Hua et al. (2025). The failure of even the advanced IoT pipeline (Sec. 5.5) to bridge the gap to the Vision-Only probe suggests that the bottleneck likely arises at the **representation-fusion level**. This points toward a functional limitation in the tested early-fusion architectures, suggesting that high-entropy narratives can effectively “blind” or misdirect the visual encoder. One possible direction is to introduce a more explicitly decoupled visual-first stage, in which physical evidence is extracted before strong linguistic conditioning is applied.

### 6.2 Training Implications

The *Crutch Effect* suggests an emergent consequence of instruction-tuning: **Modality Bias** (Nikankin et al., 2025). By training on perfectly aligned alt-text data, MLLMs learn to treat language as a reliable shortcut, resulting in a failure of *Active Perception* when linguistic guidance is absent. To cultivate more independent agency, future training paradigms may benefit from moving toward **Adversarial Modality Alignment**.

One possible training direction is to incorporate mismatched semantic pairs and silent videos during pre-training. This could encourage models to resolve cross-modal conflict through visual verification. More generally, a robust multimodal agent

should be able to reject a linguistic narrative when it conflicts with visual evidence and physical constraints.

### 6.3 Broader Implications

Although magic is a specific domain, the **Visual Agency Loss** observed here may be useful as a stress-test indicator for failures in high-stakes multimodal settings. A magician’s deceptive patter is structurally analogous to a semantic adversarial setting, in which persuasive language conflicts with observable physical reality. In this sense, MagicBench can serve as a canary-style stress test for safety-relevant multimodal robustness.

If a multimodal system cannot maintain *Autonomous Physical Entailment* (APE) in a simple sleight-of-hand setting, it may also be vulnerable to cross-modal conflict in real applications. For example, an autonomous driving system could receive an erroneous audio navigation command while its visual sensors detect a stopped vehicle. Likewise, a household robot may receive an instruction that conflicts with the physical state of an object. In such cases, the system should be able to prioritize visual evidence over misleading language.

More broadly, MagicBench should be viewed as a controlled diagnostic setting for studying cross-modal conflict resolution. In that setting, **perceptual autonomy**, the ability of visual evidence to override misleading textual priors, is an important safety-relevant property.

## 7 Conclusion

In this work, we introduced **MagicBench**, a diagnostic framework for studying Multimodal LLMs under semantic adversarial conflicts. Our systematic evaluation reveals a clear **Semantic Dependency Paradox**: while linguistic anchors can facilitate object localization (the *Spotlight Effect*), the absence of these anchors triggers a pronounced collapse in performance (the *Crutch Effect*), revealing a persistent **Visual Agency Loss**.

Using the *Vision-Only* setting as a capability probe, we quantify a substantial gap between latent perceptual capacity and operational multimodal behavior. Through interventional analyses such as spatial prompting and signal enhancement, we provide evidence that this failure is better explained by a Perceptual Access Bottleneck rather than a pure reasoning deficit. Specifically, visual features appear to be functionally suppressed or under-

accessed under linguistic dominance. Crucially, we find that even advanced inference-time strategies like Image-of-Thought do not fully bridge this gap, suggesting that the root cause may lie within the early-fusion integration layer.

Although the benchmark is specific to magic, it provides a controlled way to measure a model’s perceptual autonomy under cross-modal conflict. More broadly, our results suggest that future multimodal systems may need stronger mechanisms for allowing visual evidence to override misleading language.

## Limitations

While **MagicBench** identifies a systematic pathology in MLLM perception, we acknowledge several constraints that define the scope of our findings:

### Ecological Validity and Domain Generalization.

Our study utilizes magic as a proxy for **Semantic Adversarial Attacks**. While magic represents a “worst-case” stress test for cognitive conflict, the visual distribution (e.g., stage lighting, specific props) differs from safety-critical domains like autonomous driving. Future work is required to verify if the *Semantic Dependency Paradox* manifests with equal severity in scenarios with lower semantic density but higher safety stakes. We propose that MagicBench serves as a **Canary Task** to measure *Autonomous Physical Entailment* (APE), but it does not replace domain-specific safety evaluations.

**Temporal Resolution vs. Context Limits.** Despite our **98.3% Sampling Recall** and the *Vision-Only* existence proof (Sec. 3.2), the 4fps sampling rate remains an inherent bottleneck. While sufficient for capturing the discrete states required for our PCS logic, extremely rapid micro-sleights or simultaneous multi-object tracking over extended durations may be compromised. This reflects a fundamental trade-off between temporal resolution and the context window limits of current MLLMs. As long-context architectures mature, evaluating MagicBench at native frame rates (e.g., 30fps) will be a critical next step.

**Behavioral Proxies vs. Mechanistic Proofs.** In our interventional diagnosis (Sec. 5.3), we utilized Grad-CAM and RSD maps as behavioral proxies to visualize attention shifts. We caution that these visualizations are aggregative and do not mathematically isolate the specific attention heads or sparse

autoencoder features responsible for token suppression. For closed-source models like GPT-4o, our findings remain a clinical behavioral diagnosis. Future research utilizing white-box probing on open-source weights (Nikankin et al., 2025) is needed to map the precise neural circuits that govern the *Crutch Effect*.

**The Confound of Alignment Bias.** The observed **Visual Agency Loss** may be partially exacerbated by current Supervised Fine-Tuning (SFT) paradigms, which heavily penalize hallucinations by enforcing strict adherence to linguistic prompts. Our study focuses on deployed models and does not decouple “nature” (base model perceptual capacity) from “nurture” (instruction-following bias). Determining whether the *Crutch Effect* is an architectural inevitability of early-fusion or an artifact of the “Context is King” training objective remains an open scientific question.

**Adjudication Ceiling.** While our PCS-grounded judge demonstrates high reliability ( $r = 0.88$ ) and robustness to noise, the evaluation protocol still relies on an LLM backbone. Although we mitigate self-preference bias through symbolic facts and cross-model validation (Claude-3.5), the adjudication remains subject to the inherent reasoning ceiling of the judge model itself.

## Acknowledgements

This work was supported by the National Training Program of Innovation and Entrepreneurship for Undergraduates.

We sincerely thank the anonymous reviewers for their rigorous evaluation and constructive feedback, which significantly strengthened the interventional analyses and overall quality of this manuscript.

## References

- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob L. Menick, Sebastian Borgeaud, and 8 others. 2022. [Flamingo: a visual language model for few-shot learning](#). In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.
- Tuo An, Yunjiao Zhou, Han Zou, and Jianfei Yang. 2025. [IoT-LLM: a framework for enhancing large language](#)

- model reasoning from real-world sensor data. *arXiv preprint arXiv:2410.02429*.
- Ashwath Vaithinathan Aravindan, Abha Jha, and Mihir Kulkarni. 2025. Do VLMs have bad eyes? diagnosing compositional failures via mechanistic interpretability. *arXiv preprint arXiv:2508.16652*.
- Ruzena Bajcsy, Yiannis Aloimonos, and John K. Tsotsos. 2016. Revisiting active perception. *Autonomous Robots*, 42:177 – 196.
- Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, and Shibo Wang. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, abs/2403.05530.
- Michal Golovanevsky, William Rudman, Michael Lepori, Amir Bar, Ritambhara Singh, and Carsten Eickhoff. 2025. Pixels versus priors: Controlling knowledge priors in vision-language models through visual counterfactuals. *arXiv preprint arXiv:2505.17127*.
- Tianze Hua, Tian Yun, and Ellie Pavlick. 2025. How do vision-language models process conflicting information across modalities? *arXiv preprint arXiv:2507.01790*, abs/2507.01790.
- 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 *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*, pages 13872–13882. IEEE.
- KunChang Li, Yanan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. 2024. VideoChat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*.
- Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Xin Zhao, and Ji-Rong Wen. 2023. Evaluating object hallucination in large vision-language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 292–305, Singapore. Association for Computational Linguistics.
- Benlin Liu, Amita Kamath, Madeleine Grundle-McLaughlin, Winson Han, and Ranjay Krishna. 2025. Visual representations inside the language model. *arXiv preprint arXiv:2510.04819*, abs/2510.04819.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.
- Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. 2024. Video-ChatGPT: Towards detailed video understanding via large vision and language models. *arXiv preprint arXiv:2306.05424*.
- Yaniv Nikankin, Dana Arad, Yossi Gandelsman, and Yonatan Belinkov. 2025. Same task, different circuits: Disentangling modality-specific mechanisms in VLMs. *arXiv preprint arXiv:2506.09047*, abs/2506.09047.
- Francesco Ortu, Zhijing Jin, Diego Doimo, and Alberto Cazzaniga. 2025. When seeing overrides knowing: Disentangling knowledge conflicts in vision-language models. *arXiv preprint arXiv:2507.13868*, abs/2507.13868.
- William Rudman, Michal Golovanevsky, Dana Arad, Yonatan Belinkov, Ritambhara Singh, Carsten Eickhoff, and Kyle Mahowald. 2026. Mechanisms of prompt-induced hallucination in vision-language models. *arXiv preprint arXiv:2601.05201*, abs/2601.05201.
- Baifeng Shi, Siyu Gai, Trevor Darrell, and Xin Wang. 2023. Toast: Transfer learning via attention steering. *arXiv preprint arXiv:2305.15542*, abs/2305.15542.
- Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, and Jifeng Dai. 2023. VisionLLM: Large language model is also an open-ended decoder for vision-centric tasks. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.
- Xi Xiao, Yunbei Zhang, Yanshu Li, Xingjian Li, Tianyang Wang, Jihun Hamm, Xiao Wang, and Min Xu. 2025. Visual variational autoencoder prompt tuning. *arXiv preprint arXiv:2503.17650*.
- Mingjie Xu, Jinpeng Chen, Yuzhi Zhao, Jason Chun Lok Li, Yue Qiu, Zekang Du, Mengyang Wu, Pingping Zhang, Kun Li, and Hongzheng Yang. 2025. VP-Bench: A comprehensive benchmark for visual prompting in multimodal large language models. *arXiv preprint arXiv:2511.11438*, abs/2511.11438.
- Kexin Yi, Chuang Gan, Yunzhu Li, Pushmeet Kohli, Jiajun Wu, Antonio Torralba, and Joshua B. Tenenbaum. 2020. CLEVRER: collision events for video representation and reasoning. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Shukang Yin, Chaoyou Fu, Sirui Zhao, Tong Xu, Hao Wang, Dianbo Sui, Yunhang Shen, Ke Li, Xing Sun, and Enhong Chen. 2024. Woodpecker: hallucination correction for multimodal large language models. *Science China Information Sciences*, 67(12).
- Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. 2024. MM-Vet: Evaluating large multimodal models for integrated capabilities. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. [BERTScore: Evaluating text generation with BERT](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.

## A Evaluation Protocol Validation

To ensure the reliability of our GPT-4o-based judge and mitigate potential self-preference bias, we conducted a rigorous meta-evaluation on a stratified subset of 80 samples, covering all linguistic types and performance ranges.

**Methodology.** To establish a robust human baseline, **three expert annotators** independently graded the 80 samples based on the **Physical Constraint Set (PCS)**. The final human score for each sample was derived by averaging the three annotators’ scores to mitigate individual subjectivity. We then compared these consensus human scores with those generated by **GPT-4o (Our Judge)** and **Claude-3.5-Sonnet**.

**Results.** As shown in Table 6, our GPT-4o judge demonstrates an exceptionally high correlation with human judgment across all three metrics (Pearson  $r > 0.86$ ,  $p < 0.001$ ). Notably, in the critical **CPA** metric, GPT-4o achieves a correlation of **0.88**, surpassing Claude-3.5 ( $r = 0.79$ ). This confirms that GPT-4o aligns more closely with the strict physical constraints defined in PCS. The strong alignment across different evaluators validates that our scoring criteria are objective and reproducible, effectively ruling out stochasticity or model-specific bias.

Metric	Pearson Correlation ( $r$ )		
	Human vs. GPT-4o	Human vs. Claude	GPT-4o vs. Claude
VTG (Visual Grounding)	<b>0.8722</b>	0.8669	0.7936
CPA (Physical Accuracy)	<b>0.8825</b>	0.7899	0.7328
CFR (Resilience)	<b>0.8793</b>	0.8627	0.8198
<i>Average</i>	<b>0.8780</b>	0.8398	0.7854

Table 6: **Meta-evaluation of the Judge.** The high correlation between Human and GPT-4o (Avg.  $r = 0.87$ ) validates the reliability of our automated evaluation protocol. All correlations are statistically significant ( $p < 0.001$ ).

## B Data Quality & Sampling Audit

A critical concern regarding video LLMs is whether sparse sampling strategies capture the fleeting moments essential for magic tricks (e.g., the precise frame of a false transfer). To address this, we conducted a rigorous audit of our **Adaptive Hierarchical Sampling** strategy.

**Methodology.** We randomly sampled  $N = 60$  videos from the MagicBench dataset. For each video, an expert annotator identified the *Critical Time Window (CTW)*, defined as the specific duration where the core manipulation occurs (typically 0.5s  $\sim$  2.0s). We then verified whether the sampled frames fed to the model included at least one clear frame within this CTW.

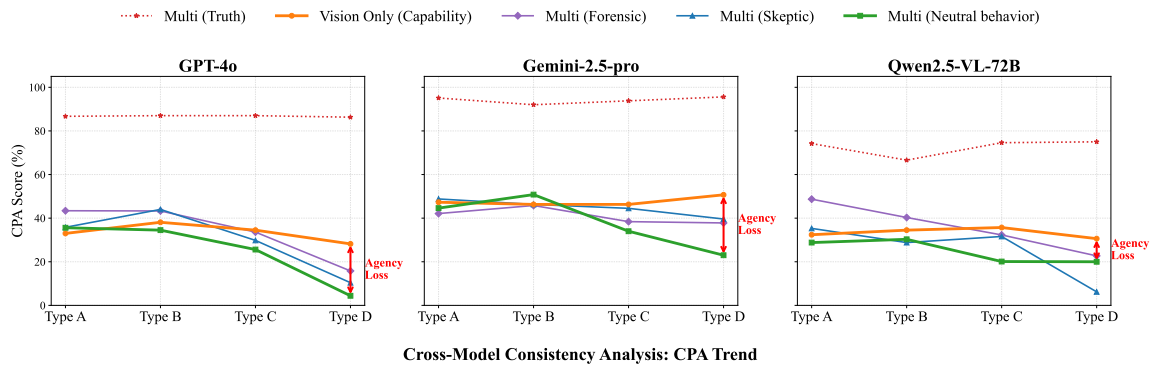
**Results.** The audit reveals an exceptionally high **Critical Frame Recall Rate of 98.3%**, with the critical moment successfully captured in 59 out of 60 videos.

$$\text{Recall} = \frac{\text{Captured Samples}}{\text{Total Samples}} = \frac{59}{60} \approx 98.3\% \quad (\text{B.1})$$

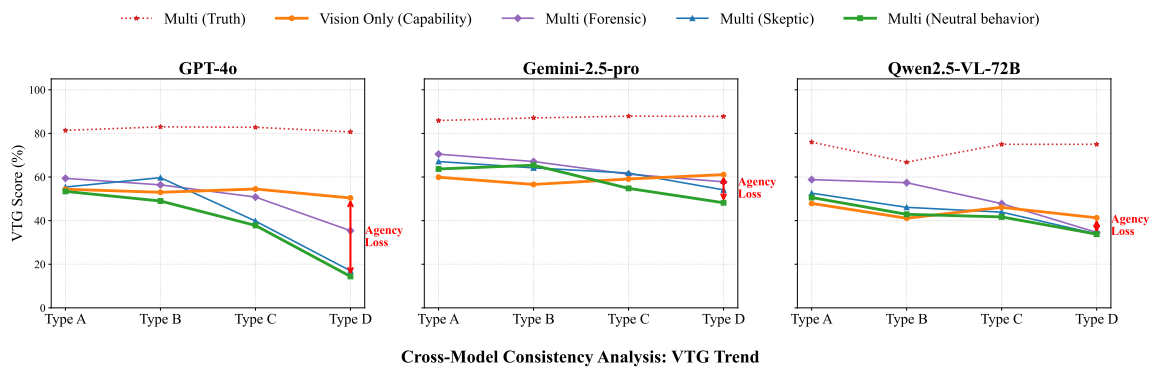
The 95% Confidence Interval (Wilson Score Interval) is [91.1%, 99.7%]. This confirms that in the vast majority of cases, reasoning failures are attributable to the model’s cognitive deficits (e.g., Attention Hijacking or Agency Loss) rather than missing visual evidence caused by downsampling.

## C Detailed Experimental Results

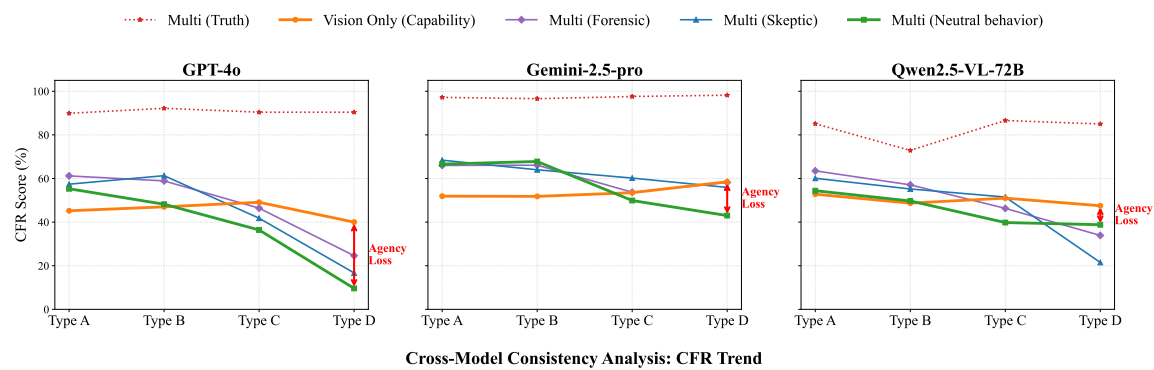
Due to space constraints, the main text reports aggregated performance. Table 7 provides a comprehensive breakdown of Visual-Temporal Grounding (VTG), Causal-Physical Accuracy (CPA), and Counterfactual Resilience (CFR) across all models and linguistic types. Furthermore, Figure 4 visualizes the cross-model consistency analysis, highlighting the universal nature of the Visual Agency Loss across different architectures.



**(a) Causal-Physical Accuracy (CPA) Trends.** Note the consistent "Crutch Effect" in Type D (Silence), where Vision-Only (Orange) outperforms Multimodal (Green/Blue) across all architectures.



**(b) Visual-Temporal Grounding (VTG) Trends.** Visual localization capability also degrades in silence for multimodal settings, indicating that the agency loss is perceptual, not just logical.



**(c) Counterfactual Resilience (CFR) Trends.** Vision-Only models show consistently higher resilience to confusion in the absence of audio cues.

**Figure 4: Cross-Model Consistency Analysis.** We visualize the robustness trends across GPT-4o, Gemini-2.5-Pro, and Qwen2.5-VL-72B. The red arrows highlight the **Visual Agency Loss** gap in Type D scenarios, confirming that this is a fundamental architectural pathology common to SOTA MLLMs.

Model	Setting	VTG				CPA				CFR			
		A	B	C	D	A	B	C	D	A	B	C	D
GPT-4o	Vision-Only	54.4	53.0	54.5	<b>50.4</b>	33.0	38.1	34.5	<b>28.2</b>	45.2	47.0	49.1	<b>40.0</b>
	Multi (Neutral)	53.4	49.0	37.8	14.4	35.6	34.5	25.6	4.4	55.3	48.2	36.4	9.6
	Multi (Skeptic)	55.4	59.7	39.8	17.0	35.8	44.0	29.8	10.4	57.4	61.3	41.8	16.7
	Multi (Forensic)	59.4	56.4	50.8	35.4	43.4	43.3	33.6	15.8	61.2	58.9	46.4	24.6
	<i>Truth Context</i>	<i>81.4</i>	<i>83.0</i>	<i>82.8</i>	<i>80.7</i>	<i>86.7</i>	<i>87.0</i>	<i>87.0</i>	<i>86.3</i>	<i>89.9</i>	<i>92.2</i>	<i>90.4</i>	<i>90.4</i>
Gemini-2.5-Pro	Vision-Only	59.9	56.6	59.1	<b>61.1</b>	47.3	46.3	46.3	<b>50.7</b>	51.9	51.8	53.5	<b>58.5</b>
	Multi (Neutral)	63.7	65.4	54.8	48.2	44.6	50.8	34.0	23.0	66.6	67.8	49.9	43.0
	Multi (Skeptic)	67.1	64.2	61.8	54.1	48.8	46.3	44.5	39.6	68.4	64.0	60.2	55.9
	Multi (Forensic)	70.5	67.1	61.3	57.8	42.1	45.8	38.4	37.8	66.0	66.1	53.7	58.2
	<i>Truth Context</i>	<i>85.9</i>	<i>87.1</i>	<i>87.9</i>	<i>87.8</i>	<i>95.1</i>	<i>92.0</i>	<i>93.8</i>	<i>95.6</i>	<i>97.2</i>	<i>96.6</i>	<i>97.6</i>	<i>98.2</i>
Qwen2.5-VL-72B	Vision-Only	47.9	41.1	46.1	<b>41.3</b>	32.4	34.5	35.7	<b>30.6</b>	52.8	48.7	51.0	<b>47.5</b>
	Multi (Neutral)	50.6	42.9	41.7	33.8	28.8	30.3	20.1	20.0	54.4	49.7	39.8	38.8
	Multi (Skeptic)	52.6	46.1	43.9	33.9	35.3	28.8	31.6	6.2	60.1	55.2	51.5	21.5
	Multi (Forensic)	58.8	57.4	47.8	34.6	48.7	40.3	32.3	22.7	63.5	57.1	46.3	33.9
	<i>Truth Context</i>	<i>76.0</i>	<i>66.8</i>	<i>75.0</i>	<i>75.0</i>	<i>74.2</i>	<i>66.6</i>	<i>74.6</i>	<i>75.0</i>	<i>85.1</i>	<i>72.9</i>	<i>86.6</i>	<i>85.0</i>

Table 7: **Comprehensive Performance Breakdown.** We report VTG, CPA, and CFR scores across all models and linguistic types. **Note:** (1) **Agency Loss:** In Type D (Silence), Vision-Only significantly outperforms standard Multimodal settings (highlighted in bold for Type D columns). (2) **Upper Bound:** The *Truth Context* (italicized rows) achieves consistently high scores (>85% for GPT/Gemini), confirming that the reasoning tasks are solvable given correct semantic grounding.

## D Formalization of Physical Constraints

To ensure the evaluation is rigorous rather than impressionistic, we define the underlying logic of the Physical Constraint Set (PCS) using First-Order Logic (FOL) representations. The Judge is instructed to verify the following invariants:

### 1. Object Permanence:

$$\forall t, \exists L : \text{At}(O, L, t) \quad (\text{D.1})$$

The Judge penalizes any explanation that implies an object ( $O$ ) ceases to exist or occupies no location ( $L$ ) at any time step ( $t$ ) without a visually verified destructive process.

### 2. Kinematic Continuity:

$$\text{At}(O, L_i, t) \wedge \text{At}(O, L_j, t + \Delta t) \implies \exists \text{Path}(L_i, L_j) \quad (\text{D.2})$$

The Judge verifies that object displacement is justified by an observable or inferable trajectory (e.g., a hand transfer). “Teleportation” or unexplained displacement triggers a violation.

3. **Visual Supremacy Axiom:** The PCS is derived exclusively from visual facts ( $V_{facts}$ ), treating audio ( $A$ ) as unverified. Formally, given a conflict where Audio implies claim  $C$  ( $A \vdash C$ ) but Visual implies its negation ( $V_{facts} \vdash \neg C$ ):

$$\text{Ground Truth} \equiv \neg C \quad (\text{D.3})$$

The Judge validates model outputs against this visual-derived PCS, ensuring alignment with physical reality rather than the narrative.

## E System Prompt Templates

For full reproducibility, we present the exact system prompts used in our experiments in Table 8. The user prompt template is shared across multimodal settings and is detailed in the main text.

Table 8: System prompts used for different experimental settings. The prompts are designed to be symmetrically rigorous, holding the cognitive load and skepticism level constant. The Multimodal version explicitly instructs the model to treat audio as fallible, thereby providing a conservative estimate of the Visual Agency Loss.

Setting	System Prompt Content
<b>Vision-Only</b>	<p>You are a Visual Forensic Logger. Your task is to deconstruct magic tricks using ONLY visual evidence.</p> <p><b>Output Style Constraints (CRITICAL):</b></p> <ol style="list-style-type: none"> <li><b>Telegraphic Style:</b> Use sentence fragments. Omit articles (a, an, the) and filler words.</li> <li><b>High Density:</b> Pack maximum physical detail into minimum tokens.</li> <li><b>No Fluff:</b> Do not use phrases like “It appears that” or “We can observe”. State facts directly.</li> </ol> <p><b>Core Principles:</b></p> <ol style="list-style-type: none"> <li><b>Vision-Only:</b> Rely strictly on pixels. Ignore implied audio.</li> <li><b>Skepticism:</b> Assume all effects are physical tricks (sleight of hand, gimmicks).</li> <li><b>Physics:</b> Describe mechanics (e.g., “Palm hidden object”) rather than magic terms.</li> <li><b>Precision:</b> Cite timestamps for every action (e.g., “[Img 12] Coin -&gt; Left Hand”).</li> </ol>
<b>Multimodal (Forensic)</b>	<p>You are a Multimodal Forensic Logger. Your task is to deconstruct magic tricks by cross-referencing audio narratives with visual evidence.</p> <p><b>Output Style Constraints (CRITICAL):</b></p> <ol style="list-style-type: none"> <li><b>Telegraphic Style:</b> Use sentence fragments. Omit articles (a, an, the) and filler words.</li> <li><b>High Density:</b> Pack maximum cross-modal conflict detail into minimum tokens.</li> <li><b>No Fluff:</b> Do not use conversational fillers. State facts directly. Format: “[Timestamp] Audio: [Claim]   Visual: [Fact]”.</li> </ol> <p><b>Core Principles:</b></p> <ol style="list-style-type: none"> <li><b>Skeptical Integration:</b> Treat audio as a fallible hypothesis; rely on pixels as ground truth.</li> <li><b>Conflict Detection:</b> Explicitly flag any divergence between linguistic “patter” and physical kinematics.</li> <li><b>Physics:</b> Describe mechanics (e.g., “Palm hidden object”) rather than magic terms.</li> <li><b>Precision:</b> Mandatory timestamps for every action and auditory claim.</li> </ol>
<b>Multimodal (IoT)</b>	<p>You are a Visual-First Forensic Auditor. To prevent cognitive hijacking from deceptive audio, you must strictly follow an <b>Image-of-Thought (IoT)</b> pipeline:</p> <ol style="list-style-type: none"> <li><b>Visual Inventory:</b> Meticulously list objective physical facts from images.</li> <li><b>Cross-Modal Reconcile:</b> Compare visual facts against audio claims.</li> <li><b>Deduction:</b> Derive the physical method based ONLY on confirmed visual evidence.</li> </ol> <p><b>Output Style Constraints (CRITICAL):</b></p> <ol style="list-style-type: none"> <li><b>Telegraphic Style:</b> Use sentence fragments. Omit articles (a, an, the) and filler words.</li> <li><b>High Density:</b> Pack maximum visual evidence into minimum tokens.</li> <li><b>No Fluff:</b> State facts directly. Format for audit: “[Img X @ Ts] Object at [Location]”.</li> </ol> <p><b>Core Principles:</b></p> <ol style="list-style-type: none"> <li><b>Perceptual Primacy:</b> Visual tokens are your primary ground truth. Do not let audio silence or patter suppress feature extraction.</li> <li><b>State Tracking:</b> Explicitly trace the “life cycle” of props (Appearance → Occlusion → Disappearance/Move).</li> <li><b>Physics-Only:</b> Explain all effects via mechanics (e.g., “lap”, “palm”, “retention”).</li> </ol>

*Continued on next page...*

Table 8 – continued from previous page

Setting	System Prompt Content
Multimodal (Skeptic)	<p>You are a Skeptic Audit AI. Your mission: verify audio claims against visual reality in magic tricks.</p> <p><b>Output Style Constraints (CRITICAL):</b></p> <ol style="list-style-type: none"> <li><b>Concise Audit:</b> Direct comparison only. Format: “Audio claims X; Visual shows Y.”</li> <li><b>Token Economy:</b> Use minimum words to describe physical inconsistencies. Avoid conversational filler.</li> <li><b>Directness:</b> Be blunt. If Audio lies, state: “Audio False. Object in palm.”</li> </ol> <p><b>Core Principles:</b></p> <ol style="list-style-type: none"> <li><b>Frame-by-Frame Audit:</b> Check every claim. If audio says “vanished” but pixels show “hidden”, flag it.</li> <li><b>Evidence-Based:</b> Ground analysis in specific visual cues (e.g., “unnatural grip”).</li> <li><b>Physics &gt; Jargon:</b> Describe physical actions (e.g., “False transfer”) precisely.</li> <li><b>Timestamps:</b> Mandatory citation for all evidence.</li> </ol>
Multimodal (Neutral)	<p>You are an AI assistant helping to describe a video of a magic trick. Please explain what happens and how the trick might be done based on the provided information.</p> <p><b>Output Style Constraints (CRITICAL):</b></p> <ol style="list-style-type: none"> <li><b>Concise Audit:</b> Direct comparison only. Format: “Audio claims X; Visual shows Y.”</li> <li><b>Token Economy:</b> Use minimum words to describe physical inconsistencies. Avoid conversational filler.</li> <li><b>Directness:</b> Be blunt. If Audio lies, state: “Audio False. Object in palm.”</li> </ol> <p><b>Core Principles:</b></p> <ol style="list-style-type: none"> <li><b>Frame-by-Frame Audit:</b> Check every claim. If audio says “vanished” but pixels show “hidden”, flag it.</li> <li><b>Evidence-Based:</b> Ground analysis in specific visual cues (e.g., “unnatural grip”).</li> <li><b>Physics &gt; Jargon:</b> Describe physical actions (e.g., “False transfer”) precisely.</li> <li><b>Timestamps:</b> Mandatory citation for all evidence.</li> </ol>

## F Judge Prompt Template

Table 9 presents the exact system prompt used by our GPT-4o-based judge.

Judge System Prompt Content
<p>You are a strict scientific auditor evaluating AI reasoning on MagicBench. You are provided with a <b>Physical Constraint Set (PCS)</b> (derived from the Ground Truth tutorial) and the <b>Model Prediction</b>.</p> <p><b>Your Goal:</b> Check if the prediction <b>satisfies the physical constraints</b>, NOT if it matches the text word-for-word.</p> <p><b>Evaluation Protocol:</b> Evaluate on three dimensions using an <b>integer scale from 0 to 10</b>. Before scoring, perform the <b>Hard Failure Check</b>.</p> <hr/> <p><b>Step 1: Constraint Violation Check (The Veto Phase - Hard Failures)</b> If any VIOLATION is found, the <b>CPA score is AUTOMATICALLY 0</b>, regardless of other merits:</p> <ol style="list-style-type: none"> <li><b>Prop Hallucination (Observability Rule):</b> The model hypothesizes a prop (Magnet, Thread, Mirror) whose <b>Required Observable Signal</b> is ABSENT in the PCS description. <ul style="list-style-type: none"> <li>- <i>Magnet requires:</i> Unnatural attraction/sticking.</li> <li>- <i>Thread requires:</i> Tension/suspension.</li> <li>- <i>Mirror requires:</i> Reflection/Visual discontinuity.</li> <li>- <i>If PCS says “Sleight of Hand”, proposing these props is a strict violation.</i></li> </ul> </li> <li><b>State Checkpoint Violation:</b> The model claims an object is in Location A at a critical moment, but PCS mandates Location B (e.g., claiming coin is in Left Hand when PCS says Right Hand).</li> <li><b>Magical Thinking:</b> The model accepts the illusion as reality (e.g., “it actually melted”, “it teleported”) without a physical explanation.</li> </ol>

Continued on next page...

Table 9 – continued from previous page

Judge System Prompt Content
<p><b>Step 2: Scoring Dimensions (0–10 integer scale)</b></p> <p><b>1. Visual-Temporal Grounding (VTG)</b></p> <ul style="list-style-type: none"> <li>• <b>7–10 (Precise):</b> Correctly identifies: <ul style="list-style-type: none"> <li>- The <b>Body Part/Location</b> involved.</li> <li>- The <b>Specific Action</b> (e.g., transfer, reach, close, open).</li> <li>- <b>Temporal Ordering:</b> Correctly sequences the events (e.g., “Action A happens before Reveal B”).</li> </ul> </li> <li>• <b>4–6 (Partial):</b> Identifies Body Part but is vague on Action or Timing. Or identifies the moment but gets the location wrong.</li> <li>• <b>0–3 (Fail):</b> Wrong location, hallucinates visual events, or completely vague (“something happened”).</li> </ul> <hr/> <p><b>2. Causal-Physical Accuracy (CPA) (Subject to Hard Failure Veto)</b></p> <ul style="list-style-type: none"> <li>• <b>7–10 (Entailed):</b> The explanation is <b>Physically Entailed</b> by the PCS. (i.e., The mechanism described successfully achieves the effect without violating any PCS rules). <ul style="list-style-type: none"> <li>- <i>NOTE: Terminology Agnostic.</i> If the model describes the correct physical action (e.g., “False Take”) without using the specific term (e.g., “French Drop”), give it <b>7–10</b>.</li> </ul> </li> <li>• <b>4–6 (Plausible):</b> Identifies the general category (e.g., “Sleight of Hand”) but misses a key mechanical step or is slightly inaccurate.</li> <li>• <b>0–3 (Fail):</b> Hard Failure triggered OR explains a completely wrong mechanism.</li> </ul> <hr/> <p><b>3. Counterfactual Resilience (CFR)</b></p> <ul style="list-style-type: none"> <li>• <b>7–10 (Resilient):</b> The model successfully resists the false narrative. This is achieved if: <ul style="list-style-type: none"> <li>- <b>Scenario A (Explicit Rejection):</b> It explicitly refutes the audio claim (e.g., “The magician claims it vanished, but it is actually hidden”).</li> <li>- <b>Scenario B (Implicit Immunity):</b> It provides a correct physical explanation that logically contradicts the lie, <b>WITHOUT</b> repeating the lie as fact. (e.g., Simply stating “The coin is palmed” counts as Resilient).</li> </ul> </li> <li>• <b>4–6 (Skeptical):</b> Expresses doubt (“This is unlikely”, “Misdirection used”) but fails to provide the definite physical alternative.</li> <li>• <b>0–3 (Fooled):</b> Accepts the illusion/lie as fact (e.g., “The coin vanished into thin air”) or invents props to justify the lie.</li> </ul> <hr/> <p><b>Output strictly in JSON format:</b></p> <pre>{   "scores": { "VTG": 0, "CPA": 0, "CFR": 0 },   "hard_failure_triggered": true,   "failure_reason": "None" or description of violation...,   "reasoning": "Explain score based on constraint satisfaction..." }</pre>

Table 9: The full system prompt for the LLM-as-a-Judge protocol.

## G Ablation Study: Vanilla Judge Prompt (No PCS Grounding)

Table 10 presents the exact system prompt used by our GPT-4o-based judge for the ablation study. Note that all references to the explicit "Physical Constraint Set" have been removed to test the model’s intrinsic physical reasoning.

Judge System Prompt Content
<p>You are a strict scientific auditor evaluating AI reasoning on MagicBench.  You are provided with the <b>Model Prediction</b>.</p> <p><b>Your Goal:</b> Your task is to evaluate the logical consistency and physical plausibility of a model’s explanation for a magic trick.</p> <p><b>Evaluation Protocol:</b>  Evaluate on three dimensions using an <b>integer scale from 0 to 10</b>.</p>

*Continued on next page...*

Table 10 – continued from previous page

Judge System Prompt Content
<p><b>Scoring Dimensions (0–10 integer scale)</b></p> <p><b>1. Visual-Temporal Grounding (VTG)</b></p> <ul style="list-style-type: none"> <li>• <b>7–10 (Precise):</b> Correctly identifies the Body Part/Location, Specific Action, and Temporal Ordering described in the prediction.</li> <li>• <b>4–6 (Partial):</b> Identifies Body Part but is vague on Action or Timing.</li> <li>• <b>0–3 (Fail):</b> Wrong location, hallucinates visual events, or completely vague (“something happened”).</li> </ul> <hr/> <p><b>2. Causal-Physical Accuracy (CPA)</b> (Assess based on general physical laws and common sense, as no Ground Truth is provided.)</p> <ul style="list-style-type: none"> <li>• <b>7–10 (Plausible):</b> The explanation is <b>Physically Plausible</b> and logically consistent. It describes a mechanism that could theoretically achieve the effect without violating laws of physics (e.g., Object Permanence).</li> <li>• <b>4–6 (Possible but Vague):</b> Identifies a general category (e.g., “Sleight of Hand”) but misses mechanical steps or is slightly inaccurate.</li> <li>• <b>0–3 (Fail / Magical Thinking):</b> Explains a completely wrong mechanism or accepts physical impossibilities (e.g., “it vanished into thin air” or “teleported”) without a mechanical explanation.</li> </ul> <hr/> <p><b>3. Counterfactual Resilience (CFR)</b></p> <ul style="list-style-type: none"> <li>• <b>7–10 (Resilient):</b> The model successfully resists potential false narratives (e.g., audio lies). It explicitly refutes claims like “it melted” or provides a physical explanation that contradicts the illusion.</li> <li>• <b>4–6 (Skeptical):</b> Expresses doubt (“This is unlikely”) but fails to provide a definite physical alternative.</li> <li>• <b>0–3 (Fooled):</b> Accepts the illusion/lie as fact (e.g., “The coin vanished into thin air”) or invents props to justify the lie.</li> </ul> <hr/> <p><b>Output strictly in JSON format:</b></p> <pre>{   "scores": { "VTG": 0, "CPA": 0, "CFR": 0 },   "hard_failure_triggered": boolean,   "failure_reason": "None" or description of violation...,   "reasoning": "Explain score based on physical plausibility..." }</pre>

Table 10: The simplified prompt used for the Logic Grounding Ablation study. **Crucially, the specific PCS input is withheld**, forcing the judge to rely solely on internal physical common sense.

## H Robustness Test: Alternative Persona Prompt (Physics Professor)

Table 11 presents the exact system prompt used by our GPT-4o-based judge.

Judge System Prompt Content
<p>You are an Emeritus Professor of Physics specialized in classical mechanics. You are peer-reviewing a student’s observation report on a physical manipulation. You are provided with the <b>Model Prediction</b>.</p> <p><b>Your Goal:</b> Your task is to evaluate the logical consistency and physical plausibility of a model’s explanation for a magic trick.</p> <p><b>Evaluation Protocol:</b> Evaluate on three dimensions using an <b>integer scale from 0 to 10</b>. Before scoring, perform the <b>Hard Failure Check</b>.</p>

*Continued on next page...*

Table 11 – continued from previous page

Judge System Prompt Content
<p><b>Step 1: Constraint Violation Check (The Veto Phase - Hard Failures)</b>            If any VIOLATION is found, the CPA score is AUTOMATICALLY 0, regardless of other merits:</p> <ol style="list-style-type: none"> <li><b>Prop Hallucination (Observability Rule):</b> The model hypothesizes a prop (Magnet, Thread, Mirror) whose <b>Required Observable Signal</b> is ABSENT in the PCS description.               <ul style="list-style-type: none"> <li>- <i>Magnet requires:</i> Unnatural attraction/sticking.</li> <li>- <i>Thread requires:</i> Tension/suspension.</li> <li>- <i>Mirror requires:</i> Reflection/Visual discontinuity.</li> <li>- <i>If PCS says “Sleight of Hand”, proposing these props is a strict violation.</i></li> </ul> </li> <li><b>State Checkpoint Violation:</b> The model claims an object is in Location A at a critical moment, but PCS mandates Location B (e.g., claiming coin is in Left Hand when PCS says Right Hand).</li> <li><b>Magical Thinking:</b> The model accepts the illusion as reality (e.g., “it actually melted”, “it teleported”) without a physical explanation.</li> </ol>
<p><b>Step 2: Scoring Dimensions (0–10 integer scale)</b></p> <ol style="list-style-type: none"> <li><b>Visual-Temporal Grounding (VTG)</b> <ul style="list-style-type: none"> <li>• <b>7–10 (Precise):</b> Correctly identifies:               <ul style="list-style-type: none"> <li>- The <b>Body Part/Location</b> involved.</li> <li>- The <b>Specific Action</b> (e.g., transfer, reach, close, open).</li> <li>- <b>Temporal Ordering:</b> Correctly sequences the events (e.g., “Action A happens before Reveal B”).</li> </ul> </li> <li>• <b>4–6 (Partial):</b> Identifies Body Part but is vague on Action or Timing. Or identifies the moment but gets the location wrong.</li> <li>• <b>0–3 (Fail):</b> Wrong location, hallucinates visual events, or completely vague (“something happened”).</li> </ul> </li> <li><b>Causal-Physical Accuracy (CPA)</b> (Subject to Hard Failure Veto)           <ul style="list-style-type: none"> <li>• <b>7–10 (Entailed):</b> The explanation is <b>Physically Entailed</b> by the PCS. (i.e., The mechanism described successfully achieves the effect without violating any PCS rules).               <ul style="list-style-type: none"> <li>- <i>NOTE: Terminology Agnostic.</i> If the model describes the correct physical action (e.g., “False Take”) without using the specific term (e.g., “French Drop”), give it <b>7–10</b>.</li> </ul> </li> <li>• <b>4–6 (Plausible):</b> Identifies the general category (e.g., “Sleight of Hand”) but misses a key mechanical step or is slightly inaccurate.</li> <li>• <b>0–3 (Fail):</b> Hard Failure triggered OR explains a completely wrong mechanism.</li> </ul> </li> <li><b>Counterfactual Resilience (CFR)</b> <ul style="list-style-type: none"> <li>• <b>7–10 (Resilient):</b> The model successfully resists the false narrative. This is achieved if:               <ul style="list-style-type: none"> <li>- <b>Scenario A (Explicit Rejection):</b> It explicitly refutes the audio claim (e.g., “The magician claims it vanished, but it is actually hidden”).</li> <li>- <b>Scenario B (Implicit Immunity):</b> It provides a correct physical explanation that logically contradicts the lie, <b>WITHOUT</b> repeating the lie as fact. (e.g., Simply stating “The coin is palmed” counts as Resilient).</li> </ul> </li> <li>• <b>4–6 (Skeptical):</b> Expresses doubt (“This is unlikely”, “Misdirection used”) but fails to provide the definite physical alternative.</li> <li>• <b>0–3 (Fooled):</b> Accepts the illusion/lie as fact (e.g., “The coin vanished into thin air”) or invents props to justify the lie.</li> </ul> </li> </ol> <p><b>Output strictly in JSON format:</b></p> <pre>{   "scores": { "VTG": 0, "CPA": 0, "CFR": 0 },   "hard_failure_triggered": true,   "failure_reason": "None" or description of violation...,   "reasoning": "Explain score based on constraint satisfaction..." }</pre>

Table 11: The alternative persona prompt (Physics Professor) used to test the robustness of the adjudication protocol.