

When Seeing Overrides Knowing: Disentangling Knowledge Conflicts in Vision-Language Models

Francesco Ortu^{1,2} Zhijing Jin^{3,4,5} Diego Doimo^{2†} Alberto Cazzaniga^{2†}

¹University of Trieste ²AREA Science Park

³MPI ⁴University of Toronto ⁵Vector Institute

Abstract

Vision-language models (VLMs) increasingly combine visual and textual information to perform complex tasks. However, conflicts between their internal knowledge and external visual input can lead to hallucinations and unreliable predictions. In this work, we investigate the mechanisms that VLMs use to resolve cross-modal conflicts by introducing WHOOPS-AHA!, a dataset of multimodal counterfactual queries that deliberately contradict internal common-sense knowledge. Through logit inspection, we identify a small set of attention heads that mediate this conflict. By intervening in these heads, we can steer the model towards its internal parametric knowledge or the visual information. Our results show that attention patterns on these heads effectively locate image regions that influence visual overrides, providing a more precise attribution compared to gradient-based methods.

1 Introduction

Vision-language models (VLMs) (Alayrac et al., 2022; Li et al., 2022; Liu et al., 2023; Team, 2024; Deitke et al., 2024) have shown remarkable versatility in various multimodal tasks, from image understanding to image generation. They draw on their ability to combine two key sources of information: a rich world knowledge acquired during pre-training, and contextual cues provided in the input prompts. However, these two sources can sometimes contradict each other, for example when the pretraining knowledge becomes outdated (Lazari-dou et al., 2021; Luu et al., 2022) or when prompts include intentionally misleading visual information (Liu et al., 2024d). Such conflicts often lead to hallucinations in model responses (Cui et al., 2023;

Correspondence: {francesco.ortu, diego.doimo, alberto.cazzaniga}@areasciencepark.it

Code and Dataset: [francescortu/Seeing-Knowing](https://github.com/francescortu/Seeing-Knowing)

[†] Equal supervision.

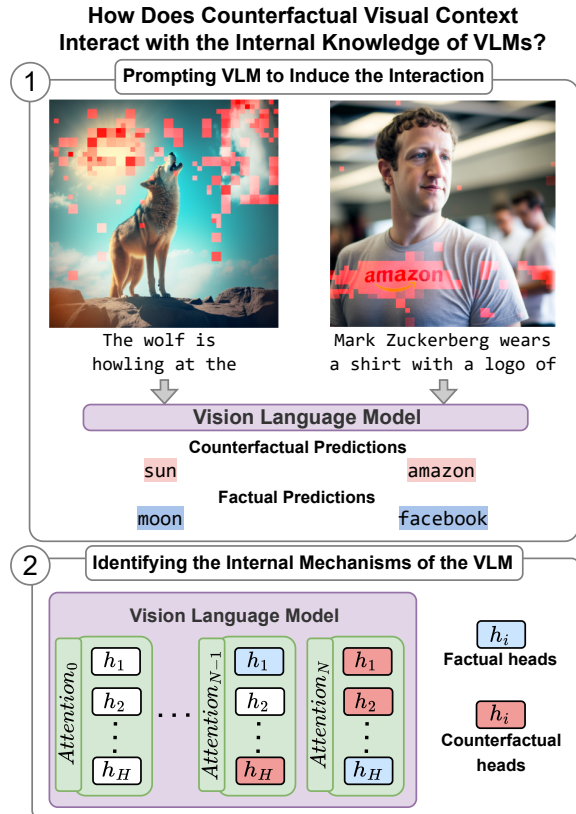


Figure 1: **Overview of our approach.** (Top) We construct prompts that induce a conflict between a vision-language model’s internal factual knowledge and counterfactual visual context. (Bottom) We then analyze which components in the model mediate this tension, identifying attention heads and visual patches that favor factual or visually grounded predictions.

Liu et al., 2024a; Guan et al., 2024), yet the internal mechanisms responsible for these errors remain poorly understood (Xu et al., 2024).

In this work, we analyze how VLMs resolve conflicts between visual input and internal knowledge by framing the problem through counterfactual image-text pairs. We prompt the VLMs with images depicting unusual or absurd scenes taken from the WHOOPS! dataset (Guetta et al., 2023), followed by a sentence encouraging a typical knowledge-based continuation. As shown in Fig. 1, each in-

put prompt is associated with a counterfactual pair of completions. For instance, the model may be shown an image of a wolf howling at the sun, a scene that contradicts commonsense knowledge, and asked to complete the prompt accordingly (see top-left panel). We construct the dataset such that VLMs generate commonsense responses when prompted with text alone, but override them in the presence of an image to align with the visual context, even when it contradicts their internal knowledge. Building on the approach of [Ortu et al. \(2024\)](#), we identify which internal components of the model contribute the most to factual (inner knowledge) versus counterfactual (image context) predictions. We find that a small subset of attention heads mediates this competition, and targeted interventions on these heads can reliably alter the model’s outputs. We also show that these heads are more effective than gradient-based methods in identifying the most important parts of an image to resolve multimodal conflicts in VLMs. In summary, our contributions are as follows:

1. We construct WHOOPS-AHA!, a dataset that combines images containing counterfactual scene elements and commonsense textual queries, designed to analyze conflicts between visual context and internal knowledge (Section 5.1);
2. We identify the attention heads that promote factual and counterfactual responses, ranking their importance with logit attribution (Section 5.2);
3. By reweighting these heads, we show that we can control the tendency of the model to rely on the visual evidence or its internal knowledge and vice versa (Section 5.3);
4. We demonstrate that direct attention attribution from conflict-resolution heads provides more accurate identification of counterfactual image regions than traditional gradient-based attribution methods (Section 5.4).

2 Related Work

Knowledge conflicts in VLMs. VLMs frequently encounter situations where visual input contradicts their internal parametric knowledge, yet the mechanisms governing conflict resolution remain poorly understood ([Xu et al., 2024](#)). Early work on multimodal conflicts focused primarily

on behavioral evaluation through benchmark construction. [Han et al. \(2024\)](#) introduced datasets that probe contextual knowledge conflicts with deceptive visual elements, while [Liu et al. \(2024c\)](#) developed ConflictVis to evaluate conflicts between visual input and parametric knowledge. [Le et al. \(2023\)](#) created COCO-Counterfactuals using minimally edited counterfactual image pairs to study model behavior under visual contradictions. More recently, [Zhu et al. \(2024\)](#) formally characterize cross-modality parametric knowledge conflict in VLMs, showing that conflict rates remain persistently high regardless of model size. However, these studies limit their analysis to evaluating model outputs and prompt structures without investigating the internal mechanisms by which models resolve conflicts.

Mechanistic interpretability in VLMs. Mechanistic interpretability, which seeks to reverse engineer deep neural networks, has made significant strides in text-only models ([Elhage et al., 2021](#); [Geva et al., 2023, 2021](#); [Hanna et al., 2023](#), *inter alia*). Recently, attention has shifted to VLMs. Early work adapted tools from the language setting to multimodal architectures: [Schwettmann et al. \(2023\)](#) identified MLP neurons that convert visual representations into language concepts, while [Palit et al. \(2023\)](#) applied causal tracing to BLIP for visual question answering. Subsequent work extended these methods to generative VLMs: [Neo et al. \(2025\)](#) and [Yu and Ananiadou \(2024\)](#) analyze how LLaVA processes visual information and VQA mechanisms, respectively, [Basu et al. \(2024\)](#) examines knowledge retrieval, and [Jiang et al. \(2025\)](#) use the logit lens on image representations to detect and edit out hallucinations. At the attention-head level, [Basile et al. \(2025\)](#) shows that a small subset of heads can be ranked by relevance to semantic or visual concepts and edited to suppress or enhance targeted concepts, [Yang et al. \(2025b\)](#) identify “hallucination heads” whose attention patterns mirror the base LLM, and [Nikankin et al. \(2025\)](#) use circuit discovery to show that visual and textual tasks recruit largely disjoint computational subgraphs. Despite these advances, the mechanistic investigation of how VLMs resolve conflicts between modalities remains underexplored.

Internal dynamics of multimodal conflicts. Although interest in VLM interpretability is growing, mechanistic studies of how these models process conflicting information remain limited. In the con-

text of LLMs, research has focused on understanding how models resolve conflicts between contextual and internal knowledge (Ortu et al., 2024; Yu et al., 2023; Jin et al., 2024). Recent work has begun exploring internal mechanisms in VLMs: Golovanevsky et al. (2025b) study attention heads in LLaVA and BLIP through semantically corrupted image pairs, Hua et al. (2025) analyze how modality preference under explicit caption-image conflict is reflected in internal representations and modulated by specific heads, and Golovanevsky et al. (2025a) use steering vectors to control the competition between visual input and parametric knowledge on simple visual attributes. Unlike these, we focus on implicit commonsense conflicts and verify that the identified heads are specifically recruited under conflict rather than for general visual processing.

3 Dataset

3.1 Requirements for Mechanistic Analysis of Multimodal Conflicts

Mechanistic interpretability of VLMs requires datasets that enable precise analysis of internal information flow. To support this goal, we identified four key requirements for a suitable dataset:

- **Controlled conflict induction:** Conflicts between visual input and internal knowledge must be systematically induced and verifiable, enabling causal analysis.
- **Token-level precision:** The dataset should allow token-level inspection and interventions, with prompts designed to elicit specific, predictable continuations.
- **Commonsense knowledge grounding:** Scenarios must rely on the model’s internal parametric knowledge, providing strong, consistent priors that can be challenged by visual input. Consistent with the type of commonsense violations studied in WHOOPS! (Guetta et al., 2023), we treat commonsense as robust parametric knowledge encoded in the model’s weights, which serves as a uniform internal prior for analysis.
- **Topical generality:** To test broad knowledge and contextual understanding, the dataset should cover a wide range of topics rather than narrow or highly specific domains.

To meet these requirements, we construct WHOOPS-AHA!, a dataset specifically designed to support mechanistic interpretability techniques for VLMs. To the best of our knowledge, no existing dataset combines these characteristics, making WHOOPS-AHA! a necessary resource for studying controlled knowledge conflicts in multimodal models. Although designed for our experiments, it may also benefit the broader community interested in mechanistic analysis of multimodal conflicts.

3.2 Dataset Construction

WHOOPS-AHA! addresses these requirements by building on the WHOOPS! collection (Guetta et al., 2023), which features 500 visually implausible, semantically rich scenes annotated with textual descriptions and explanations of their underlying anomalies. Each example in WHOOPS-AHA! consists of (i) a counterfactual image depicting an unusual scene, (ii) a sentence referring to the image, and (iii) two sets of plausible continuations: (S_{fact}) reflecting common sense knowledge, and (S_{cofa}) consistent with the counterfactual scene represented in the image. To align with previous work (Ortu et al., 2024), we refer to predictions consistent with internal commonsense (S_{fact}) as *factual*, and those driven by the contradictory visual evidence (S_{cofa}) as *counterfactual*.

Construction pipeline. For each image in WHOOPS!, we use GPT-4o to generate a sentence that references the anomaly, while remaining consistent with commonsense (factual) completion without visual input. GPT-4o is also prompted to produce a set of plausible factual tokens S_{fact} and visually-grounded counterfactual continuations S_{cofa} . For instance, given an image representing a wolf howling at the sun (see Fig. 1), the sentence proposed by GPT-4o is “The wolf is howling at the”, $S_{\text{fact}} = \{\text{“moon”}, \text{“night”}, \dots\}$ $S_{\text{cofa}} = \{\text{“sun”}, \text{“daylight”}, \text{“morning”}, \dots\}$. Full prompt details are provided in Appendix J.

Quality control and validation. To ensure dataset quality, we implemented an LLM-as-a-judge approach (Zheng et al., 2023), using GPT-4.1 (OpenAI, 2025) and Gemini-2.5-Flash (Comanici et al., 2025). Models evaluated each completion for grammatical correctness (1–3 scale) and for alignment with common knowledge or visual anomalies (1–5 scale). Across the dataset, the average grammatical score was 2.94 ± 0.25 for completions of inner knowledge and 2.93 ± 0.28 for completions

aligned with the image. Alignment with knowledge or visual anomalies received a mean score of 4.43 ± 0.97 and 4.69 ± 0.92 , respectively.

To validate this setup, we compared LLM ratings with those of 2 human evaluators on a 20-item subset. Full details, including prompts, scoring instructions, and agreement results, are provided in Appendix B.

4 Background and Methods

4.1 Model Architectures

A VLM encodes image-text tokens with a visual encoder and text embeddings, propagating the resulting residual stream through layers with attention and MLP blocks. The final output is projected to the vocabulary space. We focus our analysis on the residual stream, attention, MLP blocks, and individual attention heads. We focus on two models: LLaVA-NeXT-7B (Liu et al., 2024b) and Gemma3-12B (Kamath et al., 2025). LLaVA-NeXT has 32 layers with 32 attention heads per layer, while Gemma3 has 48 layers with 16 attention heads per layer. Both models use a visual encoder to process image features, but generate only textual output.

4.2 Analytical Tools

Logit inspection. To identify the internal components of VLMs responsible for the competition between inner knowledge and conflicting visual context, we apply the *Logit Lens* technique (Nostalgebraist, 2020), which projects intermediate hidden representations into the vocabulary space. This approach has been used in previous work to analyze token-level information flow (Nanda et al., 2023; Halawi et al., 2023; Yu et al., 2023; Ortu et al., 2024) in LLMs. In our setting, we apply the Logit Lens to the last token of the input and extract the logits corresponding to the tokens in S_{fact} and S_{cofa} in the output of the MLP, Attention block, and across all attention heads, to identify components that favor one mechanism over the other.

Targeted intervention on attention heads. To test the causal role of specific attention heads in promoting predictions aligned with either factual inner knowledge or counterfactual visual context, we intervene on their attention patterns during inference. We define two groups of heads based on Logit Inspection: factual heads ($\mathcal{H}_{\text{fact}}$), which favor predictions based on inner knowledge, and counterfactual heads ($\mathcal{H}_{\text{cofa}}$), which favor visually grounded alternatives. We apply a multiplicative intervention to

their attention weights at the final token position (i.e., the last row of the attention matrix), after the softmax operation. Let $\mathbf{A}_{\text{last}}^{hl} = [\mathbf{A}_{\text{last,img}}^{hl}, \mathbf{A}_{\text{last,text}}^{hl}]$ denote the last row of the attention weights for head h at layer l , divided between image and text tokens. The intervention is defined as

$$\mathbf{A}_{\text{last,img}}^{hl} \leftarrow (1 - \lambda) \cdot \mathbf{A}_{\text{last,img}}^{hl} \quad (1)$$

if $(h, l) \in \mathcal{H}_{\text{cofa}}$, and

$$\mathbf{A}_{\text{last,text}}^{hl} \leftarrow (1 + \lambda) \cdot \mathbf{A}_{\text{last,text}}^{hl} \quad (2)$$

if $(h, l) \in \mathcal{H}_{\text{fact}}$.

This targeted and bidirectional intervention alters the flow of information in a controlled way, allowing us to test whether modulating the influence of these heads changes the model predictions toward the factual or counterfactual outcome.

Identification of conflict-inducing visual tokens.

To isolate the visual tokens responsible for introducing counterfactual information that conflicts with the inner knowledge of the model, we apply two methods. Both are based on a threshold parameter $\tau \in [0, 1]$, which controls the sensitivity of token selection.

1. **Most-Attended Visual Tokens:** Given a set of attention heads, we select the visual tokens that receive at least τ times the maximum attention weight within each head. We then take the union of these tokens across all heads.
2. **Gradient-Based Token Importance:** We compute the gradient of the logit associated with a target token (e.g., from S_{fact} or S_{cofa}) with respect to the input visual token embeddings. Visual tokens whose gradient magnitudes exceed τ times the maximum are selected as influential.

By varying τ , we control how many image patches are selected, from none when τ is 1, to all when τ is 0. This allows us to ablate different image portions and analyze how they affect the model predictions.

Quantifying visual attribution. To evaluate how precisely attribution methods locate counterfactual visual elements, we measure their overlap with ground-truth object regions. We segment the counterfactual object (e.g., the “Amazon” in the t-shirt) and compute an Attribution Ratio: the average attribution value (attention weight or gradient magnitude) within the segmented mask divided by the

average attribution across background patches only. A ratio greater than 1 indicates that the method assigns higher attribution intensity to the counterfactual object than to the background.

5 Experimental Results

5.1 Inducing the Conflict between Inner Knowledge and Visual Context

Given a model, we select t_{fact} as the highest probability token in S_{fact} using text-only prompts, and t_{cofa} as the highest probability token from S_{cofa} using multimodal input. Selecting from these sets ensures that we capture the completions most aligned with the model’s internal knowledge (text-only) or most influenced by visual information (multimodal), allowing us to reliably study the interaction and potential conflicts between the two sources of information. For example, “The wolf is howling at the” yields $t_{\text{fact}} = \text{“moon”}$ (with probabilities of 78% and 100% in LLaVA-NeXT and Gemma3, respectively) in text-only mode, but shift to $t_{\text{fact}} = \text{“sun”}$ (26% LLaVA-NeXT, 44% Gemma3) when the image is included, while the probability of moon drops to 17% and 0.02%. After filtering ambiguous cases where counterfactual tokens dominate in text-only scenarios, we retain 436 examples for LLaVA-NeXT and 432 for Gemma3. The systematic shift from factual to counterfactual predictions (factual accuracy drops to 27% and 24%, respectively) confirms that visual input successfully overrides internal knowledge. This setup ensures that the image introduces a counterfactual signal that conflicts with the model’s inner knowledge, allowing us to analyze how visual input alters the model’s prediction compared to its default behavior based on factual knowledge alone.

5.2 The Tension Between Inner Knowledge and Visual Context is Localized

Building on the controlled knowledge conflict, we apply Logit Lens to identify which model components mediate the competition between t_{fact} and t_{cofa} . For attention and MLP blocks, we report factual preference strength—the deviation from random baseline (0.5) in factual accuracy—where positive values indicate bias toward internal knowledge and negative values toward visual context. For individual attention heads, we report raw factual accuracy (the fraction of examples where factual logits exceed counterfactual logits) to identify heads with strong directional preferences.

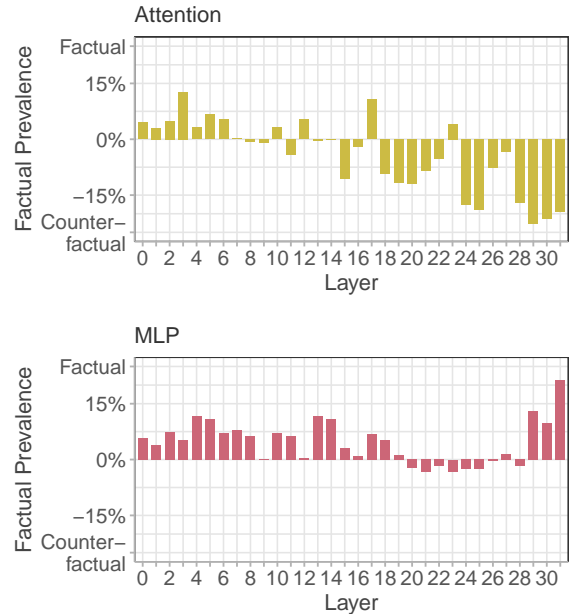


Figure 2: **Factual prevalence in attention and MLP blocks.** Factual prevalence of LLaVA-NeXT shows whether each block favors predictions aligned with factual knowledge (positive) or counterfactual visual context (negative). The results reveal a functional distinction: attention blocks tend to support counterfactual information (**top**), whereas MLP blocks frequently promote the model’s internal knowledge (**bottom**).

Functional separation between attention and MLP layers. We first compare attention and MLP contributions to predicting t_{fact} and t_{cofa} (Fig. 2 for LLaVA-NeXT; see Appendix D for Gemma3). Attention blocks exhibit a stronger tendency to favor the counterfactual visual context, whereas MLP blocks are more aligned with the internal factual knowledge. In particular, the influence of attention blocks increases from the middle layers (around layer 15), peaking in the final four layers. MLP blocks similarly show their strongest alignment to factual knowledge in the upper layers, with a peak at the final layer, consistent with prior findings on upper-layer MLPs retrieving factual knowledge (Geva et al., 2021; Meng et al., 2022; Dai et al., 2022).

Localization of the modality conflict to individual attention heads. We next examine the role of individual attention heads. Figure 3-left shows the tendency for each attention head to promote or suppress the factual token in LLaVA-NeXT (see Fig. 9 for Gemma3). The distribution shows that only a small subset of heads exhibit a strong, consistent alignment with t_{fact} or t_{cofa} . Moreover, con-

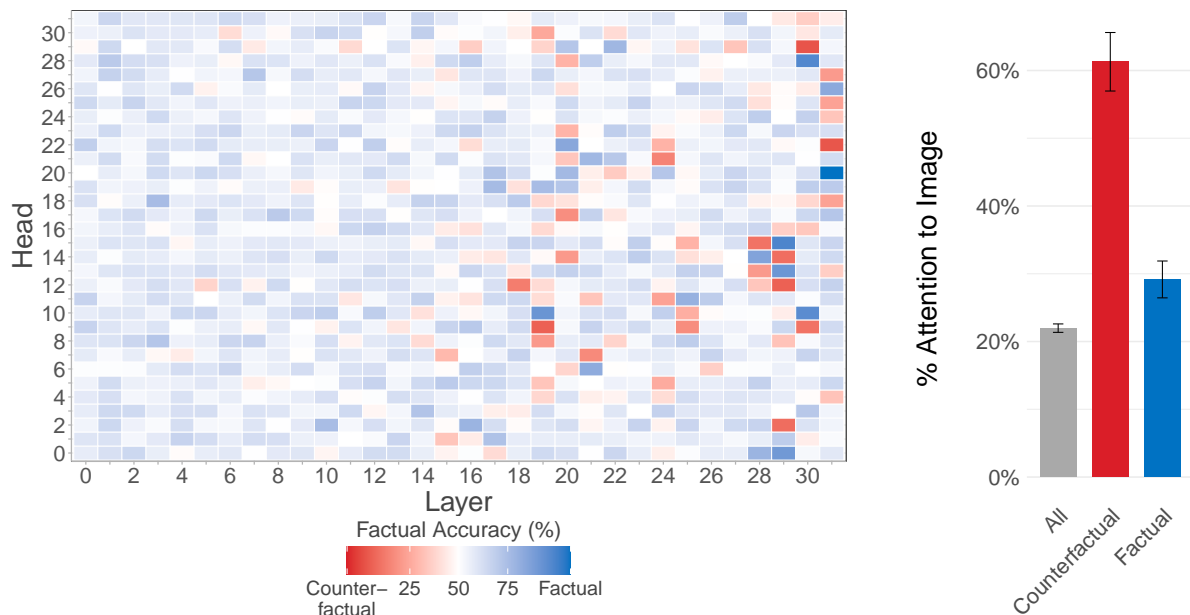


Figure 3: **Contribution of attention heads to factual and counterfactual predictions.** (Left) Factual accuracy of individual attention heads in LLaVA-NeXT, based on Logit Lens projections at the final token position. Blue indicates heads that tend to favor the factual token (reflecting inner knowledge), while red indicates heads that favor the counterfactual token (introduced by the visual context). (Right) Mean attention to image tokens at the final generation step for heads in each group. Each group contains 20 attention heads. Counterfactual heads attend significantly more to the image (60%) than factual heads (28%) or the model-wide average (22%), indicating that visual information is directly propagated to the output and plays a key role in counterfactual predictions.

sistent with the results at the block level, these factual and counterfactual heads are concentrated in the final layers of the model, indicating that the conflict between inner knowledge and visual context is resolved late in the forward pass. In the analyses that follow, we focus on the 20 attention heads that most strongly promote the factual and counterfactual tokens. We chose 20 heads as this provides an optimal balance: it maximizes factual accuracy while minimizing potential disruptions to model stability that could arise from intervening in too many heads (see Appendix E). On average, the factual heads favor the t_{fact} 85% of the time, and the counterfactual ones t_{cofa} 15% of the time, indicating strong alignment with their respective information sources.

Factual and counterfactual heads exhibit distinct visual attention patterns. We then investigate whether heads associated with the factual mechanism or the counterfactual visual context exhibit distinct attention patterns – specifically, whether they attend to different token modalities (image or text). Since the counterfactual information is introduced through the image, a natural hypothesis is that counterfactual heads attend more

strongly to visual tokens, while factual heads rely more on textual content. To test this hypothesis, for each group of heads, we sum the attention weights assigned to visual tokens in the last row of each head and average across the dataset. Figure 3-right reports the average amount of attention to the image for the two groups of heads. Heads favoring the counterfactual token t_{cofa} attend to image tokens significantly more (61%) than those aligned with inner knowledge (29%) or the model-wide average (22%). Although the counterfactual signal originates in the image, it is not a priori clear that this information is transmitted directly to the final token. The model could, in principle, diffuse or encode this signal in different positions across intermediate layers. However, the observed attention patterns suggest that the visual context influences the output token directly in late layers of the model, with limited intermediate processing. These findings are consistent for Gemma3 (see Appendix D).

5.3 Targeted Intervention on Selected Attention Heads Causally Shifts Model Behavior

Having identified attention heads aligned with either factual knowledge or counterfactual visual con-

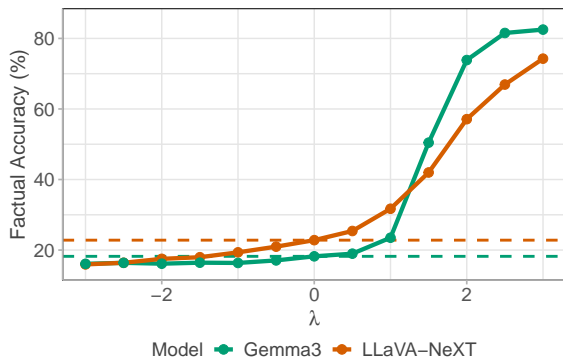


Figure 4: **Intervention on target attention heads.** Change in factual accuracy under different levels of intervention strength (λ). For $\lambda < 0$, we boost the counterfactual heads (on image tokens) and weaken the factual heads (on text tokens); for $\lambda > 0$, we do the opposite. The intervention is applied at the final token position, modifying only the relevant attention values in the last row.

text, we next examine whether these components play a causal role in shaping model predictions. To this end, guided by our earlier observation that counterfactual heads attend more to visual tokens, we apply the targeted bidirectional intervention strategy described in Section 4.2 that selectively adjusts attention values based on head type and token modality, modifying the attention weights to steer the output of the model towards one mechanism or the other. As a control experiment to isolate the effect of targeted interventions, we randomly select 100 attention heads and apply the same intervention for varying λ values. This manipulation does not produce a substantial deviation from the baseline. The complete results for the control experiment are reported in Appendix E.

Figure 4 shows the results of our intervention for LLaVA-NeXT (orange profile) and Gemma3 (green profile). When we increase attention from factual heads and decrease it from counterfactual heads using LLaVA-NeXT, the factual accuracy increases to 74%, indicating a strong shift towards predictions of inner knowledge. Conversely, reversing the intervention reduces the accuracy to 16%, confirming that these heads causally influence whether the model favors factual or counterfactual content. A similar trend can be observed for Gemma3, with an even stronger relative effect driven by its lower baseline factual accuracy of 18% and a peak of 83%. Comprehensive details about the choice of the parameter λ are reported in Appendix F. To verify that $\mathcal{H}_{\text{cofa}}$ are not generic

vision-centric heads, we evaluate them on POPE (Li et al., 2023), a standard binary VQA benchmark designed to probe general visual grounding. Suppressing their image attention leaves accuracy completely unchanged in both models, while removing the image entirely collapses performance to chance. This dissociation confirms that the identified heads are not required for routine visual processing, but are selectively recruited when parametric knowledge and visual evidence conflict. Full results are in Appendix H. To assess whether the identified mechanisms generalize beyond WHOOPS-AHA!, we replicate the same head identification and intervention procedure on an independently constructed visual counterfactual dataset with more photo-realistic images. The results closely mirror our main findings and are reported in Appendix G.

5.4 Counterfactual Predictions Depend on Localized Image Regions

The previous analysis reveals that specific attention heads at the final token position mediate the conflict between contextual information and internal knowledge, with heads aligned with the visual context strongly attending to image tokens, injecting visually grounded information into the generation process. However, two key questions remain open. (i) Is the counterfactual visual signal localized to specific image regions or spread across the input? (ii) Is the visual signal passed directly to the last token position, or is it mediated by successive layers and tokens before reaching the output in the upper layers? To address these, we conduct two analyses: (i) using attention and gradient-based attribution to identify the image patches driving counterfactual predictions, as described in Section 4.2; and (ii) ablating these patches by setting their visual token embeddings to zero and measuring the change in factual accuracy. A control experiment is also performed where an equivalent number of randomly selected patches are ablated.

Patch attribution and ablation effects. The results (Fig. 5) show that ablating patches identified through attention-based attribution leads to a sharp and consistent increase in factual accuracy as more pixels are removed (green profiles). For LLaVA-NeXT, factual accuracy improves markedly with the ablation of just 10–30% of the top-ranked patches and eventually plateaus around 80%. Gradient-based attribution (shown in red) also yields a substantial increase in factual

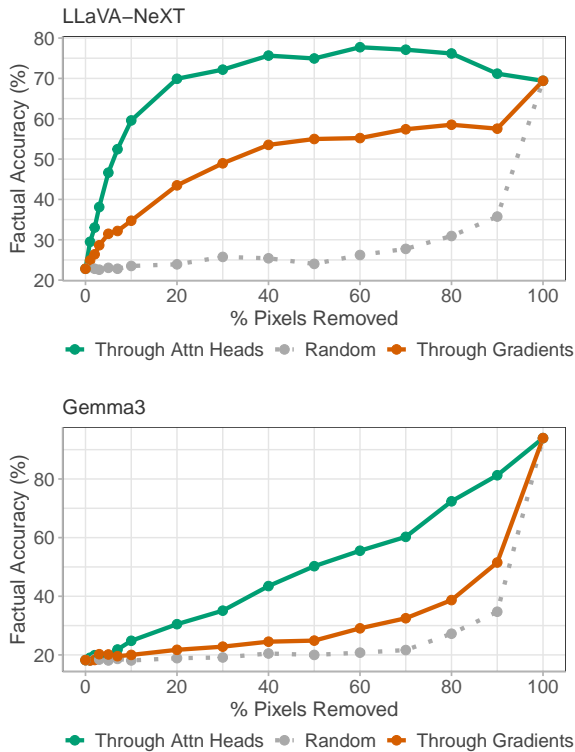


Figure 5: **Ablation of relevant pixels.** The plot shows the effect of ablating different percentages of image pixels in LLaVA-NeXT. The green line corresponds to pixels selected based on the highest attention from counterfactual heads, while the orange line corresponds to pixels with the highest gradient magnitude with respect to the counterfactual token. The gray line shows a random baseline where pixels are removed uniformly at random.

accuracy, but with less pronounced effects, suggesting lower precision in identifying counterfactual-driving regions. In contrast, ablating an equivalent number of randomly selected patches results in only minor fluctuations in accuracy. These findings confirm the causal role of the identified regions and support the hypothesis that counterfactual signals are spatially localized and semantically specific.

Localization of counterfactual visual evidence.

To assess the semantic coherence of the identified visual regions, we qualitatively inspect examples where attribution highlights patches responsible for counterfactual predictions (Fig. 6). In many cases, these patches correspond to intuitive visual elements that directly contradict commonsense expectations, such as implausible objects or substitutions. For example, when the model predicts “rainbow” instead of “black” for a bearskin hat, the highlighted regions focus on the hat’s unrealistic coloring, and when “fruit” replaces “tissue”

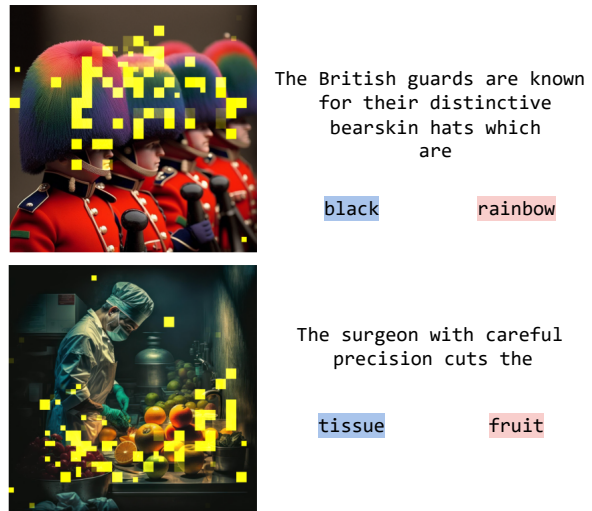


Figure 6: **Visual regions driving counterfactual predictions.** Highlighted image regions, identified through attention-based attribution, show the most influential visual patches for counterfactual predictions. The highlighted areas align with semantically meaningful, visually anomalous content, indicating that counterfactual outputs are grounded in localized image features.

in a surgical scene, attention concentrates on the unexpected oranges on the operating table.

We complement this analysis with a quantitative experiment on 20 randomly selected images, using SAM-3 (Carion et al., 2025) to segment the counterfactual objects. We compare attention from the identified counterfactual heads to both gradient-based attribution and a random-head baseline. Using the *Average Ratio* metric (Section 4.2), we find that for Gemma3 the counterfactual heads strongly concentrate attention on the segmented regions (median 4.41, IQR 1.09–7.85), significantly outperforming both gradients (median 1.74) and random heads (median 0.92). These results quantitatively validate that the heads identified by our method act as precise pointers to the visual evidence driving the counterfactual override (see Appendix I for full details).

6 Discussion

The competition between parametric knowledge and conflicting external evidence, identified in text-only LLMs (Ortu et al., 2024; Jin et al., 2024; Wu et al., 2025; Yu et al., 2023), extends to the multimodal setting: parametric–visual conflicts are resolved by a small, identifiable subset of attention heads in the upper layers of the model. The functional dissociation between attention and MLP blocks, with attention favoring visual context and

MLPs reinforcing parametric priors, aligns with the established view that upper-layer MLPs store factual associations (Geva et al., 2023; Meng et al., 2022) while attention heads integrate contextual signals. Subsequent work has converged on related findings: Golovanevsky et al. (2025a) trace the layer-wise evolution of modality preference via early decoding and propose steering vectors to control it, while Hua et al. (2025) show that specific attention heads modulate modality preference under explicit caption–image conflict. Relative to both, our setting is distinctive in studying implicit commonsense conflicts rather than explicit modality instructions or simple visual attributes, and in establishing that the identified heads are specifically recruited under conflict rather than being generic vision-processing components. Golovanevsky et al. (2025b) similarly identify functionally specialized heads in LLaVA and BLIP, though without focusing on conflict resolution.

The additional finding that conflict-resolution heads encode *where* in the image the contradicting evidence originates suggests that mechanistically identified circuits can serve as a byproduct attribution tool, complementing approaches that use internal representations for visual grounding (Jiang et al., 2025; Phukan et al., 2025) and functionally specialized heads (Yang et al., 2025a; Sarkar et al., 2025; He et al., 2025).

Our analysis is bounded by its focus on late-fusion architectures and the logit lens, which introduces known approximation errors (Belrose et al., 2023); extending this framework to free-form generation, earlier-fusion models, and broader conflict types remains important future work.

7 Conclusion

We introduced WHOOPS-AHA! and a mechanistic pipeline to study how VLMs resolve conflicts between visual evidence and internal parametric knowledge. A small set of upper-layer attention heads mediates this competition with clear functional specialization: counterfactual heads amplify visual evidence while factual heads reinforce parametric priors, and targeted interventions on them causally shift model behavior in both directions. These heads are conflict-specific rather than generically visual, and their attention patterns localize counterfactual image regions more precisely than gradient-based attribution. We hope these findings provide a foundation for building multimodal sys-

tems that more reliably calibrate reliance on vision and knowledge.

Limitations

Methodological limitations. The analysis relies on the Logit Lens technique to project intermediate hidden states into token logits. Although this method has been widely adopted for interpretability, it is known to introduce distortions due to projection from non-final residual states (Belrose et al., 2023), and should be interpreted as an approximate diagnostic rather than a precise decoding proxy. In our setting, we use a representative factual and counterfactual token per example to enable controlled comparisons. Although this simplifies the generative landscape of the model, it offers a practical and interpretable probe of the underlying mechanisms. Future work could explore more model behavior across full generations to complement this approach. Our attribution and intervention methods focus on attention heads and target the final token position. This design isolates interpretable causal signals while remaining tractable, though it does not capture the possible contributions of other components, such as MLP layers or visual encoders. Extending this framework to broader architectural elements is a promising direction.

Scope and generalizability. We focus on late-fusion, LLaVA-style architectures, which are particularly well-suited for controlled image-understanding tasks. These models are among the best open-source architectures for visual understanding, making them ideal for the interpretability methods employed in our study. Our interest is specifically in how visual input interacts with internal knowledge during textual generation, so we chose models that are designed with a focus on image understanding. While early or mid-fusion models also use attention to integrate visual features into the language stream, they may differ significantly in how information is communicated between the modalities (Serra et al., 2024). The point of injection of visual features varies, but the underlying mechanism of cross-modal communication through attention remains consistent across these models. By focusing on late-fusion models, we ensure a more controlled and traceable examination of visual-to-text interactions in widely used open source multimodal models, though this choice limits the generalizability of our findings to models with different fusion strategies.

Ethical Considerations

This work aims to improve our understanding of how VLMs resolve conflicts between internal factual knowledge and contradictory visual context. Our analysis is intended to contribute to foundational research in model interpretability, with the broader goal of developing more transparent and controllable multimodal systems. The techniques presented are diagnostic and exploratory in nature, designed to support responsible development and analysis of multimodal systems. We believe that studying the dynamics of conflicting information sources is essential for anticipating model failure modes, mitigating unintended behaviors, and building more robust AI systems. All models and data are used in accordance with their intended research licenses, and WHOOPS-AHA! is released solely for non-commercial, research purposes under compatible terms. We used AI assistants (e.g., GitHub Copilot) to support code completion during experiment implementation; all generated code was manually reviewed and supervised by the authors.

Acknowledgments

The authors acknowledge the AREA Science Park supercomputing platform ORFEO made available for conducting the research reported in this paper and the technical support of the Laboratory of Data Engineering staff. Francesco Ortu, Diego Doimo and Alberto Cazzaniga were supported by the project “Supporto alla diagnosi di malattie rare tramite l’intelligenza artificiale” CUP: F53C22001770002 and “Valutazione automatica delle immagini diagnostiche tramite l’intelligenza artificiale”, CUP: F53C22001780002. A. C. acknowledges financial support under the National Recovery and Resilience Plan (NRRP), mission 4, component 2, investment 1.1, and call for tender no. 1409 published on 14 September 2022 by the Italian Ministry of University and Research (MUR), funded by the European Union — NextGenerationEU — CUP J53D23015070001.

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, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Miłkołaj Bińkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karén Simonyan. 2022.

[Flamingo: a visual language model for few-shot learning](#). In *Advances in Neural Information Processing Systems*, volume 35, pages 23716–23736. Curran Associates, Inc.

Lorenzo Basile, Valentino Maiorca, Diego Doimo, Francesco Locatello, and Alberto Cazzaniga. 2025. [Head pursuit: Probing attention specialization in multimodal transformers](#). In *Advances in Neural Information Processing Systems (NeurIPS)*.

Samyadeep Basu, Martin Grayson, Cecily Morrison, Bsmira Nushi, Soheil Feizi, and Daniela Massiceti. 2024. [Understanding Information Storage and Transfer in Multi-modal Large Language Models](#). *arXiv preprint*. ArXiv:2406.04236.

Nora Belrose, Zach Furman, Logan Smith, Danny Halawi, Igor Ostrovsky, Lev McKinney, Stella Biderman, and Jacob Steinhardt. 2023. [Eliciting latent predictions from transformers with the tuned lens](#). *CoRR*, abs/2303.08112.

Nicolas Carion, Laura Gustafson, Yuan-Ting Hu, Shoubhik Debnath, Ronghang Hu, Didac Suris, Chaitanya Ryali, Kalyan Vasudev Alwala, Haitham Khedr, Andrew Huang, Jie Lei, Tengyu Ma, Baishan Guo, Arpit Kalla, Markus Marks, Joseph Greer, Meng Wang, Peize Sun, Roman Rädle, Triantafyllos Afouras, Effrosyni Mavroudi, Katherine Xu, Tsung-Han Wu, Yu Zhou, Liliane Momeni, Rishi Hazra, Shuangrui Ding, Sagar Vaze, Francois Porcher, Feng Li, Siyuan Li, Aishwarya Kamath, Ho Kei Cheng, Piotr Dollár, Nikhila Ravi, Kate Saenko, Pengchuan Zhang, and Christoph Feichtenhofer. 2025. [Sam 3: Segment anything with concepts](#). *Preprint*, arXiv:2511.16719.

Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit S. Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, Luke Maris, Sam Petulla, Colin Gaffney, Asaf Aharoni, Nathan Lintz, Tiago Cardal Pais, Henrik Jacobsson, Idan Szepetor, Nan-Jiang Jiang, Krishna Haridasan, Ahmed Omran, Nikunj Saunshi, Dara Bahri, Gaurav Mishra, Eric Chu, Toby Boyd, Brad Hekman, Aaron Parisi, Chaoyi Zhang, Kornrathop Kawintiranon, Tania Bedrax-Weiss, Oliver Wang, Ya Xu, Ollie Purkiss, Uri Mendlovic, Ilai Deutel, Nam Nguyen, Adam Langley, Flip Korn, Lucia Rossazza, Alexandre Ramé, Sagar Waghmare, Helen Miller, Nathan Byrd, Ashrith Sheshan, Raia Hadsell Sangnie Bhardwaj, Pawel Janus, Tero Rissa, Dan Horgan, Sharon Silver, Ayzaan Wahid, Sergey Brin, Yves Raimond, Klemen Kloboves, Cindy Wang, Nitesh Bharadwaj Gundavarapu, Iliia Shumailov, Bo Wang, Mantas Pajarskas, Joe Heyward, Martin Nikolchev, Maciej Kula, Hao Zhou, Zachary Garrett, Sushant Kafle, Sercan Arik, Ankita Goel, Mingyao Yang, Jiho Park, Koji Kojima, Parsa Mahmoudieh, Koray Kavukcuoglu, Grace Chen, Doug Fritz, Anton Bulyenov, Sudeshna Roy, Dimitris Paparas, Hadar Shemtov, Bo-Juen Chen, Robin Strudel, David Reitter, Aurko Roy, Andrey Vlasov, Changwan Ryu, Chas Lechner, Haichuan Yang, Zeldia Mariet, Denis Vnukov, Tim Sohn, Amy Stuart, Wei Liang, Minmin Chen, Praynaa Rawlani, Christy Koh, JD Co-Reyes, Guangda Lai, Praseem Banzal, Dimitrios Vytiniotis, Jieru Mei, and Mu Cai. 2025. [Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities](#). *CoRR*, abs/2507.06261.

Chenhang Cui, Yiyang Zhou, Xinyu Yang, Shirley Wu, Linjun Zhang, James Zou, and Huaxiu Yao. 2023. [Holistic analysis of hallucination in gpt-4v\(ision\): Bias and interference challenges](#). *CoRR*, abs/2311.03287.

- Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. 2022. [Knowledge neurons in pretrained transformers](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 8493–8502. Association for Computational Linguistics.
- Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, Jiasen Lu, Taira Anderson, Erin Bransom, Kiana Ehsani, Huong Ngo, YenSung Chen, Ajay Patel, Mark Yatskar, Chris Callison-Burch, Andrew Head, Rose Hendrix, Favyen Bastani, Eli VanderBilt, Nathan Lambert, Yvonne Chou, Arnavi Chheda, Jenna Sparks, Sam Skjonsberg, Michael Schmitz, Aaron Sarnat, Byron Bischoff, Pete Walsh, Chris Newell, Piper Wolters, Tanmay Gupta, Kuo-Hao Zeng, Jon Borchardt, Dirk Groeneveld, Jen Dumas, Crystal Nam, Sophie Lebrecht, Caitlin Wittlif, Carissa Schoenick, Oscar Michel, Ranjay Krishna, Luca Weihs, Noah A. Smith, Hannaneh Hajishirzi, Ross Girshick, Ali Farhadi, and Aniruddha Kembhavi. 2024. Molmo and pixmo: Open weights and open data for state-of-the-art multimodal models. *arXiv preprint arXiv:2409.17146*.
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. 2021. A mathematical framework for transformer circuits. *Transformer Circuits Thread*. <https://transformer-circuits.pub/2021/framework/index.html>.
- Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. 2023. [Dissecting recall of factual associations in auto-regressive language models](#). *CoRR*, abs/2304.14767.
- Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2021. [Transformer feed-forward layers are key-value memories](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 5484–5495. Association for Computational Linguistics.
- Michal Golovanevsky, William Rudman, Michael Lepori, Amir Bar, Ritambhara Singh, and Carsten Eickhoff. 2025a. [Pixels versus priors: Controlling knowledge priors in vision-language models through visual counterfactuals](#). *CoRR*, abs/2505.17127.
- Michal Golovanevsky, William Rudman, Vedant Palit, Carsten Eickhoff, and Ritambhara Singh. 2025b. [What do VLMs NOTICE? a mechanistic interpretability pipeline for Gaussian-noise-free text-image corruption and evaluation](#). In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 11462–11482, Albuquerque, New Mexico. Association for Computational Linguistics.
- 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](#). *Preprint*, arXiv:2310.14566.
- Nitzan Bitton Guetta, Yonatan Bitton, Jack Hessel, Ludwig Schmidt, Yuval Elovici, Gabriel Stanovsky, and Roy Schwartz. 2023. [Breaking common sense: Whoops! A vision-and-language benchmark of synthetic and compositional images](#). In *IEEE/CVF International Conference on Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023*, pages 2616–2627. IEEE.
- Danny Halawi, Jean-Stanislas Denain, and Jacob Steinhardt. 2023. [Overthinking the truth: Understanding how language models process false demonstrations](#). *CoRR*, abs/2307.09476.
- Tianyang Han, Qing Lian, Rui Pan, Renjie Pi, Jipeng Zhang, Shizhe Diao, Yong Lin, and Tong Zhang. 2024. [The instinctive bias: Spurious images lead to illusion in MLLMs](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 16163–16177, Miami, Florida, USA. Association for Computational Linguistics.
- Michael Hanna, Ollie Liu, and Alexandre Variengien. 2023. [How does GPT-2 compute greater-than?: Interpreting mathematical abilities in a pre-trained language model](#). *CoRR*, abs/2305.00586.
- Jinghan He, Kuan Zhu, Haiyun Guo, Junfeng Fang, Zhenglin Hua, Yuheng Jia, Ming Tang, Tat-Seng Chua, and Jinqiao Wang. 2025. [Cracking the code of hallucination in vlms with vision-aware head divergence](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, pages 3488–3501. Association for Computational Linguistics.
- Tianze Hua, Tian Yun, and Ellie Pavlick. 2025. [How do vision-language models process conflicting information across modalities?](#) *CoRR*, abs/2507.01790.
- Nick Jiang, Anish Kachinthaya, Suzanne Petryk, and Yossi Gandelsman. 2025. [Interpreting and editing vision-language representations to mitigate hallucinations](#). In *The Thirteenth International Conference on Learning Representations*.
- Zhuoran Jin, Pengfei Cao, Hongbang Yuan, Yubo Chen, Jiexin Xu, Huaijun Li, Xiaojian Jiang, Kang Liu, and Jun Zhao. 2024. [Cutting off the head ends the conflict: A mechanism for interpreting and mitigating knowledge conflicts in language models](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 1193–1215, Bangkok, Thailand. Association for Computational Linguistics.
- Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, Louis Rouillard, Thomas Mesnard, Geoffrey Cideron, Jean-Bastien Grill, Sabela Ramos, Edouard Yvinec, Michelle Casbon, Etienne Pot, Ivo Penchev, Gaël Liu, Francesco Visin, Kathleen Kenealy, Lucas Beyer, Xiaohai Zhai, Anton Tsitsulin, Róbert Busa-Fekete, Alex Feng, Noveen Sachdeva, Benjamin Coleman, Yi Gao, Basil Mustafa, Iain Barr, Emilio Parisotto, David Tian, Matan Eyal, Colin Cherry, Jan-Thorsten Peter, Danila Sinopalnikov, Surya Bhupatiraju, Rishabh Agarwal, Mehran Kazemi, Dan Malkin, Ravin Kumar, David Vilar, Idan Brusilovsky, Jiaming Luo, Andreas Steiner, Abe Friesen, Abhanshu Sharma, Abheesht Sharma, Adi Mayrav Gilady, Adrian Goedeckemeyer, Alaa Saade, Alexander Kolesnikov, Alexei Bendebury, Alvin Abdagic, Amit Vadi, András György, André Susano Pinto,

- Anil Das, Ankur Bapna, Antoine Miech, Antoine Yang, Antonia Paterson, Ashish Shenoy, Ayan Chakrabarti, Bilal Piot, Bo Wu, Bobak Shahriari, Bryce Petrin, Charlie Chen, Charline Le Lan, Christopher A. Choquette-Choo, CJ Carey, Cormac Brick, Daniel Deutsch, Danielle Eisenbud, Dee Cattle, Derek Cheng, Dimitris Paparas, Divyashree Shivakumar Sreepathihalli, Doug Reid, Dustin Tran, Dustin Zelle, Eric Noland, Erwin Huizenga, Eugene Kharitonov, Frederick Liu, Gagik Amirkhanyan, Glenn Cameron, Hadi Hashemi, Hanna Klimczak-Plucinska, Harman Singh, Harsh Mehta, Harshal Tushar Lehri, Hussein Hazimeh, Ian Ballantyne, Idan Szpektor, and Ivan Nardini. 2025. [Gemma 3 technical report](#). *CoRR*, abs/2503.19786.
- Angeliki Lazaridou, Adhiguna Kuncoro, Elena Gribovskaya, Devang Agrawal, Adam Liska, Tayfun Terzi, Mai Gimenez, Cyprien de Masson d’Autume, Tomáš Kočiský, Sebastian Ruder, Dani Yogatama, Kris Cao, Susannah Young, and Phil Blunsom. 2021. [Mind the gap: Assessing temporal generalization in neural language models](#). In *Advances in Neural Information Processing Systems*.
- Tiep Le, Vasudev Lal, and Phillip Howard. 2023. [Coco-counterfactuals: Automatically constructed counterfactual examples for image-text pairs](#). 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*.
- Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. [BLIP: Bootstrapping language-image pre-training for unified vision-language understanding and generation](#). In *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 12888–12900. PMLR.
- Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne 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, EMNLP 2023, Singapore, December 6-10, 2023*, pages 292–305. Association for Computational Linguistics.
- Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. 2024a. [Mitigating hallucination in large multi-modal models via robust instruction tuning](#). In *The Twelfth International Conference on Learning Representations*.
- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024b. [Llava-next: Improved reasoning, ocr, and world knowledge](#).
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. [Visual instruction tuning](#). In *Advances in Neural Information Processing Systems*, volume 36, pages 34892–34916. Curran Associates, Inc.
- Xiaoyuan Liu, Wenxuan Wang, Youliang Yuan, Jen tse Huang, Qiuzhi Liu, Pinjia He, and Zhaopeng Tu. 2024c. [Insight over sight? exploring the vision-knowledge conflicts in multimodal llms](#). *Preprint*, arXiv:2410.08145.
- Yi Liu, Gelei Deng, Yuekang Li, Kailong Wang, Zihao Wang, Xiaofeng Wang, Tianwei Zhang, Yepang Liu, Haoyu Wang, Yan Zheng, and Yang Liu. 2024d. [Prompt injection attack against llm-integrated applications](#). *Preprint*, arXiv:2306.05499.
- Kelvin Luu, Daniel Khashabi, Suchin Gururangan, Karishma Mandyam, and Noah A. Smith. 2022. [Time waits for no one! analysis and challenges of temporal misalignment](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5944–5958, Seattle, United States. Association for Computational Linguistics.
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. [Locating and editing factual associations in GPT](#). In *NeurIPS*.
- Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. 2023. [Progress measures for grokking via mechanistic interpretability](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.
- Clement Neo, Luke Ong, Philip Torr, Mor Geva, David Krueger, and Fazl Barez. 2025. [Towards interpreting visual information processing in vision-language models](#). In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net.
- Yaniv Nikankin, Dana Arad, Yossi Gandelsman, and Yonatan Belinkov. 2025. [Same task, different circuits: Disentangling modality-specific mechanisms in VLMs](#). In *Advances in Neural Information Processing Systems*, volume 38.
- Nostalgebraist. 2020. [interpreting gpt: the logit lens](#). Accessed: Nov 2023.
- OpenAI. 2025. Introducing gpt-4.1 in the api. <https://openai.com/index/gpt-4-1/>. Accessed: 2025-09-22.
- OpenRouter, Inc. 2025. [Openrouter: The unified interface for llms](#).
- Francesco Ortu, Zhijing Jin, Diego Doimo, Mrinmaya Sachan, Alberto Cazzaniga, and Bernhard Schölkopf. 2024. [Competition of mechanisms: Tracing how language models handle facts and counterfactuals](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 8420–8436. Association for Computational Linguistics.
- Vedant Palit, Rohan Pandey, Aryaman Arora, and Paul Pu Liang. 2023. [Towards vision-language mechanistic interpretability: A causal tracing tool for BLIP](#). In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops*, pages 2856–2861.
- Anirudh Phukan, Divyansh, Harshit Kumar Morj, Vaishnavi, Apoorv Saxena, and Koustava Goswami. 2025. [Beyond logit lens: Contextual embeddings for robust hallucination detection & grounding in vlms](#). In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2025 - Volume 1: Long Papers, Albuquerque, New Mexico, USA, April 29 - May 4, 2025*, pages 9661–9675. Association for Computational Linguistics.
- Sreetama Sarkar, Yue Che, Alex Gavin, Peter Anthony Beerel, and Souvik Kundu. 2025. [Mitigating hallucinations in vision-language models through image-guided head suppression](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing, EMNLP 2025, Suzhou, China, November 4-9, 2025*, pages 12481–12500. Association for Computational Linguistics.

- Sarah Schwettmann, Neil Chowdhury, Samuel Klein, David Bau, and Antonio Torralba. 2023. Multimodal neurons in pretrained text-only transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops*, pages 2862–2867.
- Alessandro Serra, Francesco Ortu, Emanuele Panizon, Lucrezia Valeriani, Lorenzo Basile, Alessio Ansuini, Diego Doimo, and Alberto Cazzaniga. 2024. [The narrow gate: Localized image-text communication in vision-language models](#). *CoRR*, abs/2412.06646.
- Chameleon Team. 2024. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint arXiv:2405.09818*.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.
- Wenhao Wu, Yizhong Wang, Guangxuan Xiao, Hao Peng, and Yao Fu. 2025. [Retrieval head mechanistically explains long-context factuality](#). In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net.
- Rongwu Xu, Zehan Qi, Zhijiang Guo, Cunxiang Wang, Hongru Wang, Yue Zhang, and Wei Xu. 2024. [Knowledge conflicts for LLMs: A survey](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 8541–8565, Miami, Florida, USA. Association for Computational Linguistics.
- Le Yang, Ziwei Zheng, Boxu Chen, Zhengyu Zhao, Chenhao Lin, and Chao Shen. 2025a. [Nullu: Mitigating object hallucinations in large vision-language models via halluspace projection](#). In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2025, Nashville, TN, USA, June 11-15, 2025*, pages 14635–14645. Computer Vision Foundation / IEEE.
- Tianyun Yang, Ziniu Li, Juan Cao, and Chang Xu. 2025b. [Understanding and mitigating hallucination in large vision-language models via modular attribution and intervention](#). In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net.
- Qinan Yu, Jack Merullo, and Ellie Pavlick. 2023. [Characterizing mechanisms for factual recall in language models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9924–9959, Singapore. Association for Computational Linguistics.
- Zeping Yu and Sophia Ananiadou. 2024. Understanding multimodal LLMs: the mechanistic interpretability of LLaVA in visual question answering. *arXiv preprint arXiv:2411.10950*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. [Judging llm-as-a-judge with mt-bench and chatbot arena](#). 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*.
- Tinghui Zhu, Qin Liu, Fei Wang, Zhengzhong Tu, and Muhao Chen. 2024. [Unraveling cross-modality knowledge conflicts in large vision-language models](#). *CoRR*, abs/2410.03659.

A Reproducibility

We ran the experiments on one NVIDIA H100 GPU, and two GPUs for the gradient-based attribution tests. We use the HuggingFace Transformers library (Wolf et al., 2020) with public implementations of LLaVA-NeXT and Gemma3. The total compute time is 20 GPU hours. The WHOOPS! dataset was released with a CC-By 4.0 license.

B LLM-as-a-Judge: Detailed Validation and Analysis

B.1 Evaluation Setup

We used GPT-4.1 (gpt-4.1-2025-04-14) and Gemini-2.5-Flash (gemini-2.5-flash-image-preview) through OpenRouter (OpenRouter, Inc., 2025) to evaluate each dataset completion along two dimensions:

- Grammatical correctness:
 - 1 (No) All completed sentences contain grammatical errors.
 - 2 (Some do not make sense) Some completed sentences have grammatical errors.
 - 3 (Yes) All completed sentences follow proper grammar rules.
- Knowledge/Anomaly Alignment (1–5 scale)
 - 1 = Poor alignment: Completion ignores or misrepresents common knowledge or visual anomalies
 - 5 = Excellent alignment: Completion clearly reflects correct knowledge or accurately captures anomalies in the image.

B.2 Aggregate Statistics

We report the average results for all dimensions in Table 1. The results show that the judges rated the sentences and the completions as mostly grammatically correct and aligned with the common knowledge or the anomaly in the image.

To measure agreement between the two judges, we calculate the exact agreement rate. We observed a strong ceiling effect in our data, with the distribution of responses being heavily imbalanced. Overall, the results demonstrate that the LLM as a judge setup produces consistent and interpretable ratings across both factual and counterfactual scenarios, with high alignment between models (Gemini vs. GPT-4.1). These findings support the reliability

of LLM-based evaluation for large-scale dataset validation.

B.3 Human Validation

To assess how well the LLM-as-a-judge setup reflects human evaluation, we sampled 20 entries from our dataset and asked 2 human evaluators, both graduate students at a local institution, to complete the same tasks as the LLMs. To measure agreement between human evaluators and LLM judgments, we averaged the ratings of the human evaluators and rounded to the nearest integer for each category, and then measured the exact and adjacent agreement rates.

For grammatical correctness in text-based completions, we observed 60% agreement, and 70% for image-associated completions. For alignment, the agreement was 35% for text-only completions and 65% for image-associated completions. However, when considering adjacent category agreement (i.e., when responses were within one point of each other on the scale), the agreement increased to 85% for text-based completions and 100% for image-associated completions.

B.4 Prompt and Instruction

Prompt Used to Instruct LLM-as-a-Judge with Image

You will evaluate sentence completions that should describe UNUSUAL or ANOMALOUS elements visible in the provided image. Focus on identifying completions that capture strange, unexpected, or contradictory visual elements.

EVALUATION CRITERIA:

1. GRAMMATICAL CORRECTNESS: Are the completed sentences grammatically valid?

Check syntax, word order, agreement, and grammatical structure of each completed sentence.

→ "Yes" = ALL completed sentences follow proper grammar rules → "Some do not make sense" = SOME completed sentences have grammatical errors → "No" = ALL completed sentences contain grammatical errors

Focus on: subject-verb agreement, article usage, word order, tense consistency

2. KNOWLEDGE ALIGNMENT: How well do completions reflect unusual/anomalous visual elements in the image? IMAGE CONTEXT: Rate how accurately completions describe the strange/unusual elements you can SEE in the image.

SCORING RUBRIC: • 5 = Perfectly captures unusual/anomalous visual elements in the image - highly accurate and relevant • 4 = Mostly accurate - minor

Metric	Mode	Gemini-2.5-Flash	GPT-4.1	Average	Exact Agreement
Grammatical Correctness	Text-only	2.95 ± 0.23	2.93 ± 0.26	2.94 ± 0.25	95.0%
	With Image	2.92 ± 0.30	2.93 ± 0.26	2.93 ± 0.28	92.8%
Knowledge/Anomaly Alignment	Text-only	4.55 ± 0.98	4.31 ± 0.94	4.43 ± 0.97	69.5%
	With Image	4.76 ± 0.93	4.60 ± 0.91	4.68 ± 0.92	80.2%

Table 1: **LLM-as-a-judge evaluation results.** Gemini-2.5-Flash and GPT-4.1 for Text-only and Image-based Scenarios

gaps or imprecisions • 3 = Partially accurate - some correct elements, some missing/wrong • 2 = Minimally accurate - mostly incorrect with few relevant elements • 1 = Completely inaccurate - fails to reflect unusual/anomalous visual elements in the image
 SENTENCE: "{text}" COMPLETIONS TO EVALUATE: {tokens}
 INSTRUCTION: Look carefully at the image and identify what appears unusual, unexpected, or anomalous. Then evaluate how well the completions capture these strange elements.
 IMPORTANT: You must respond with valid JSON only. Do not include any text before or after the JSON.
 Required JSON format: {{"grammatically_correct": "Yes, all are grammatically correct" | "Some are not grammatically correct" | "No, they are not grammatically correct", "knowledge_reflection_score": 1-5, "reasoning": "Explain your evaluation for each criterion, citing specific completions" } }""

Prompt Used to Instruct LLM-as-a-Judge without image

You will evaluate sentence completions that should reflect normal, expected real-world knowledge real-world scenarios based on common knowledge and typical experiences.
 EVALUATION CRITERIA:
 1. GRAMMATICAL CORRECTNESS: Are the completed sentences grammatically valid?
 Check syntax, word order, agreement, and grammatical structure of each completed sentence.
 → "Yes" = ALL completed sentences follow proper grammar rules → "Some do not make sense" = SOME completed sentences have grammatical errors → "No" = ALL completed sentences contain grammatical errors
 Focus on: subject-verb agreement, article usage, word order, tense consistency
 2. KNOWLEDGE ALIGNMENT: How well do completions reflect "normal, expected real-world knowledge? REAL-WORLD CONTEXT: Rate how well completions reflect typical, widely-accepted

real-world scenarios.
 SCORING RUBRIC: • 5 = Perfectly captures normal, expected real-world knowledge - highly accurate and relevant • 4 = Mostly accurate - minor gaps or imprecisions • 3 = Partially accurate - some correct elements, some missing/wrong • 2 = Minimally accurate - mostly incorrect with few relevant elements • 1 = Completely inaccurate - fails to reflect normal, expected real-world knowledge
 SENTENCE: "{text}" COMPLETIONS TO EVALUATE: {tokens}
 Consider what would be normal, expected, and typical in real-world scenarios. Then evaluate how well the completions reflect this common knowledge.
 IMPORTANT: You must respond with valid JSON only. Do not include any text before or after the JSON.
 Required JSON format: {{"grammatically_correct": "Yes, all are grammatically correct" | "Some are not grammatically correct" | "No, they are not grammatically correct", "knowledge_reflection_score": 1-5, "reasoning": "Explain your evaluation for each criterion, citing specific completions" } }""

Instruction Given to Human Evaluator

Task Instructions You will be asked to evaluate sentences and their possible completions. Sometimes, an image will also be provided. Your job is to judge whether the completions are appropriate, whether the sentence is grammatically correct, and how well the sentence and completions reflect knowledge or the content of the image.
 Please follow these criteria carefully for each question:
 1. Is the sentence grammatically correct?
 Ignore meaning here; only focus on grammar and syntax. Mark "Yes" if the base sentence with all the possible completions is grammatically well-formed. Mark "Some do not make sense" if at least one completion create a grammatically incorrect sentence. Mark "No" if the sentence with all the possible completions contains grammar errors.
 2. How well do the sentence and completions reflect

common knowledge (or reflect strange/anomalous things in the image)?

If the question refers to common knowledge: Judge how typical, reasonable, or widely accepted the sentence + completion is.

Example: “*The sun rises in the east*” should score high (5).

Example: “*The sun rises in the north*” should score low (1).

If the question refers to strange/anomalous things in the image: Judge how accurately the sentence and completions capture unusual, odd, or unexpected elements visible in the image. Score higher if the completion clearly reflects what is strange in the image.

Score lower if it ignores or misrepresents the anomaly. Use the 1–5 scale consistently:

1 = Not at all accurate/appropriate

3 = Neutral or partially accurate

5 = Very accurate and appropriate

Please read each question carefully and provide your evaluations with attention.

C MLP Intervention

We tested whether intervening on MLP blocks could produce effects comparable to those observed with attention heads. Specifically, we applied interventions to the last three MLP blocks at the final token position in both LLaVA-NeXT and Gemma3. The results, reported in Fig. 7, show only marginal changes in factual accuracy relative to the baseline. This effect is substantially weaker than the gains obtained from targeted attention-head interventions (Fig. 4).

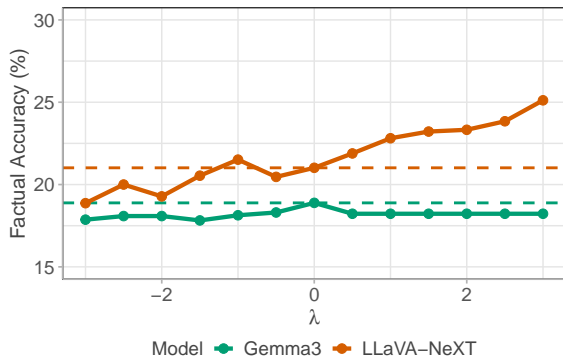


Figure 7: **Effect of MLP interventions.** Factual accuracy when intervening on the last three MLP blocks at the final token position in LLaVA-NeXT and Gemma3. The observed improvements are minor compared to targeted attention-head interventions (Fig. 4).

These findings reinforce two key conclusions.

First, MLP interventions are less precise: they modify the residual stream broadly, affecting a much larger number of parameters and acting indiscriminately across modalities. This broad influence makes it harder to isolate causal mechanisms and increases the risk of introducing unintended side effects. Second, the limited efficacy of MLP interventions indicates that factual–counterfactual conflicts are primarily mediated by attention mechanisms, not by MLP transformations. This aligns with prior evidence that late-layer MLPs often retrieve factual associations, whereas attention heads are more directly responsible for integrating conflicting cross-modal signals. For these reasons, our analysis centers on attention interventions, which provide both stronger causal leverage and more interpretable control over the balance between internal knowledge and visual input.

D Experimental Analysis for Gemma-12b

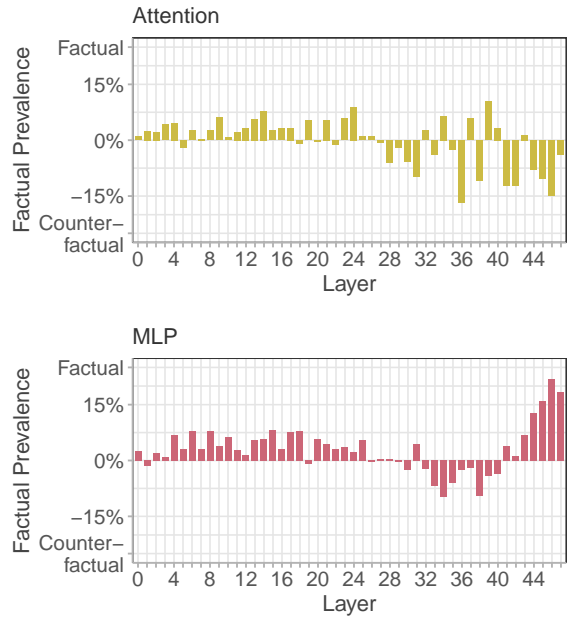


Figure 8: **Factual and counterfactual contributions of MLP and attention blocks in Gemma3.** Layer-wise deviation from 50% factual accuracy for attention and MLP blocks, as measured by the relative logits of t_{fact} and t_{cofa} via Logit Lens. Positive values indicate a bias toward the factual token, while negative values indicate preference for the counterfactual token. Consistent with trends observed in LLaVA-NeXT, attention blocks in Gemma3 increasingly support counterfactual predictions in higher layers, while MLP blocks show stronger alignment with internal factual knowledge.

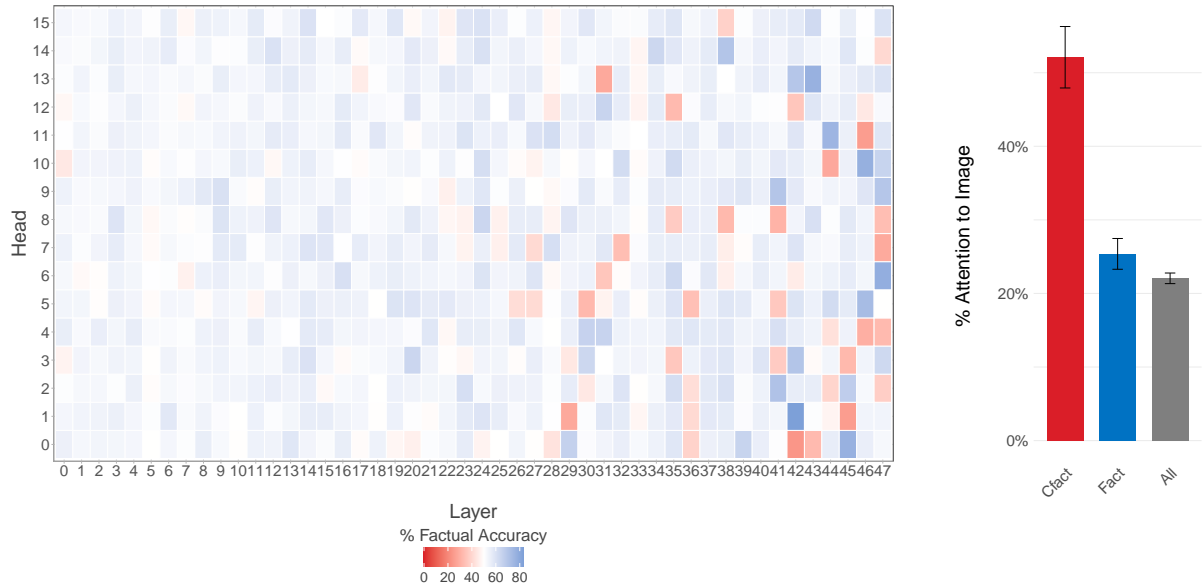


Figure 9: **Factual and counterfactual contributions of attention heads for Gemma3.** (Left) Factual accuracy of individual attention heads in Gemma3, computed using Logit Lens projections of the final token’s hidden state. Blue indicates heads that more frequently favor the factual token (t_{fact}), while red indicates those that favor the counterfactual token (t_{cofa}). As in LLaVA-NeXT, highly polarized heads are concentrated in the upper layers. (Right) Mean attention to image tokens at the final generation step. Counterfactual heads attend more strongly to image tokens (52%) than factual heads (25%) or the model-wide average (22%), highlighting the direct role of visual input in modulating counterfactual predictions.

E Details on the Number of Heads Selected and Control Experiment

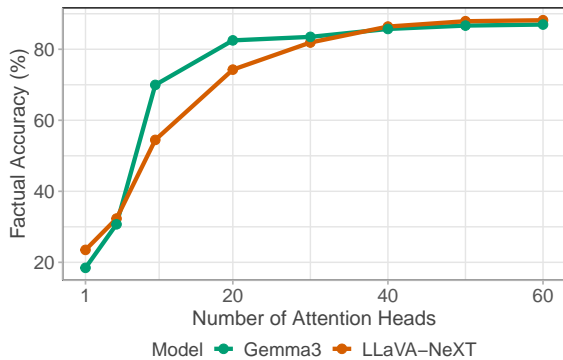


Figure 10: **Effect of intervening on varying numbers of attention heads.** Change in factual accuracy as a function of the number of attention heads involved in the intervention. Each value x indicates that x heads are selected from both the factual and counterfactual groups. Intervention strength is fixed at $\lambda = 3$. The results highlight that intervening on 20 heads provides the optimal trade-off, maximizing factual accuracy without excessively affecting model stability.

Figure 10 examines the effect of varying the number of intervened attention heads, with intervention strength fixed at $\lambda = 3$. We observe that

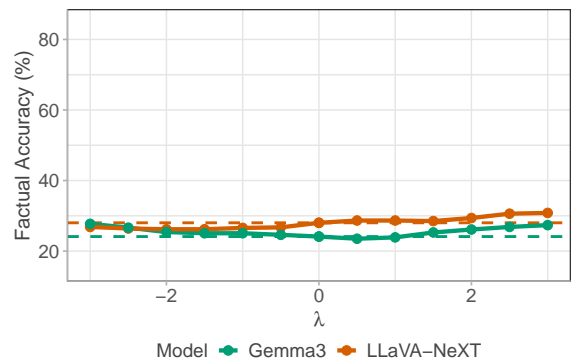


Figure 11: **Control experiment: intervention on random attention heads.** Change in factual accuracy under varying levels of intervention strength (λ) applied to 100 randomly selected attention heads. The results show no substantial deviation from baseline, confirming the specificity of the identified target heads.

factual accuracy increases as the number of heads grows, reaching its peak at 20 heads. Beyond this point, further interventions do not yield additional gains and may introduce instability. This demonstrates that intervening on 20 heads provides the best balance between accuracy improvement and model robustness. Figure 11 reports the control experiment, where the intervention was applied to

3, and then the growth slows down and stabilizes around $|\lambda| = 12$ for $\lambda < -20$ and 18 for $\lambda > 20$. When the KL is smaller than 10, for λ between -3 and 3, the output sentences have a similar quality to those generated before intervention.

G Generalization to Visual CounterFact

To evaluate whether the mechanisms for resolving factual and counterfactual conflicts identified with WHOOPS-AHA! generalize to other data distributions, we extended our analysis to the Visual CounterFact dataset (Golovanevsky et al., 2025a).

Experimental setup. We utilized the “color” split of Visual CounterFact, which consists of paired factual and counterfactual images depicting the same object with different color attributes (e.g., a standard red strawberry versus a counterfactual blue strawberry). This dataset offers a distinct visual domain characterized by more photo-realistic manipulation compared to the illustrative nature of some WHOOPS-AHA! samples.

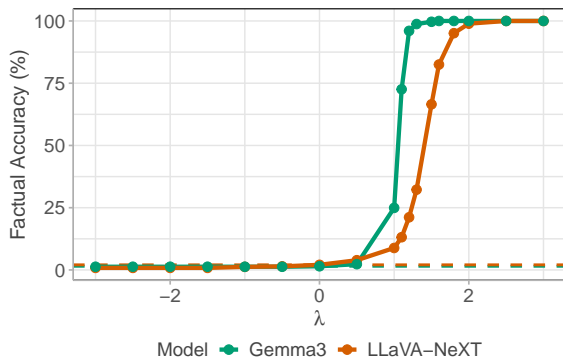


Figure 13: **Effect of intervention on factual accuracy in Visual CounterFact.** Positive values of λ effectively steer both models toward factual accuracy, confirming that the results generalize to this dataset.

Prompt construction. For each image, we constructed a completion prompt of the form: “The color of the [object] is”. We then analyzed the model’s logits for the specific factual and counterfactual color tokens provided by the dataset (e.g., “red” vs. “blue”).

Procedure. We replicated the procedure described in Section 5.2 and Section 5.3:

1. We identified the top-20 factual and counterfactual attention heads specific to this dataset by analyzing the Logit Attribution on the final token position.

2. We performed the attention intervention using the parameter λ to steer the model between parametric knowledge (factual) and visual evidence (counterfactual).

Intervention efficacy. Figure 13 illustrates the effect of intervention on LLaVA-NeXT and Gemma3. Consistent with our main results, increasing the activation of factual heads ($\lambda > 0$) significantly boosts the model’s reliance on internal parametric knowledge. Due to the naturally high success rate of the visual override in this dataset (low factual baseline), the effect of negative λ is negligible, as the model is already effectively attending to the visual evidence.

Mechanism transferability. A key question is whether the specific attention heads mediating this conflict are consistent across datasets. We compared the top-20 heads identified on Visual CounterFact with those identified on WHOOPS-AHA!. We observed a strong overlap:

- 13 out of 20 counterfactual heads in LLaVA-NeXT and 14 out of 20 counterfactual heads in Gemma3 are identical.
- 10 out of 20 factual heads are identical across the two models.

This substantial overlap suggests that the identified heads are not specific to WHOOPS-AHA!. Instead, they appear to reflect a stable and transferable mechanism for resolving conflicts between visual input and parametric knowledge across independently constructed datasets and varying degrees of photo-realism.

H Specificity of Conflict-Resolution Heads: Control on General Visual Understanding

A potential alternative explanation for the role of the identified counterfactual heads is that they function as generic vision-centric heads—that is, heads broadly responsible for processing visual information regardless of whether a knowledge conflict is present. Under this hypothesis, suppressing their image attention should degrade performance on any visually grounded task, not only on conflict-resolution ones.

To test this directly, we evaluate model performance on POPE (Li et al., 2023), a standard binary VQA benchmark consisting of 1,000 yes/no questions about object presence in images. POPE is a

heavily vision-dependent task: removing the image entirely reduces accuracy to chance level (0.50) for both models, confirming that correct answers cannot be recovered from text alone.

We apply the same intervention used in our main analysis: we zero out the image attention of the identified counterfactual heads at the final token position (equivalent to $\lambda = 1$ applied to image tokens of $\mathcal{H}_{\text{cofa}}$), and measure accuracy before and after intervention. Results are reported in Table 3.

Table 3: **Effect of suppressing counterfactual head image attention on POPE accuracy.** Removing the image entirely collapses accuracy to chance (0.50), confirming the task is vision-dependent. Suppressing image attention only in the identified counterfactual heads leaves accuracy unchanged, indicating these heads are not involved in general visual processing.

Model	No image	Baseline	After intervention
Gemma3-12B	0.50	0.84	0.84
LLaVA-NeXT-7B	0.50	0.87	0.87

The results show no measurable degradation in POPE accuracy after suppressing image attention in the identified heads. This stands in sharp contrast to the no-image condition, where accuracy collapses to chance. The absence of any effect indicates that the identified heads are not required for general visual understanding, and therefore cannot be characterized as generic vision-centric components.

This finding reinforces our interpretation: the counterfactual heads do not attend to image tokens simply because they process visual input, but specifically because a competition between parametric knowledge and visual evidence is active. Their selective involvement under conflict conditions, combined with their irrelevance to routine visual processing, supports the claim that they constitute a dedicated mechanism for cross-modal conflict resolution.

I Quantitative Analysis of Visual Attribution

To validate the qualitative observation that counterfactual attention heads act as precise visual pointers, we performed a quantitative evaluation on a subset of 20 sampled images. For each image, we obtained ground-truth object masks corresponding to the counterfactual concept (e.g., the "oranges" in the surgical scene) using the SAM3 model (Carion et al., 2025).

Metric definition. We evaluate the alignment between the attention maps and the ground-truth masks using the *Average Ratio*. This metric compares the attention intensity on the object versus the background. For a set of attention heads (e.g., the identified counterfactual heads), we first compute the average attention map across the heads for the given image. The ratio is then defined as:

$$\text{AvgRatio} = \frac{\frac{1}{|O|} \sum_{i \in O} \bar{A}_i}{\frac{1}{|B|} \sum_{j \in B} \bar{A}_j} \quad (3)$$

where O is the set of visual tokens falling within the counterfactual object mask, B is the set of background visual tokens, and \bar{A}_i represents the attention value at visual token i averaged across the selected $N = 20$ heads. A value > 1 indicates that the mechanism collectively focuses more intensely on the object than on the background.

Results. We compared our 20 identified counterfactual attention heads against two baselines: (1) a gradient baseline, using standard gradient-based attribution of visual tokens that are responsible for the counterfactual prediction, and (2) a random baseline, obtained by 20 randomly sampled attention heads.

Table 4: **Quantitative comparison of visual attribution methods on 20 sampled images.** We report the median and IQR. The p -value indicates significance against the counterfactual heads method selection (Wilcoxon signed-rank test).

Model	Method	Average Ratio (Median [IQR])	p -value (vs. Ours)
Gemma 3	Counterfactual heads	4.41 [1.09 – 7.85]	–
	Gradient	1.74 [1.09 – 2.01]	< 0.01
	Random heads	0.92 [0.71 – 1.22]	< 0.01
LLaVA-NeXT	Counterfactual heads	2.05 [1.48 – 3.33]	–
	Gradient	1.88 [1.26 – 2.44]	0.003
	Random heads	1.09 [0.86 – 1.25]	< 0.01

Table 4 summarizes the results for both Gemma3 and LLaVA-NeXT. We report the Median and Inter-Quantile Range (IQR) as the distributions are non-normal. We also perform a Wilcoxon signed-rank test to assess whether the counterfactual heads assign significantly higher attention to the segmented regions than the baselines. In both architectures, the counterfactual heads significantly outperform the baselines, achieving a much higher contrast ratio.

For Gemma3, the identified counterfactual heads are highly selective with a median contrast of 4.41, indicating that object patches receive over four times the attention intensity of background patches. For LLaVA-NeXT, the median ratio is 2.05, meaning the object is attended to with double the intensity of the background. In both cases, the counterfactual heads show a statistically significant improvement over the gradient and random baselines, confirming that the identified mechanism acts as a precise visual pointer across different VLM architectures.

J Prompts For Dataset Generation

Prompt Used to Generate Dataset Instances.

You are a helpful assistant expert in LLMs research. Counterfactual Dataset Generation Prompt

Objective: Generate captions for images that highlight a clear contrast between common (factual) and unusual (counterfactual) scenarios involving the subject depicted. Each caption must include the subject of the image and end with "___" indicating the blank space where a single-word token is placed.

Definitions: - **Factual token**: A single word that represents typical, expected behavior or attributes of the main subject shown in the image. - **Counterfactual token**: A single word introducing a surprising, unexpected, or unusual element related explicitly to the same main subject; it makes sense only if the image explicitly illustrates this twist.

Context Provided: For each image, you will receive the following textual information: - **Selected Caption**: A primary description identifying the main subject clearly. - **Crowd Captions**: Alternative descriptions from multiple annotators. - **Designer Explanation**: Explanation emphasizing the unusual or counterintuitive aspect involving the subject. - **Crowd Explanations**: Multiple explanations focusing on the unusual aspects related directly to the subject of the image.

Task Instructions:

Caption Construction: - Create exactly one neutral sentence (caption) clearly containing the main subject depicted in the image, but avoiding the description of unusual aspects contained in the image. - The sentence must end with an intentional blank ("___"). - **Critical Requirement**: The caption must compel the model to complete the blank differently based on the context: - **Without the image**: complete with a factual token (typical scenario involving the subject). - **With the image**: complete with a counterfactual token (unexpected scenario explicitly depicted). - **Important Constraint**: Use neutral language with NO textual hints indicating abnormality. The main

subject must explicitly appear in the caption to establish context clearly. Only the image content itself should disambiguate the scenario. - The caption should not contain any unusual or counterintuitive elements; the unusual aspect should be reflected solely in the image content and in the counterfactual tokens. - Make sure that if you substitute the blank with a factual or counterfactual token, the sentence is fluent and grammatically correct.

Explicit Single-Word Token Generation: - Generate exactly **ten single-word factual tokens** representing common scenarios involving the main subject that could complete in a grammatically correct way the sentence. - Generate exactly **ten single-word counterfactual tokens** representing surprising scenarios involving the same subject, justified solely by the provided image, and that could complete the sentence in a grammatically correct way. - Strictly enforce single-word tokens; no multi-word phrases or sentences. - Ensure clear differentiation without conceptual overlap between factual and counterfactual tokens.

JSON Output Format: Provide each caption and tokens following this exact schema:

```
{ "caption": "Neutral sentence explicitly containing the main subject and ending with an intentional blank ('___')", "factual_tokens": ["token1", "token2", "token3", "token4", "token5", ...], "counterfactual_tokens": ["token1", "token2", "token3", "token4", "token5", ...], "context": { "selected_caption": "Primary description clearly stating the main subject of the image", "crowd_captions": ["Caption 1", "Caption 2", "..."], "designer_explanation": "Explanation highlighting the unusual aspect directly involving the main subject", "crowd_explanations": ["Explanation 1", "Explanation 2", "..."] } }
```

Your role is to craft neutral captions explicitly containing the main subject of each image, along with precisely differentiated factual and counterfactual single-word tokens. The explicit presence of the main subject in the caption must guide factual versus counterfactual completions, relying solely on the provided image for disambiguation.

Prompt Used to Generate Factual and Counterfactual Tokens.

You are presented with an image and an incomplete sentence describing its content. The image intentionally portrays an unusual scenario that contrasts typical or factual knowledge.

Your task is to generate two lists of tokens:

1. **Factual Tokens (5 tokens)**: These tokens should represent words or concepts that accurately and typically complete the sentence based solely on common knowledge, without considering the unusual image.

2. Counterfactual Tokens (5 tokens): These tokens should represent words or concepts that correctly complete the sentence when explicitly considering the unusual content depicted in the image, even if it contradicts common factual knowledge.

Please format your response clearly as a JSON object as follows:

```
““json { "sentence": "INCOMPLETE_SENTENCE",  
"factual_tokens": ["token1", "token2", "token3",  
"token4", "token5"], "counterfactual_tokens": ["to-  
ken1", "token2", "token3", "token4", "token5"] } ““
```

Choose tokens that clearly differentiate between typical knowledge and the unusual scenario depicted by the provided image.