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

Anonymous ACL submission

1

Abstract

Vision-language models (VLMs) increasingly 002 leverage diverse knowledge sources to address complex tasks, inevitably encountering conflicts between their internal parametric knowledge and external information. Knowledge conflicts often result in hallucinations and unreliable responses, but the mechanisms governing 007 such interactions remain unknown. To address this gap, we analyze the mechanisms VLMs use to resolve cross-modal conflicts by introducing a dataset of multimodal counterfactual queries 011 that deliberately contradict internal commonsense knowledge. We localize with logit inspec-013 tion a small set of heads that control the conflict. 014 Moreover, by modifying these heads, we can steer the model towards its internal knowledge or the visual inputs. Finally, we show that atten-017 tion from such heads pinpoints localized image regions driving visual overrides, outperforming 019 gradient-based attribution in precision.¹

Introduction 1

021

030

034

037

Vision-language models (VLMs) (Alayrac et al., 2022; Li et al., 2022; Liu et al., 2023; Team, 2024; Deitke et al., 2024) have shown a remarkable versatility across various multimodal tasks, from image understanding to image generation. They draw on their ability to combine a rich set of world knowledge acquired during training, while also integrating contextual information provided in the prompts. However, these two sources of information can contradict each other, such as when the pretraining knowledge becomes outdated (Lazaridou et al., 2021; Luu et al., 2022) or when intentionally misleading visual cues are injected into the prompt (Liu et al., 2024d). These conflicts often trigger hallucinations and mistakes (Cui et al., 2023; Liu et al., 2024a; Guan et al., 2024), and little is known

1



The wolf is howling at the

╶

sur

moon

struct prompts that induce a conflict between a visionlanguage 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.

about the internal mechanisms employed by VLMs to resolve this conflict (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), which are followed by a sentence encouraging a typical, knowledge-based continuation. As shown in Fig. 1, each input prompt is associated with a

Vision Language Model

Counterfactual Predictions

Factual Predictions

Mark Zuckerberg wears

a shirt with a logo of

amazon

facebook

 h_i

Factual heads

 h_i Counterfactual heads

¹Our code and data have been uploaded to the submission system and will be open-sourced upon acceptance.

⁰⁴⁰ 042 043 044 045 047

counterfactual pair of completions. For instance, 049 the model may be shown an image of a wolf howling at the sun, a scene that contradicts common-051 sense knowledge, and asked to complete the prompt accordingly (see top-left panel). We construct the dataset such that VLMs, when prompted with text alone, generate factual responses while in the presence of the image, change their prediction to align with the visual context, even when it contradicts 057 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 versus counterfactual predictions. We find 061 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 prove more effective than gradient-based methods at identifying which parts of an image are most important for resolving 067 multimodal conflicts in VLMs.

In summary, our contributions are as follows:

- 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 (Sec. 4.1);
 - We identify the attention heads that promote factual and counterfactual responses, ranking their importance with logit attribution (Sec. 4.2);
- 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 knowl-edge and vice versa (Sec. 4.3);
- 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 (Sec. 4.4).

2 Related Work

072

077

084

086

090

096

Most prior work on knowledge conflict has focused on language models and unimodal tasks, leaving the multimodal domain underexplored (Xu et al., 2024).

The analyses of knowledge conflicts in language models have largely been behavioral, showing that when resolving conflicts between contextual and internal knowledge, language models can overrely 097 on their internal knowledge or contextual infor-098 mation, depending on factors such as model size 099 (Longpre et al., 2021) and conflicting external in-100 formation (Chen et al., 2022). Wang et al. (2024) 101 found that even SOTA language models often fail 102 to report inconsistencies between in-context infor-103 mation and their internal knowledge. Few works 104 have analyzed the internal mechanisms underlying 105 conflict resolution. Ortu et al. (2024) identified two 106 heads that mediate between factual and counterfac-107 tual information, while Jin et al. (2024) showed 108 that pruning specific heads can steer the model's 109 reliance toward internal or contextual sources. In 110 the multimodal domain, studies on VLMs have 111 primarily focused on benchmark construction and 112 black-box evaluation. Han et al. (2024) introduced 113 a dataset probing contextual knowledge conflicts 114 introduced by deceptive visual elements in prompts. 115 Golovanevsky et al. (2025) proposed NOTICE, 116 using semantically corrupted image pairs to an-117 alyze attention heads behavior in LLaVA and BLIP. 118 Liu et al. (2024c); Guan et al. (2024) developed 119 ConflictVis to study conflicts between visual in-120 put and parametric knowledge, but restricted their 121 analysis to the prompt structure rather than internal 122 mechanisms. 123

In contrast, in this work, we focus on model internals, identifying specific attention heads responsible for mediating factual and counterfactual reasoning, and validating their roles through targeted ablations.

124

125

126

127

128

129

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

3 Background and Methods

3.1 Model Architectures

This study investigates how visual input interacts with the model's internal knowledge during text generation in VLMs. Given a sequence of k imagetext tokens, a VLM encodes the image using a vision encoder and the text using an embedding matrix, producing the residual stream $\mathbf{x} \in \mathbb{R}^{d \times k}$, where d is the hidden dimension of the model. We denote the residual stream at position i and layer l as \mathbf{x}_i^l . The residual stream is processed through L layers, each composed of an attention block \mathbf{a}^l and an MLP block \mathbf{m}^l . After the final layer, it is projected to the vocabulary space via an unembedding matrix $W_U \in \mathbb{R}^{d \times |V|}$. Formally, the update of the residual stream at the l^{th} layer is:

$$\mathbf{x}^{l} = \mathbf{x}^{l-1} + \mathbf{a}^{l} + \mathbf{m}^{l} , \qquad (1)$$

146 147

- 148
- 149

150

151

153

154

155

156

157

158

160

161

162

163

166

167

169

170

171

173

174

175

177

178

179

181

183

190

191

193

where both the attention and the MLP block take as input the x after layer normalization norm:

$$\mathbf{a}^{l} = \mathbf{a}^{l}(\operatorname{norm}(\mathbf{x}^{l-1})), \qquad (2)$$

$$\mathbf{m}^{l} = \mathbf{m}^{l}(\operatorname{norm}(\mathbf{x}^{l-1} + \mathbf{a}^{l})).$$
 (3)

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.

3.2 Dataset Construction

To study how VLMs handle conflicts between visual context and internal knowledge, we introduce WHOOPS-AHA!, a new dataset designed to induce controlled competition between the two information sources. Each example in WHOOPS-AHA! consists of (i) a counterfactual image, (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. We construct our dataset on top of the WHOOPS! collection (Guetta et al., 2023), which consists of 500 images illustrating visually implausible scenes, each annotated with descriptions of the image content and the underlying anomaly. For each image in WHOOPS!, we use GPT-40 to generate a sentence that references the anomaly, while remaining consistent with commonsense (factual) completion without visual input. GPT-40 is also prompted to produce a set of plausible factual tokens S_{fact} and visually-grounded counterfactual continuations S_{cofa} . For instance, for the case of 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", $S_{\text{fact}} = \{\text{"moon", "night",...}\} S_{\text{cofa}} = \{\text{"sun",}$ "daylight", "morning", ... }. All generated content is manually verified to ensure a clear distinction between factual and counterfactual continuations. Full prompt details are provided in appendix C.

3.3 Analytical Tools

Logit Inspection To identify the internal components of VLMs responsible for the competition between inner knowledge and conflicting visual context, we trace the evolution of token logits across the model's architecture. Specifically, 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 (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 prompt and extract the logits corresponding to the tokens in S_{fact} and S_{cofa} at various layers and components of the model to identify the components that contribute to the promotion of one mechanism over the other.

194

196

197

199

200

202

203

204

205

206

207

209

210

211

212

213

214

215

216

217

218

219

220

221

223

224

226

227

229

230

231

232

233

234

235

236

237

238

239

240

241

242

Targeted Intervention on Attention Heads То test the causal role of specific attention heads in promoting predictions aligned with either factual inner knowledge or counterfactual visual context, we intervene directly on their attention patterns during inference. We define two groups of heads based on Logit Inspection: factual heads (\mathcal{H}_{fact}), which favor predictions based on inner knowledge, and counterfactual heads (\mathcal{H}_{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

$$\mathbf{A}_{\text{last,img}}^{hl} \leftarrow (1+\lambda) \cdot \mathbf{A}_{\text{last,img}}^{hl} \tag{4}$$

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

Δ

$$\mathbf{A}_{\text{last,text}}^{hl} \leftarrow (1 - \lambda) \cdot \mathbf{A}_{\text{last,text}}^{(hl)}$$
(5)

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. To determine the number of heads to include in each group, we experiment with different group sizes ranging from 5 to 60 heads. We select 20 heads of the configuration that offers the best trade-off between the effectiveness of the intervention and the stability of the model's output. Stability is measured by tracking the rank position of the two representative tokens (t_{fact} and t_{cofa}) in the model's nexttoken logit distribution, ensuring that the higherranked token remains within the top 80 positions for Gemma3 and the top 30 for LLaVA-NeXT.



Figure 2: Factual Prevalence in Attention and MLP Blocks. The plot shows the factual prevalence of attention and MLP blocks in LLaVA-NeXT across layers, indicating whether each component promotes predictions aligned with factual knowledge or counterfactual visual context. Positive values correspond to blocks favoring the factual (commonsense) continuation. Negative values indicate preference for the counterfactual continuation induced by the image. 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).

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.

244

245

246

247

248

249

250

251

253

257

261

- 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. 262

263

264

265

267

269

270

271

272

273

274

275

277

278

279

281

282

283

284

285

287

288

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

3.4 Reproducibility

We run 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 15 GPU hours. The WHOOPS! dataset was released with a CC-By 4.0 license.

4 Results

4.1 Inducing the Conflict between Inner Knowledge and Visual Context

To systematically induce competition between visual input and internal knowledge, we construct the WHOOPS-AHA! dataset as described in Sec. 3.2. Each example of WHOOPS-AHA! includes a counterfactual image, a sentence describing the image, and two sets of plausible next-word candidates proposed by GPT-40: S_{fact} , consistent with commonsense knowledge, and S_{cofa} aligned with the counterfactual visual context. For each model, first identify t_{fact} as the token in S_{fact} with the highest probability using only the textual part of the prompt. We consider only the first token if a candidate word is tokenized into multiple tokens. Then, using the full multimodal input (image and text), we select $t_{\rm cofa}$ as the token with the highest probability from S_{cofa} . For example, when prompted with the sentence "The wolf is howling at the", LLaVA-NeXT and Gemma3 predict the factual token moon with probabilities of 78% and 100%, respectively. However, when the corresponding image is included, both models shift to the counterfactual token sun, with probabilities of 26% (LLaVA-NeXT) and 44% (Gemma3), while the probability of moon drops to 17% and 0.02%. We filter out ambiguous cases in which $S_{\rm cofa}$ contains tokens with a probability higher than S_{fact} in the text-only setup, keeping 436 examples for LLaVA-NeXT and 432 for Gemma3. In the following sections, we always prompt the model with image and text using t_{fact} and t_{cofa} to assess whether different model components promote internal knowledge or contextual information. Notably, introducing the image reduces the preference of the



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.

model for the commonsense token: the prediction of the factual token t_{fact} drops to 27% for LLaVA-NeXT and 24% for Gemma3. 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.

311

312

316

317

319

322

324

325

331

335

4.2 The Tension Between Inner Knowledge and Visual Context is Localized

Building on the controlled knowledge conflict induced by WHOOPS-AHA!, we now study how the competition between factual and counterfactual continuations is resolved internally and which components mediate it. To do this, we use the Logit Lens technique to analyze the hidden state at the *final token position* of the prompt, after each attention block and MLP, projecting it into the vocabulary space (see Sec. 3.3). We then compute, across the dataset, how often the logit of the factual token t_{fact} is larger than that of the counterfactual token t_{cofa} . This gives an accuracy score for each component that reflects whether it tends to promote the factual or counterfactual mechanism. To measure the strength of this tendency, we compute the factual preference strength, which is defined as the difference between the fraction of examples for which $t_{\text{fact}} > t_{\text{cofa}}$ and 0.5, the random baseline. A value near zero indicates no consistent tendency to favor factual versus contextual information across the dataset, while higher values reflect stronger, more polarized behavior. This method allows us to localize the components that modulate the interaction between visual inputs and internal knowledge.

336

337

338

340

341

342

344

345

346

347

348

349

350

352

354

356

360

Functional Separation Between Attention and MLP Layers. We first compare the contributions of attention and MLP blocks to the prediction of t_{fact} and t_{cofa} . Figure 2 shows the results for LLaVA-NeXT (see appendix A for similar results on 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. This pattern is consistent with previous findings on the role of upper-layer MLPs in retrieving factual knowledge

(Geva et al., 2021; Meng et al., 2022; Dai et al., 2022).

361

390

391

400

401

402

403

404

405

406

407

408

409

410

411

Localization of the Modality Conflict to Individual Attention Heads. We next examine the role 364 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 Appendix, Fig. 8 for Gemma3). The distribution shows that only a small subset of heads exhibit a strong, consistent alignment with t_{fact} or t_{cofa} . Moreover, consistently with the results at the block level, these factual and counterfactual heads are 372 concentrated in the final layers of the model, indicating that the conflict between inner knowledge 374 and visual context is resolved late in the forward pass. In the analyses of the next sections, we focus 376 on the 20 heads that promote the factual and counterfactual tokens more strongly. On average, the factual heads favor the t_{fact} 85% of the time, and the counterfactual ones 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 directlyin late layers of the model, with limited intermediate processing. These findings are



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.

consistent for Gemma3, and we report the analysis in appendix A.

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

4.3 Targeted Intervention on Selected Attention Heads Causally Shifts Model Behavior

Having identified attention heads aligned with either factual knowledge or counterfactual visual context, we next examine whether these components play a causal role in shaping model predictions. To this end, we apply the targeted intervention strategy introduced in section 3.3, modifying the attention weights to steer the output of the model towards one mechanism or the other. Guided by our earlier observation that counterfactual heads attend more to visual tokens, we design a bidirectional intervention that selectively adjusts attention values based on head type and token modality. For counterfactual heads, we modify their attention to image tokens; for factual heads, we target their attention to text tokens. In both cases, we apply a multiplicative adjustment at the final token position. Each intervention simultaneously enhances the attention of one group to its relevant modality while suppressing the other group's attention, for instance increasing the attention to image tokens for counterfactual heads while reducing attention to text tokens for factual heads, and vice versa. This approach enables us to modulate the relative influence of factual and counterfactual mechanisms on the model prediction.

Figure 4 shows the results of our intervention

for LLaVA-NeXT (orange profile) and Gemma3 443 (green profile). For LLaVA-NeXT, the baseline 444 accuracy, defined as the proportion of examples in 445 which the factual token t_{fact} receives a logit higher 446 than the counterfactual token t_{cofa} , is 27%. When 447 we increase attention from factual heads and de-448 crease it from counterfactual heads, the factual ac-449 curacy increases to 82%, indicating a strong shift to-450 wards predictions of inner knowledge. Conversely, 451 reversing the intervention reduces the accuracy to 452 20%, confirming that these heads causally influ-453 ence whether the model favors factual or counter-454 factual content. A similar trend can be observed 455 for Gemma3, with an even stronger relative effect 456 driven by its lower baseline factual accuracy of 457 24% and a peak of 85%. To ensure plausible in-458 terventions, we constrain the scaling parameter to 459 $\lambda \in [-3,3]$ and monitor the position of the higher-460 logit token in the full next-token distribution. For 461 example, using LLaVA-NeXT, the average rank of 462 the token t_{fact} shifts from 3 at $\lambda = 0$ to 31 at $\lambda = 3$, 463 indicating that while the intervention is highly ef-464 fective, it introduces some deviation in the overall 465 logit distribution, an expected effect when strongly 466 467 modulating internal components. As a control, we randomly select 100 attention heads and apply the 468 same intervention for varying λ values. This ma-469 nipulation does not produce a substantial deviation 470 from the baseline, confirming that the observed ef-471 fects are specific to the heads identified as aligned 472 with factual or counterfactual mechanisms. The 473 complete results are reported in Appendix Fig. 9. 474

4.4 Counterfactual Predictions Depend on Localized Image Regions

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

In the previous sections, we analyzed the conflict between contextual information and internal knowledge using WHOOPS-AHA! prompts, which induce a competition between counterfactual visual cues and factual model knowledge. This analysis revealed that specific attention heads at the final token position mediate this conflict, with heads aligned with the visual context attending strongly to image tokens and thereby 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 questions, we perform two



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.

complementary analyses: (i) we identify the image patches most responsible for driving counterfactual predictions using attention and gradient-based attribution methods, as described in section 3.3; and (ii) we ablate the identified patches by setting the corresponding visual token embeddings to zero at the input of the transformer, and measure the resulting change in inner knowledge accuracy. In addition to the quantitative analysis, we inspect the selected image patches to assess whether they correspond to intuitive counterfactual regions or visually salient objects contradicting the model's internal knowledge. To test the specificity of our findings, we also perform a control experiment in which we randomly sample an equivalent number of image patches for ablation. This allows us to assess whether the identified regions are uniquely responsible for triggering counterfactual predictions or whether any removal of visual input affects the model's behavior.

494

495

496

497

498

499

500

501

502

503

505

506

507

508

509

510

511

512



Figure 6: Qualitative Examples of Visual Regions Driving Counterfactual Predictions. Highlighted image regions correspond to visual patches identified as most responsible for counterfactual predictions using attention-based attribution. In both examples, the model generates a visually grounded but factually incorrect token (e.g., rainbow, fruit) instead of the commonsense alternative (black, tissue). The highlighted areas align with semantically meaningful and visually anomalous content, demonstrating that counterfactual outputs are grounded in localized, interpretable image features.

Quantitative Analysis of Patch Attribution and 514 Ablation The results of the experiments are 515 shown in Figure 5. We observe that the ablation of 516 visual patches identified through attention-based at-517 tribution leads to a sharp and consistent increase in 518 factual accuracy as more pixels are removed (green 519 profiles). For instance, in the case of LLaVA-NeXT, factual accuracy improves markedly with the ablation of just 10-30% of the top-ranked patches 522 and eventually plateaus around 80%. This sug-523 gests that counterfactual predictions are primarily driven by a small, localized subset of visually 525 salient regions. Gradient-based attribution (shown in red) also yields a substantial increase in fac-527 tual accuracy, though the effect is less pronounced 528 and saturates earlier, suggesting lower precision in identifying counterfactual-driving regions. In contrast, ablating an equivalent number of randomly selected patches results in only minor fluctuations 532 in accuracy, never approaching the improvements 534 achieved through targeted attribution. These findings confirm the causal role of the identified regions and support the hypothesis that counterfactual signals are spatially localized and semantically specific. 538

Qualitative Analysis of Visual Attribution To assess the semantic coherence of the identified visual regions, we qualitatively examine examples where attribution methods highlight specific patches as responsible for counterfactual predictions (see Fig. 6. In many cases, these regions correspond to intuitive scenes that directly contradict commonsense knowledge, such as unusual objects, implausible substitutions, or visual features that override typical textual expectations. For instance, when the model predicts "rainbow" instead of "black" for a bearskin hat, the highlighted patches focus on the hat's unrealistic coloring (Fig. 6-top). Similarly, when "fruit" replaces "tissue" in a surgical scene, the attention centers on the bright, unexpected presence of oranges on the operating table (Fig. 6-bottom). These observations confirm that the model's counterfactual outputs are not arbitrary but grounded in semantically meaningful and localized image features.

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

587

5 Conclusion

In this work, we investigated how counterfactual visual inputs interact with the internal knowledge representations of VLMs during generation. To this end, we introduced WHOOPS-AHA!, a dataset that pairs visually anomalous scenes with textual prompts designed to elicit either a commonsense (factual) continuation or one grounded in the visual counterfactual. This setup enables fine-grained analysis of how conflicting visual and textual cues influence model behavior. We showed that a small set of attention heads mediate this competition. These heads also exhibit distinct modality preferences and play a causal role in determining the model's output. By intervening on their attention weights, we were able to shift predictions in a controlled way, favoring either the internal knowledge or the visual context. Finally, we demonstrated that these heads provide accurate attribution of the visual regions responsible for counterfactual completions, outperforming standard gradient-based attribution techniques. These findings contribute to a deeper mechanistic understanding of multimodal reasoning in VLMs and offer a foundation for developing more interpretable and controllable systems under conflicting input conditions.

Limitations

The analysis relies on the Logit Lens technique to project intermediate hidden states into token logits.

88	Although this method has been widely adopted for
89	interpretability, it is known to introduce distortions
90	due to projection from non-final residual states
91	(Belrose et al., 2023), and should be interpreted
92	as an approximate diagnostic rather than a precise
93	decoding proxy. In our setting, we use a represen-
94	tative factual and counterfactual token per example
95	to enable controlled comparisons. Although this
96	simplifies the generative landscape of the model,
97	it offers a practical and interpretable probe of the
98	underlying mechanisms. Future work could ex-
99	plore more model behavior across full generations
00	to complement this approach. Our attribution and
01	intervention methods focus on attention heads and
02	target the final token position. This design iso-
03	lates interpretable causal signals while remaining
04	tractable, though it does not capture the possible
05	contributions of other components, such as MLP
06	layers or visual encoders. Extending this frame-
07	work to broader architectural elements is a promis-
08	ing direction. Finally, the WHOOPS-AHA! dataset
09	is constructed from synthetic and curated inputs,
10	which allow precise manipulation of visual-textual
11	conflict. Although this setting facilitates analysis,
12	future extensions to more naturalistic data could
13	further validate the findings in less constrained con-
14	texts.

Ethical Considerations

5

5

5

5

5

5

5

615

This work aims to improve our understanding of 616 how VLMs resolve conflicts between internal fac-617 tual knowledge and contradictory visual context. 618 Our analysis is intended to contribute to founda-619 tional 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, 623 designed to support responsible development and 624 analysis of multimodal systems. We believe that 625 studying the dynamics of conflicting information sources is essential for anticipating model failure 627 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 631 for non-commercial, research purposes under compatible terms. We used AI assistants (e.g., GitHub Copilot) to support code completion during experi-634 ment implementation; all generated code was manually reviewed and supervised by the authors. 636

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, Mikoł aj 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. 637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653 654

655

656

657

658

659

660

661 662

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

682

683

684

686

687

688 689

690

691

692

693

694

695

696

697

- 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.
- Hung-Ting Chen, Michael Zhang, and Eunsol Choi. 2022. Rich knowledge sources bring complex knowledge conflicts: Recalibrating models to reflect conflicting evidence. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2292–2307, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- 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.
- 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, Vedant Palit, Carsten Eickhoff, and Ritambhara Singh. 2025. 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*

- 702 703 704 706 710 711 714 715 717 718 719 720 721 722 726 727 728 729 730 731 733 735 736 737 739 740 741
- 742 743 745 746 747 748 749 750 751 752 753 754 757 759 761 762

763

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.
- 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 Petrini, 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. 764

765

766

767

768

769

770

773

777

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

- 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.
- 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.
- 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 Repre*sentations.
- 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.
- Shayne Longpre, Kartik Perisetla, Anthony Chen, Nikhil Ramesh, Chris DuBois, and Sameer Singh. 2021. Entitybased knowledge conflicts in question answering. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7052–7063, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- 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*.
- Nostalgebraist. 2020. interpreting gpt: the logit lens. Accessed: Nov 2023.
- 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*

- 823 824 825 826 827 828 830 831 832 833 835 836 837 838 840 842 843 844 845 846 847 849 850 851 852
- 854

- 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.
- Chameleon Team. 2024. Chameleon: Mixed-modal early-fusion foundation models. arXiv preprint arXiv:2405.09818.
- Yike Wang, Shangbin Feng, Heng Wang, Weijia Shi, Vidhisha Balachandran, Tianxing He, and Yulia Tsvetkov. 2024. Resolving knowledge conflicts in large language models. In First Conference on Language Modeling.
- 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.
- 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.
- 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.



Figure 7: Factual and Counterfactual Contributions of MLP and Attention Blocks in Gemma3. Layerwise deviation from 50% factual accuracy for attention and MLP blocks, as measured by the relative logits of $t_{\rm fact}$ and $t_{\rm cofa}$ 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 8: 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.

B Additional Results



Figure 9: 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.

C Prompts For Dataset Generation

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": ["token1", "token2", "token3", "token4", "token5"]
}
```

Figure 10: **Prompt Used to Generate Factual and Counterfactual Tokens.** Given a fixed input sentence and an image from the WHOOPS! dataset, GPT-40 is prompted to propose candidate next-token completions. The prompt guides the model to return two sets of tokens: one reflecting commonsense completions consistent with world knowledge (factual), and one aligned with the visually depicted but counterfactual scene.

You are an helpfull 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 aspect

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.

Figure 11: Prompt Used to Generate Dataset Instances. We provide GPT-40 with an image, a set of captions, and an explanation of the visual anomaly, and instruct it to generate a sentence that implicitly refers to the anomaly while remaining commonsense-compatible. The model is then asked to propose plausible factual and counterfactual next-token completions, reflecting typical knowledge-based and visually grounded interpretations, respectively.