# A Comprehensive Analysis for Visual Object Hallucination in Large Vision-Language Models

Anonymous ACL submission

#### Abstract

Large Vision-Language Models (LVLMs) demonstrate remarkable capabilities in multimodal tasks, but visual object hallucination remains a persistent issue. It refers to scenarios where models generate inaccurate visual object-related information based on the query input, potentially leading to misinformation and concerns about safety and reliability. Previous works focus on the evaluation and mitigation of visual hallucinations, but the underlying causes have not been comprehensively investigated. In this paper, we analyze each component of LLaVA-like LVLMs-the large lan-014 guage model, the vision backbone, and the projector, to identify potential sources of error and 016 their impact. Based on our observations, we propose methods to mitigate hallucination for 017 each problematic component. Additionally, we developed two hallucination benchmarks: QA-VisualGenome, which emphasizes attribute and relation hallucinations, and QA-FB15k, which focuses on cognition-based hallucinations.

#### 1 Introduction

034

040

Large Language Models (LLMs), such as GPT-3 (Brown, 2020) and ChatGPT (OpenAI, 2022), have showcased remarkable proficiency in language tasks, yet they encounter significant challenges when it comes to processing multimodal inputs. This limitation has driven a shift in research towards Large Vision-Language Models (LVLMs) (Liu et al., 2023e; Ye et al., 2023; Sun et al., 2023b), which integrate advanced LLMs (Touvron et al., 2023; Chiang et al., 2023) with Vision Foundation Models (VFMs) (Dosovitskiy et al., 2021; Bommasani et al., 2021) to enhance multimodal understanding. LVLMs have demonstrated impressive capabilities across various tasks that require visual and textual integration, including Visual Question Answering (Antol et al., 2015), Image Captioning (Lin et al., 2014), and Visual Entailment (Zhang et al., 2025).



Figure 1: An overview of our paper. We first investigate the sources of hallucination from a component-level perspective within the LVLM architecture. Based on the identified causes, we then design targeted methods to mitigate hallucinations effectively.

042

043

044

045

046

047

049

051

052

058

060

061

062

063

064

065

Despite these advances, visual hallucination remains a persistent issue in LVLMs (Rohrbach et al., 2018; Liu et al., 2023b,a; Yin et al., 2023; Zhang et al., 2024b). This phenomenon occurs when models generate inaccurate or misleading information unrelated to the actual visual input, potentially leading to misinformation and raising concerns about safety and reliability in real-world applications (Li et al., 2023e). Visual object hallucination, including object existence, attribute, and relation, has garnered significant attention due to its widespread occurrence in images. Current works on visual object hallucination mainly focus on evaluation and mitigation. For example, Li et al. (2023e) extends CHAIR (Rohrbach et al., 2018) and proposes POPE, a polling-based query technique for probing object-level hallucination. For hallucination mitigation, Sun et al. (2023a) introduce new alignment algorithm called Factually Augmented RLHF that augments the reward model with additional factual information such as image captions and groundtruth multi-choice options, which alleviates the reward hacking phenomenon in RLHF and further improves the performance.

While existing works have achieved notable suc-

cess in visual object hallucination, they lack a 067 comprehensive component-level analysis of the 068 model architecture to pinpoint where and how hal-069 lucinations occur. In this work, we focus on visual object-related hallucination and LLaVA-like LVLMs, which typically consist of three modules: the large language model (LLM), the vision back-073 bone, and the projector. Errors in any of these modules can lead to issues in the overall performance or functionality of the model. Therefore, we conduct an independent analysis of each component to 077 identify potential sources of error and their impact. From our study, we have the following findings. 1) The LLM in LVLM is able to generate faithful content when captions of images are provided as input. 2) Hallucinations exist in the perception process of the vision backbone. 3) Projector is able to preserve visual features, but has trouble aligning between visual and textual spaces.

> Based on our observations, we propose methods for the two problematic components to mitigate their hallucination issue. To improve the **vision backbone**, we propose to finetune CLIP with fine-grained data and fine-grained perceptionbased visual instruction tuning, and find that both of them can reduce hallucination caused by the vision backbone. For the **projector**, we propose a contrastive alignment objective with three variations, which can all be integrated into the original training pipeline with minimal additional costs.

087

097

100

101

102

103

104

106

107

108

109

110

111

112

113

114

115

116

117

To conduct a comprehensive hallucination evaluation, we develop a fine-grained hallucination benchmark named QA-VisualGenome, which is built upon the Visual Genome dataset (Krishna et al., 2017). Unlike existing objectoriented hallucination benchmarks (e.g., POPE), QA-VisualGenome emphasizes the detailed attribute and relationship hallucinations. Furthermore, existing hallucination benchmarks primarily focus on perception-based hallucinations for general objects, neglecting cognition-based hallucinations such as the names of people and famous buildings. To address this gap, we construct a cognitionbased hallucination benchmark named QA-FB15K, which is based on the FB-15K dataset (Bordes et al., 2013), a multimodal knowledge graph with textual entities, image entities, and textual relations. QA-FB15K presents challenges for models in leveraging world knowledge to solve the questions.

Our main content is shown in Figure 1. Our contributions can be summarized as follows: 1) We analyze the hallucination caused by each component in LVLMs and provide component-wise takeaway messages. 2) Based on our observation, we propose several methods to improve each hallucinated component. 3) We construct a fine-grained hallucination benchmark based on Visual Genome and a cognition-based hallucination benchmark based on FB15k for evaluation. 4) We extensively evaluate our proposed methods on various benchmarks, and provide in-depth analysis<sup>1</sup>.

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

## 2 Hallucination Analysis

LVLMs consist of three components: language decoder  $\mathcal{D}$ , projector vision encoder  $\mathcal{V}$ , and  $\mathcal{P}$ . We first introduce the datasets for evaluation and then provide in-depth analysis for each component.

#### 2.1 Settings

We select two benchmarks to benchmark the performance of each component. 1) POPE (Li et al., 2023e). POPE is a benchmark designed for evaluating object existence hallucinations in LVLMs, incorporating three sampling methods for generating negative samples: random, popular, and adversarial. In the random setting, objects not present in the image are randomly selected. In the popular setting, negative samples are drawn from a pool of frequently occurring objects. In the adversarial setting, the sampling focuses on objects that frequently co-occur with present objects but do not exist in the image. 2) QA-VisualGenome. To further investigate the hallucination issue on relations and attributes of objects, we construct a fine-grained evaluation benchmark based on the VisualGenome dataset (Krishna et al., 2017), which collects dense annotations of attributes and relationships of objects for each image. Specifically, we design two types of Yes-or-No questions to evaluate models: attributes and relations. For example, an attribute question could be "Is the dog red in the image?" A relational question would ask, "Is the dog standing on the table?". Similar to previous work (Wang et al., 2020), we exclude uncommon relations and attributes. We randomly select relations/attributes to generate negative samples.

#### 2.2 Language Decoder

Conjecture 1. LLM in LVLM is able to generate faithful content when image captions are

<sup>&</sup>lt;sup>1</sup>All benchmark datasets, code, and models will be released.

Table 1: Performance (%) of LLMs across different datasets when visual information is provided in textual format. LLaVA: image+text query as input on original LLaVA model; Vicuna: caption+text query as input on Vicuna-1.5; Vicuna<sub>LLaVA</sub>: caption+text query as input on the Vicuna model in LLaVA (LLM undergone visual instruction tuning).

		POPE							QA-VisualGenome			
Model	Random		Рор	ular	Adver	rsarial	Attr	ibute	Rela	ation		
	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1		
LLaVA-7B	87.42	86.36	86.63	85.25	85.13	83.88	64.67	66.60	67.57	74.81		
Vicuna-7B	92.67	92.09	92.67	92.09	93.00	92.47	57.23	69.83	79.50	80.79		
$Vicuna-7B_{LLaVA}$	100.00	100.00	100.00	100.00	99.67	99.67	68.29	75.92	63.2	73.06		
LLaVA-13B	91.33	91.72	88.33	89.16	84.33	85.97	55.99	68.86	56.40	69.38		
Vicuna-13B	87.90	89.15	95.00	95.24	90.00	90.91	87.90	89.15	87.90	89.25		
Vicuna-13B <sub>LLaVA</sub>	99.67	99.67	99.67	99.60	99.33	99.33	75.41	80.10	84.30	84.29		

provided as input. To validate this conjecture, 164 we use the POPE dataset to evaluate the perfor-165 166 mance of LLMs. Instead of providing images to the LVLMs, we only input text descriptions of the images. For POPE, we obtain objects from 168 the MSCOCO (Lin et al., 2014) dataset and feed the LVLM with objects in the image and the tex-170 tual query from POPE to generate the response. For QA-VisualGenome, we feed the LVLM with 172 objects, object attributes, and relations presented 173 in the image to replace visual information. This 174 helped assess the model's ability to hallucinate 175 when provided with accurate textual descriptions 176 of the image. In addition, we also test the original Vicuna as a baseline. 178

167

171

179

183

184

185

190

191

192

193

194

195

196

197

199

200

We show the performance of LLMs in Table 1. From the results, we found that the performance will be improved largely if we provide the correct visual information in a textual format. This indicates the current main reason for hallucination is caused by a vision encoder or projector. Specifically, the model could achieve an accuracy of 99.67% when provided with complete object descriptions for the random setting of POPE, which shows the LLM is robust when given the correct information about the whole image. In addition, we also found that the LLM after the pertaining and instruction tuning of LLaVA performs better than the original LLM. LLaVA fine-tuning likely enhances the model's object recognition, memory of object-specific features, instruction-following ability, and contextual understanding of visual descriptions, enabling it to accurately identify common objects within text descriptions even without actual images.

#### 2.3 Vision Encoder

Conjecture 2. There are hallucinations in the perception process of the vision encoder. To verify this factor, we conducted experiments using Table 2: Performance of CLIP in the text-image matching across different datasets measured by Accuracy (%).

	POPE	QA-Visua	QA-VisualGenome		
Random	Popular	Adversarial	Attribute	Relation	
83.33	87.30	86.00	61.57	60.22	

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

223

224

225

226

227

229

230

231

232

233

CLIP on a text-image matching task. Specifically, we designed a template of the form "There is a/an {object} in the image," where {object} corresponds to various objects in the input images. For each image, we assigned one ground-truth object and a hallucinated object for the template. We use accuracy as the evaluation metric. We show all the experimental results in Table 2. Overall, we found that the performance of CLIP on the text-matching task is not good. For example, the performance of CLIP on the text-image matching task is 83.33% accuracy on the random setting of POPE, indicating the presence of hallucinations within the vision encoder's perception process.

Another interesting phenomenon is that the accuracy of CLIP in recognizing objects is worse than LLaVA, even the LLaVA adopts CLIP as the vision encoder. Specifically, the accuracy of LLaVA is 91.33% on the random setting of POPE, but CLIP only achieves 83.33% accuracy. This indicates that the hallucination caused by CLIP can be alleviated to a certain extent after the pre-training feature alignment and instruction tuning. The potential reason may be that LLaVA's training uses diverse questions aligned with specific image features, optimizing for generative loss. This fine-grained alignment helps the model better understand and describe visual content with greater accuracy and detail.

#### 2.4 Projector

We analyze the projector module from two perspectives corresponding to its two roles in the LVLM: preserving visual information and aligning visual

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

280

281

241 242

236

237

240

and textual spaces.

projector representations as

label).

Conjecture 3. The projector should not re-

sult in significant visual information loss. We

formalize the hypothesis using the notion of V-

information (Hewitt et al., 2021). Let  $\Phi_{pre}(X)$ 

and  $\Phi_{post}(X)$  represent the pre-projector and post-

projector representations, respectively. We com-

pare the V-information between these representa-

tions and a target property Y (e.g., a classification

 $I_{\mathcal{V}}(\Phi_{pre}(X) \to Y) = H_{\mathcal{V}}(Y) - H_{\mathcal{V}}(Y|\Phi_{pre}(X))$ 

 $I_{\mathcal{V}}(\Phi_{post}(X) \to Y) = H_{\mathcal{V}}(Y) - H_{\mathcal{V}}(Y|\Phi_{post}(X))$ 

where  $H_{\mathcal{V}}$  is the V-entropy (Hewitt et al., 2021).

 $H_{\mathcal{V}}(Y)$  is the entropy of Y, which reflects the in-

herent uncertainty of Y without any conditioning

on the representations.  $H_{\mathcal{V}}(Y|\Phi(X))$  represents

the uncertainty we have in predicting Y after ob-

serving the representation  $\Phi(X)$ , using functions from the family  $\mathcal{V}$ . It is formally defined as:

 $H_{\mathcal{V}}(Y|\Phi(X)) = \inf_{f \in \mathcal{V}} \mathbb{E}_{\Phi(X),Y} \left[ -\log f(\Phi(X))(Y) \right]$ 

This expression measures the best performance that a function f from the function family  $\mathcal{V}$  can achieve

when predicting Y given the representation  $\Phi(X)$ .

The lower this value, the more predictive power the

loss occurs in the projection layer. If the projection

layer introduces no information loss, then the V-

information of the pre-projector and post-projector representations should be approximately equal:

 $I_{\mathcal{V}}(\Phi_{pre}(X) \to Y) = I_{\mathcal{V}}(\Phi_{post}(X) \to Y)$ 

both the pre-projector and post-projector represen-

tations. The performance of a probe (e.g., classi-

fier) trained on  $\Phi_{pre}(X)$  and  $\Phi_{post}(X)$  provides an

 $Perf_{pre} = \max_{a} \mathbb{E}[\log P(Y|f_{\theta}^{pre}(\Phi_{pre}(X)))]$ 

 $\textit{Perf}_{\textit{post}} = \max_{\theta} \mathbb{E}[\log P(Y | f_{\theta}^{\textit{post}}(\Phi_{\textit{post}}(X)))]$ 

To determine if information loss occurs, we com-

empirical estimate of these quantities:

pute the difference in performance:

We compare the V-information accessible from

The goal is to determine whether information

representation  $\Phi(X)$  has regarding Y.

We define the V-information for pre- and post-

244 245

246

247 248

\_

25

25

25

25

257

25

260

26

262

264 265

20

267

268

21

269

270

272

273

274

275

275 276

0

278

279

$$\Delta Perf = Perf_{pre} - Perf_{post}$$

If  $\Delta Perf = 0$ , this implies that no information loss has occurred and the information available in  $\Phi_{pre}(X)$  is fully retained in  $\Phi_{post}(X)$ . However, if  $\Delta Perf > 0$ , this indicates that the post-projector representation has lost some information present in the pre-projector representation, leading to a decrease in predictive power for Y.

With the hypothesis grounded to V-information, we conduct a probing experiment on LLaVA-7B to verify it. We linear-probe the pre- and post-projector feature with image classification tasks on CIFAR10 (Krizhevsky et al., 2009), CI-FAR100 (Krizhevsky et al.) and ImageNet (Deng et al., 2009). Results in Table 3 shows that for the 13B LLaVA model, performance percentage drop of post-projection features is less than 2%, indicating that the visual features are well preserved by the projectors in both models.

Table 3: Performance of linear probing using pre- and post-projector image features on CIFAR10, CIFAR100 and ImageNet. Accuracy% is used as the metric.

Deternt	LLaVA-13B				
Dataset	Perf <sub>pre</sub>	Perf <sub>post</sub>			
CIFAR10	96.27	96.15 <u>-0.12</u> %			
CIFAR100	81.78	81.02 <u>-0.93%</u>			
ImageNet	71.97	70.83 <u>-1.58%</u>			

Conjecture 4. The projector should align the visual and textual spaces. As its name suggests, the projector should be able to project the source (visual) space to the target (textual) space. To probe the alignment between two spaces, we collect caption data from MSCOCO (Lin et al., 2014), LLaVA-Caption (Liu et al., 2023d), ALLaVA (Chen et al., 2024a) and compute the similarity between a projected image feature and the textual embedding of its caption. The rationale of using cosine similarity is that, based on the findings in Section 2.2, a large performance boost is observed if we replace an image with its caption. Therefore, if the projected image feature is similar enough to its caption embedding (*i.e.* cosine similarity=1), then an LVLM should gain similar performance to the case where an image is replaced by its caption as input.

Results in Table 4 show that the cosine similarities of the two features are fairly low, indicating nearly independent relationships. This finding is consistent with the existing work (Huang et al., 2024b; Li et al., 2025), which reveals that visual and textual representations are apart from each other in the embedding space. Therefore, the

Table 4: Cosine similarity between projected image features and the caption embedding. Captions are processed by Vicuna (Chiang et al., 2023) tokenizer.

Deteert	Talaan Lanath	Image Dec	Cos. Sim.		
Dataset	Token Length	image Res.	7B	13B	
MSCOCO	15.16	(575, 488)	0.03	0.04	
LLaVA Caption	15.09	(412, 366)	0.03	0.04	
ALLaVA	222.83	(1020, 923)	0.05	0.06	

projector in LLaVA models may not function as an alignment module as well as expected, which could be one of the causes of hallucination for the entire model.

# 3 Mitigating Object Hallucination Caused by Different Modules

323

331

333

334

337

339

341

342

343

344

346

347

351

354

359

361

Based on the analysis in Section 2, we further devised different methods to mitigate the object hallucination in different components in LVLMs.

# **3.1** How to alleviate the hallucination caused by CLIP?

As previously noted, the vision backbone within LVLMs also contributes to hallucinations. The CLIP model, as the vision encoder of LLaVA, is trained on massive image-caption pairs from the internet with a contrastive loss objective. However, these captions are typically brief and noisy, and negative pairs often differ substantially from positive ones. Therefore, it is likely that the model can distinguish them without needing to capture the finer details in the images. Consequently, the model may achieve high accuracy while lacking a nuanced understanding of the visual content (Liu et al., 2024b). To address this issue, we propose two methods to reduce hallucination caused by the vision backbone, as shown in Figure 2.

**Tuning CLIP with fine-grained data** A direct method to improve CLIP is to post-train CLIP with more fine-grained samples. This is because the CLIP is trained with massive images paired with brief captions. In this method, we leverage GPT-4 (OpenAI, 2022) to generate negative examples, which are then used in a contrastive learning setup to improve the discriminative ability of CLIP.

*Generate Negative Examples:* Inspired by prior work indicating that LVLMs are more likely to generate hallucinatory responses for frequently occurring objects (Liu et al., 2024b), we devise two strategies: inserting hallucinatory objects and removing existing ones.



Figure 2: Tuning CLIP with fine-grained data (left) and fine-grained perception-based instruction tuning (right).

362

363

364

365

366

367

368

369

370

371

373

374

375

376

377

378

379

380

381

382

383

384

386

388

389

390

391

392

393

394

395

396

397

For the insertion strategy, we categorize objects in images into three types—random, popular, and adversarial—each containing three objects. Random objects are sampled randomly, popular objects are the top frequent objects in the whole dataset, and adversarial objects are the top frequent objects with the current objects. By inserting one to three objects from each category into the correct captions with the assistance of GPT-4, we create examples with varying levels of hallucinations (*i.e.*, negative samples). For the removal strategy, we randomly select one or two segmented objects from the caption and instruct GPT-4 to eliminate them from the caption.

*Contrastive Learning*: We use these generated negative examples in a contrastive learning framework where CLIP is trained to correctly distinguish between the positive and negative pairs. By exposing the model to these fine-grained differences, CLIP becomes better at understanding nuanced visual features.

First, let I represent an image embedding and Ta text embedding. Let  $T^+$  be the text vector that correctly matches I, and let  $T^-$  denote a collection of negative texts not semantically aligned with I. We also introduce  $\beta$  as a temperature parameter.

The fundamental image-to-text contrastive objective can be expressed as:

$$\mathcal{L}_{i2t} = -\log(\frac{\exp(I \cdot T^{+}/\beta)}{\sum_{T^{*} \in \{T^{+}, T^{-}\}} \exp(I \cdot T^{*}/\beta)}).$$
(1)

The symmetric term  $\mathcal{L}_{t2i}$  can be constructed for text-to-image alignment. Combining them yields the image-text contrastive loss:

$$\mathcal{L}_{itc} = \frac{1}{2} (\mathcal{L}_{i2t} + \mathcal{L}_{t2i}). \tag{2}$$

Next, consider that we introduce an additional set of artificially generated negative texts  $\{T^{neg}\}$ . Incorporating these into the image-to-text objective

433

434

435

436

437

438

439

440

441

442

gives:

$$\mathcal{L}_{i2t} = -\log(\frac{\exp(I \cdot T^{+}/\beta)}{\sum_{T^{*} \in \{T^{+}, T^{-}, T^{neg}\}} \exp(I \cdot T^{*}/\beta)})$$
(3)

To further refine the separation between correct matches and all classes of negative samples (both standard and synthetic), we introduce a marginbased term. Let  $\tau_1$  be the margin threshold enforcing that a positive pair's similarity should exceed that of any negative pair by at least  $\tau_1$ :

$$\mathcal{L}_1 = \max(0, \tau_1 - (I \cdot T^+) + (I \cdot T^*)),$$
 (4)

where  $T^{\star} = \{T^{-}, T^{neg}\}$  is the union of standard and synthetic negatives.

Additionally, to encourage the model to distinguish synthetic negatives from standard negatives-thus capturing subtle semantic cues-we introduce another margin loss. Let  $\tau_2$  control the required margin between these two types of negative samples:

$$\mathcal{L}_2 = \max(0, \ \tau_2 - (I \cdot T^{neg}) + (I \cdot T^{-})). \ (5)$$

Finally, assigning weighting factors  $\lambda_1$  and  $\lambda_2$ to the margin terms allows adaptive emphasis on these constraints. The complete objective function is:

> $\mathcal{L} = \mathcal{L}_{itc} + \lambda_1 \mathcal{L}_1 + \lambda_2 \mathcal{L}_2.$ (6)

This integrated loss framework guides the model to better discriminate correct image-text pairs from both standard and refined negative samples.

Fine-grained perception-based visual instruction tuning As we mentioned, CLIP may not capture the finer details in the visual representation from the vision encoder. Therefore, we attempt to enable the LLM to perceive the fine-grained information within the CLIP vision encoder. Meanwhile, the method of enhancing CLIP and then replacing it is time-consuming, as it requires additional steps for feature alignment and instruction tuning after replacing the vision encoder of LVLMs. As a result, we explore a more efficient approach by directly enabling the LLM to perceive the detailed visual features during visual instruction tuning.

To achieve this, we propose fine-grained perception-based visual instruction tuning. Specifically, we randomly select two bounding boxes from the image, and then use the object attributes corresponding to these bounding boxes and their relationships to generate the corresponding captions. We then create instruction tuning data  $(I_f, T_f, R_f)$ , where  $T_f$  is the textual prompt: "Please caption the content in the bounding box",  $I_f$  is the image with bounding boxes, and  $R_f$  is the corresponding caption. This approach allows the model to perceive fine-grained information, such as region-level details, within the image.

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

# **3.2** How to reduce hallucination caused by the projector?

In Section 2.4, we reveal that hallucination introduced by the projector may be due to the inability of aligning visual and textual spaces, manifested by the low cosine similarity of caption embeddings and projected image features. Therefore, a straightforward remedy would be to explicitly bridge the image and caption representation during LLaVA's alignment stage.

#### 3.2.1 Loss Objectives

Besides autoregressive image-text generation loss  $\mathcal{L}_{itq} = -p(R|I,T)$ , we introduce an in-batch contrastive alignment loss  $\mathcal{L}_{itc}$  similar to Equation 2, where we maximize the similarity between a projected image feature and the corresponding text embedding for its caption. We only focus on the alignment stage and design three settings that involve the contrastive loss in different fashions.

**Integrated Alignment Loss** <sup>6</sup> The training process consists of two stages: alignment and visual instruction tuning. The contrastive loss is integrated to the *alignment* stage with a *learnable* (<sup>6</sup>) weight  $\lambda$ . The alignment objective is given by:  $\min_{\mathcal{P},\lambda} \mathcal{L}_{itg} + \lambda \mathcal{L}_{itc}$ . The visual instruction tuning stage is identical to LLaVA's.

Integrated Alignment Loss <sup>300</sup> All settings are the same as above except that the weight  $\lambda$  is *fixed* (<sup>30</sup>). The alignment objective is given by:  $\min_{\mathcal{P}} \mathcal{L}_{itg} + \lambda \mathcal{L}_{itc}.$ 

Separate Contrastive Alignment Loss We prepend a contrastive alignment stage solely for the projector  $\mathcal{P}$ . Namely, the first stage objective is given by:  $\min_{\mathcal{P}} \mathcal{L}_{itc}$ . The second stage and third stage correspond to the original autoregressive alignment and visual instruction tuning stage.

#### 4 **Results and Analysis**

We first introduce the benchmarks on which our methods to be evaluated, which are shown as follows. 1) Object-based benchmarks: testing the object perception of LVLMs. POPE and POPE-NoCaps (Liu et al., 2024b) are adopted, where the

Table 5: Performance of different methods across different benchmarks. The best results in each column are made **bold**. *w-ECLIP*: LLaVA with enhanced CLIP trained on fine-grained data; *w-FineIns*: LLaVA trained on fine-grained visual instruction tuning data.

	POPE						POPE-	NoCaps			Q	QA-VisualGenome				
Method	Random		Pop	ular	Adver	sarial	Ran	dom	Pop	ular	Adver	rsarial	Attr	ibute	Rela	tion
	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1
LLaVA-7B	85.40	86.36	86.63	85.25	85.13	83.88	84.80	82.97	79.40	78.30	74.77	74.69	64.67	66.60	67.57	74.81
w-ECLIP	87.80	86.87	87.30	86.04	85.87	84.70	85.27	83.50	81.00	79.69	75.77	75.46	67.67	68.79	67.00	74.11
w-FineIns	87.77	86.78	86.80	85.51	85.53	84.33	85.53	84.00	81.73	80.61	76.50	76.37	69.01	70.12	69.75	76.17

Table 6: Performance of different projector alignment methods across different benchmarks. The best results in each column are made **bold**. *Int. Align*.: Integrated Alignment Loss with trainable () / frozen() weighting parameter; *Sep. Ctrs. Align*.: Separate Contrastive Alignment Loss.

			РО	PE					POPE-	NoCaps			QA-VisualGenome			
Method	Random		Рор	ular	Adver	rsarial	Ran	dom	Pop	ular	Adver	sarial	Attr	ibute	Rela	ation
	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1
LLaVA-7B	87.42	86.36	86.63	85.25	85.13	83.88	84.80	82.97	79.40	78.30	74.77	74.69	64.67	66.60	67.57	74.81
Int. Align. 🔥	88.21	87.41	86.70	85.65	84.27	83.46	85.57	84.46	77.27	77.58	72.23	73.91	60.95	61.97	66.67	74.60
Int. Align. 🏶	88.04	87.20	86.67	85.56	84.50	83.60	84.90	83.28	79.37	78.47	74.57	74.76	63.84	65.21	66.73	74.26
Sep. Ctrs. Align.	88.56	87.86	87.33	86.38	84.57	83.88	85.57	84.24	80.07	79.42	75.13	75.54	64.26	64.77	69.60	76.06

latter is built on NoCaps (Agrawal et al., 2019) following a similar manner as in POPE. 2) Attributeand relation-based benchmark: *QA-VisualGenome* is constructed and adopted (detailed in Sec. 2.1). We provide an in-depth analysis of our methods for improving the vision encoder and the projector. We call object-, attribute- and relation-based benchmarks as perception-based benchmarks.

For a fair comparison, we only use the LLaVA-Caption dataset for alignment. All experiments are conducted on 4\*A100 GPUs. For the alignment stage, we set per-GPU batch size to 64, which is also the batch size contrastive alignment. We choose the well-known LLaVA-v1.5-7B model as our baseline. All three settings introduce no extra learnable parameters (except for the weighting parameter  $\lambda$  in **Integrated Alignment Loss** 6 setting). Under our setting, both the original and integrated alignment stage take 6 hours, and visual instruction tuning stage takes 24 hours. Notably, the prepended contrastive alignment stage takes only 12 minutes to train since only the vision encoder  $\mathcal{V}$ , projector  $\mathcal{P}$  and the embedding layer of LLM  $\mathcal{D}$  are involved in the forward process. For the two integrated loss settings, we empirically initialize  $\lambda$  with 5, make it learnable for 6 while keep it fixed for  $\overset{\text{w.}}{\approx}$ .  $\lambda_1$  and  $\lambda_2$  are set to 1.

Can our methods reduce hallucination caused by the vision encoder? Table 5 presents the comprehensive experimental results of various settings across different testing benchmarks. From this table, several key observations can be drawn: 1) Our proposed w-ECLIP method demonstrates superior performance compared to LLaVA-7B on perception-based benchmarks. This result underscores the effectiveness of our approach in reducing visual object hallucinations by enhancing the fine-grained perception capabilities of CLIP. 2) w-FineIns exhibits better performance than baseline on perception-based benchmarks. This finding suggests that our fine-grained instruction data can augment the fine-grained perception abilities of LLaVA by leveraging region-level captions during training. 3) Compared to w-FineIns, w-**ECLIP** demonstrates comparable or even better performance on perception-based benchmarks. Notably, w-FineIns offers efficiency advantages as it only requires the final training stage-instruction tuning-for the LVLM, simplifying the overall training process.

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

**Can our methods reduce hallucination caused by the projector?** We benchmark our methods in Table 6. For object-oriented benchmarks POPE and POPE-NoCaps, the model trained with *Separate Contrastive Alignment Loss* outperforms others on most splits of benchmarks, though the improvement over baseline seems marginal. For QA-VisualGenome benchmark, we only observe improvement on the "Relation" split with *Separate Contrastive Alignment Loss*, whereas slight performance drops are observed for others. These observations provide insights for the alignment process.

519

520

_	_	_
5	5	7
5	5	8
5	5	9
5	6	0
5	6	1
5	6	2
5	6	3
5	6	4
5	6	5
5	6	6
5	6	7
5	6	8
5	6	9
5	7	0
5	7	1
5	7	2
5	7	3
5	7	4
5	7	5
5	7	6
5	7	7
5	7	8
5	7	9
5	8	0
5	8	1
5	8	2
5	2 2	2
5	0 0	о Л
5	0	-4
5	0	ວ ດ
с -	ğ	0
5	8	7

590

591

554

555

556

#### Table 7: Performance on QA-FB15K.

Method	En En	tity	Relation		
	Acc	F1	Acc	F1	
LLaVA-7B	78.39	73.14	56.79	48.79	
Int. Align. 🔥	84.28	83.03	59.16	58.07	
Int. Align. 🌼	84.05	81.76	59.16	56.97	
Sep. Ctrs. Align.	83.94	81.65	59.39	57.41	
LLaVA-7B	78.39	73.14	56.79	48.70	
w-ECLIP	77.60	71.47	56.79	45.58	
w-FineIns	76.47	69.86	55.45	49.10	

Firstly, **object hallucinations may not be directly related to alignment in LVLM**, where vision encoder is mostly responsible for the perception process. Secondly, **perception-based attribute and relation hallucination can hardly be mitigated by contrastive training of projector**. Similar to object hallucination, better visual representations may be needed as a remedy.

Can our method influence other hallucinations? To further investigate the influence of our method on other kinds of hallucination, we introduced the Cognition-based benchmark: necessitating world knowledge in LVLMs for problem solving. We construct a cognition-based benchmark QA-FB15k based on the knowledge graph FB15K (Bordes et al., 2013). We show the results in Table 7. Contrastive alignment objective is beneficial for cognition-based knowledge, as evidenced by the performance boost on QA-FB15K. By better aligning between vision encoder and LLM, the LVLM is able to leverage the ability of LLM to answer the question that requires world knowledge, which is typically stored in LLMs pretrained on mountains of data. Nevertheless, performance boosts are found on QA-FB15K for all three settings over baselines. Neither w-FineIns nor w-ECLIP shows any improvement on the cognition-based benchmark. This may be attributed to the fact that, unlike perception-based benchmarks, cognition-based benchmarks necessitate not only the ability to identify objects but also the comprehension and application of relevant associated knowledge. The two methods primarily focus on improving perception, may not cater for the knowledge-intensive requirements of cognition-based benchmarks.

**More Analysis:** In addition, we add more experimental results on the hallucination benchmark and general benchmark, ablation study, and performance comparison with more baselines in Appendix B, E, F, and D.

## 5 Related Work

Large Vision-Language Model. The multimodal learning field has recently pivoted its focus towards Large Vision-Language Models (LVLMs) (Awadalla et al., 2023; Li et al., 2023a). Current advanced LVLMs primarily comprise three essential components: a language encoder, a visual encoder, and a cross-modal alignment mechanism (Rohrbach et al., 2018). To achieve comprehensive visual understanding, LVLMs generally undergo a series of training stages (Gong et al., 2023; Zhu et al., 2023; Liu et al., 2023d,e; Ye et al., 2023; Dai et al., 2023; Liu et al., 2023e). Despite significant advancements, LVLMs still face challenges with hallucination, which significantly affects performance across various multimodal applications.

594

596

597

598

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

Hallucinations in Large Vision-language Models. Since hallucination issues and mitigation techniques have been extensively explored in text generation (Ji et al., 2023; Min et al., 2023), research on hallucinations in LVLMs (Dai et al., 2023; Liu et al., 2023e; Jing and Du, 2024) attracts more attention. To evaluate the hallucination in the LVLMs, several researchers propose metrics and benchmarks (Rohrbach et al., 2018; Li et al., 2023e; Lovenia et al., 2023; Lu et al., 2023; Jing et al., 2024). Recently, various methods have been proposed to mitigate hallucinations in LVLMs, leveraging a range of techniques including decoding strategies (Leng et al., 2023; Huang et al., 2023), post-processing methods (Zhou et al., 2023; Chang et al., 2024; Yin et al., 2023), the development of higher-quality datasets (Liu et al., 2023c; Li et al., 2023d), and modality alignment(Li et al., 2023c; Yu et al., 2023; Zhou et al., 2024; Jing and Du, 2024; Sun et al., 2023a; Gunjal et al., 2023). Despite the success of the existing works, there lacks a comprehensive study of what causes visual hallucinations in LVLMs.

#### 6 Conclusion

Our study delves into the visual hallucination problem in LVLMs, identifying its sources within the model's components. By independently analyzing the LLM, vision backbone, and projector, we propose targeted mitigation strategies. We introduce fine-grained hallucination benchmarks, QA-VisualGenome and QA-FB15k, to comprehensively evaluate hallucinations. Our methods demonstrate effectiveness in reducing hallucinations, contributing to the reliability and accuracy of LVLMs.

746

747

748

749

750

751

752

753

698

699

700

# 4 Limitations

645Our work primarily focuses on analyzing and im-646proving hallucinations of general objects, such as647tables and people, while neglecting the research648topic of how to mitigate cognition-level hallucina-649tions, such as the names of individuals and famous650buildings.

#### References

651

653

654

660

666

670

671

672

673

674

675

677

679

696

- Harsh Agrawal, Peter Anderson, Karan Desai, Yufei Wang, Xinlei Chen, Rishabh Jain, Mark Johnson, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. nocaps: novel object captioning at scale. In 2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 November 2, 2019, pages 8947–8956. IEEE.
- Wenbin An, Feng Tian, Sicong Leng, Jiahao Nie, Haonan Lin, QianYing Wang, Ping Chen, Xiaoqin Zhang, and Shijian Lu. 2025. Mitigating object hallucinations in large vision-language models with assembly of global and local attention.
- Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. 2015. VQA: visual question answering. In *IEEE International Conference on Computer Vision*, pages 2425–2433. IEEE Computer Society.
- Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Yitzhak Gadre, Shiori Sagawa, Jenia Jitsev, Simon Kornblith, Pang Wei Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. 2023. Openflamingo: An open-source framework for training large autoregressive vision-language models. *CoRR*, abs/2308.01390.
- Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ B. Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri S. Chatterji, Annie S. Chen, Kathleen Creel, Jared Quincy Davis, Dorottya Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren Gillespie, Karan Goel, Noah D. Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani, Omar Khattab, Pang Wei Koh, Mark S. Krass, Ranjay Krishna, Rohith Kuditipudi, and et al. 2021. On the opportunities and risks of foundation models. CoRR, abs/2108.07258.
  - Antoine Bordes, Nicolas Usunier, Alberto García-Durán, Jason Weston, and Oksana Yakhnenko.

2013. Translating embeddings for modeling multirelational data. In Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States, pages 2787–2795.

- Tom B Brown. 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.
- Yue Chang, Liqiang Jing, Xiaopeng Zhang, and Yue Zhang. 2024. A unified hallucination mitigation framework for large vision-language models. *CoRR*, abs/2409.16494.
- Guiming Hardy Chen, Shunian Chen, Ruifei Zhang, Junying Chen, Xiangbo Wu, Zhiyi Zhang, Zhihong Chen, Jianquan Li, Xiang Wan, and Benyou Wang. 2024a. Allava: Harnessing gpt4v-synthesized data for a lite vision-language model. *arXiv preprint arXiv:2402.11684*.
- Liang Chen, Haozhe Zhao, Tianyu Liu, Shuai Bai, Junyang Lin, Chang Zhou, and Baobao Chang. 2024b. An image is worth 1/2 tokens after layer 2: Plug-andplay inference acceleration for large vision-language models. In *Computer Vision - ECCV 2024 - 18th European Conference, Milan, Italy, September 29-October 4, 2024, Proceedings, Part LXXXI*, volume 15139 of *Lecture Notes in Computer Science*, pages 19–35. Springer.
- Zhaorun Chen, Zhuokai Zhao, Hongyin Luo, Huaxiu Yao, Bo Li, and Jiawei Zhou. 2024c. Halc: Object hallucination reduction via adaptive focal-contrast decoding.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. vicuna: An opensource chatbot impressing gpt-4 with 90
- Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James R. Glass, and Pengcheng He. 2024. Dola: Decoding by contrasting layers improves factuality in large language models. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024.* OpenReview.net.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven C. H. Hoi. 2023. Instructblip: Towards general-purpose visionlanguage models with instruction tuning. *CoRR*, abs/2305.06500.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In 2009 IEEE conference on computer vision and pattern recognition, pages 248–255. Ieee.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias

859

860

861

863

864

865

809

810

811

Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An image is worth 16x16 words: Transformers for image recognition at scale. In ICLR. OpenReview.net.

754

755

763

764

770

778

781

787

789

790

791

795

797

798

804

807

- Tao Gong, Chengqi Lyu, Shilong Zhang, Yudong Wang, Miao Zheng, Qian Zhao, Kuikun Liu, Wenwei Zhang, Ping Luo, and Kai Chen. 2023. Multimodal-gpt: A vision and language model for dialogue with humans. CoRR, abs/2305.04790.
- Anisha Gunjal, Jihan Yin, and Erhan Bas. 2023. Detecting and preventing hallucinations in large vision language models. arXiv preprint arXiv:2308.06394.
- John Hewitt, Kawin Ethayarajh, Percy Liang, and Christopher D Manning. 2021. Conditional probing: measuring usable information beyond a baseline. arXiv preprint arXiv:2109.09234.
- Qidong Huang, Xiaoyi Dong, Pan Zhang, Bin Wang, Conghui He, Jiaqi Wang, Dahua Lin, Weiming Zhang, and Nenghai Yu. 2023. OPERA: alleviating hallucination in multi-modal large language models via over-trust penalty and retrospection-allocation. CoRR, abs/2311.17911.
- Qidong Huang, Xiaoyi Dong, Pan Zhang, Bin Wang, Conghui He, Jiaqi Wang, Dahua Lin, Weiming Zhang, and Nenghai Yu. 2024a. Opera: Alleviating hallucination in multi-modal large language models via over-trust penalty and retrospection-allocation.
- Qidong Huang, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Jiaqi Wang, Dahua Lin, Weiming Zhang, and Nenghai Yu. 2024b. Deciphering cross-modal alignment in large vision-language models with modality integration rate. arXiv preprint arXiv:2410.07167.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. ACM Computing Surveys, 55(12):1-38.
- Liqiang Jing and Xinya Du. 2024. FGAIF: aligning large vision-language models with fine-grained AI feedback. CoRR, abs/2404.05046.
- Liqiang Jing, Ruosen Li, Yunmo Chen, and Xinya Du. 2024. Faithscore: Fine-grained evaluations of hallucinations in large vision-language models.
- Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Li Fei-Fei. 2017. Visual genome: Connecting language and vision using crowdsourced dense image annotations. Int. J. Comput. Vis., 123(1):32-73.
- Alex Krizhevsky, Geoffrey Hinton, et al. 2009. Learning multiple layers of features from tiny images.
- Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-100 (canadian institute for advanced research).

- Sicong Leng, Hang Zhang, Guanzheng Chen, Xin Li, Shijian Lu, Chunyan Miao, and Lidong Bing. 2023. Mitigating object hallucinations in large visionlanguage models through visual contrastive decoding. CoRR, abs/2311.16922.
- Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. 2023a. Otter: A multi-modal model with in-context instruction tuning. CoRR, abs/2305.03726.
- Kenneth Li, Oam Patel, Fernanda B. Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023b. Inferencetime intervention: Eliciting truthful answers from a language model. 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.
- Lei Li, Zhihui Xie, Mukai Li, Shunian Chen, Peiyi Wang, Liang Chen, Yazheng Yang, Benyou Wang, and Lingpeng Kong. 2023c. Silkie: Preference distillation for large visual language models. CoRR, abs/2312.10665.
- Lei Li, Yuwei Yin, Shicheng Li, Liang Chen, Peiyi Wang, Shuhuai Ren, Mukai Li, Yazheng Yang, Jingjing Xu, Xu Sun, Lingpeng Kong, and Qi Liu. M<sup>3</sup>it: A large-scale dataset towards 2023d. multi-modal multilingual instruction tuning. CoRR, abs/2306.04387.
- Qing Li, Jiahui Geng, Derui Zhu, Zongxiong Chen, Kun Song, Lei Ma, and Fakhri Karray. 2025. Internal activation revision: Safeguarding vision language models without parameter update. In AAAI-25, Sponsored by the Association for the Advancement of Artificial Intelligence, February 25 - March 4, 2025. Philadelphia, PA, USA, pages 27428-27436. AAAI Press.
- Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. 2023e. 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.
- Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. 2014. Microsoft COCO: common objects in context. In ECCV, volume 8693 of Lecture Notes in Computer Science, pages 740-755. Springer.
- Fuxiao Liu, Tianrui Guan, Zongxia Li, Lichang Chen, Yaser Yacoob, Dinesh Manocha, and Tianyi Zhou. 2023a. Hallusionbench: You see what you think? or you think what you see? an image-context reasoning benchmark challenging for gpt-4v (ision), llava-1.5, and other multi-modality models. arXiv preprint arXiv:2310.14566.

Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser

Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser

Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser

Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae

Yufang Liu, Tao Ji, Changzhi Sun, Yuanbin Wu, and

Aimin Zhou. 2024b. Investigating and mitigating

object hallucinations in pretrained vision-language (CLIP) models. In *Proceedings of the 2024 Con-*

ference on Empirical Methods in Natural Language

Processing, EMNLP 2024, Miami, FL, USA, Novem-

ber 12-16, 2024, pages 18288-18301. Association

Holy Lovenia, Wenliang Dai, Samuel Cahyawijaya, Zi-

Jiaying Lu, Jinmeng Rao, Kezhen Chen, Xiaoyuan

Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike

Lewis, Wen-tau Yih, Pang Wei Koh, Mohit Iyyer,

Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023.

Factscore: Fine-grained atomic evaluation of fac-

tual precision in long form text generation. CoRR,

Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns,

Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu,

Chunyuan Li, Yikang Shen, Chuang Gan, Liang-Yan

Gui, Yu-Xiong Wang, Yiming Yang, Kurt Keutzer,

and Trevor Darrell. 2023a. Aligning large mul-

timodal models with factually augmented RLHF.

Trevor Darrell, and Kate Saenko. 2018. Object hal-

lucination in image captioning. In EMNLP, pages

Guo, Yawen Zhang, Baochen Sun, Carl Yang, and

Jie Yang. 2023. Evaluation and mitigation of agnosia in multimodal large language models. *CoRR*,

wei Ji, and Pascale Fung. 2023. Negative object

presence evaluation (nope) to measure object halluci-

for Computational Linguistics.

nation in vision-language models.

Lee. 2023e. Visual instruction tuning. CoRR,

tion tuning. arXiv preprint arXiv:2310.03744.

Lee. 2023d. Improved baselines with visual instruc-

Yacoob, and Lijuan Wang. 2024a. Mitigating hal-

lucination in large multi-modal models via robust

Yacoob, and Lijuan Wang. 2023c. Aligning large

multi-modal model with robust instruction tuning.

arXiv preprint arXiv:2306.14565.

CoRR, abs/2306.14565.

instruction tuning.

abs/2304.08485.

abs/2309.04041.

abs/2305.14251.

4035-4045. ACL.

CoRR, abs/2309.14525.

OpenAI. 2022. Chatgpt blog post.

Yacoob, and Lijuan Wang. 2023b. Aligning large

multi-modal model with robust instruction tuning.

- 875 876 877
- 878 879
- 88
- 88
- 88
- 8
- 88 88
- 8
- 89
- 05
- 89
- 894 895
- 89
- 89

901 902

903 904 905

906

907

908

908 909

910 911

912

913 914

914 915

916 917 Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Yikang Shen, Chuang Gan, Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, Kurt Keutzer, and Trevor Darrell. 2023b. Aligning large multimodal models with factually augmented RLHF. *CoRR*, abs/2309.14525. 918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models. *CoRR*, abs/2302.13971.
- Junyang Wang, Yuhang Wang, Guohai Xu, Jing Zhang, Yukai Gu, Haitao Jia, Ming Yan, Ji Zhang, and Jitao Sang. 2023. An Ilm-free multi-dimensional benchmark for mllms hallucination evaluation. *CoRR*, abs/2311.07397.
- Tan Wang, Jianqiang Huang, Hanwang Zhang, and Qianru Sun. 2020. Visual commonsense R-CNN. In 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020, pages 10757–10767. Computer Vision Foundation / IEEE.
- Jinfeng Wei and Xiaofeng Zhang. 2024. DOPRA: decoding over-accumulation penalization and reallocation in specific weighting layer. In *Proceedings* of the 32nd ACM International Conference on Multimedia, MM 2024, Melbourne, VIC, Australia, 28 October 2024 - 1 November 2024, pages 7065–7074. ACM.
- Yun Xing, Yiheng Li, Ivan Laptev, and Shijian Lu. 2024. Mitigating object hallucination via concentric causal attention. In Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024.
- Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, Chenliang Li, Yuanhong Xu, Hehong Chen, Junfeng Tian, Qian Qi, Ji Zhang, and Fei Huang. 2023. mplug-owl: Modularization empowers large language models with multimodality. *CoRR*, abs/2304.14178.
- Shukang Yin, Chaoyou Fu, Sirui Zhao, Tong Xu, Hao Wang, Dianbo Sui, Yunhang Shen, Ke Li, Xing Sun, and Enhong Chen. 2023. Woodpecker: Hallucination correction for multimodal large language models. *CoRR*, abs/2310.16045.
- Tianyu Yu, Yuan Yao, Haoye Zhang, Taiwen He, Yifeng Han, Ganqu Cui, Jinyi Hu, Zhiyuan Liu, Hai-Tao Zheng, Maosong Sun, and Tat-Seng Chua. 2023. RLHF-V: towards trustworthy mllms via behavior alignment from fine-grained correctional human feedback. *CoRR*, abs/2312.00849.
- Zihao Yue, Liang Zhang, and Qin Jin. 2024. Less is more: Mitigating multimodal hallucination from an

EOS decision perspective. 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 11766–11781. Association for Computational Linguistics.

974

975

976

978

987

991

993

994

996

997

1000

1001

1002

1004

1006

1007 1008

- Xiaofeng Zhang, Yihao Quan, Chaochen Gu, Chen Shen, Xiaosong Yuan, Shaotian Yan, Hao Cheng, Kaijie Wu, and Jieping Ye. 2024a. Seeing clearly by layer two: Enhancing attention heads to alleviate hallucination in lvlms. *CoRR*, abs/2411.09968.
- Yue Zhang, Liqiang Jing, and Vibhav Gogate. 2025. Defeasible visual entailment: Benchmark, evaluator, and reward-driven optimization. In AAAI-25, Sponsored by the Association for the Advancement of Artificial Intelligence, February 25 - March 4, 2025, Philadelphia, PA, USA, pages 25976–25984. AAAI Press.
  - Yue Zhang, Jingxuan Zuo, and Liqiang Jing. 2024b. Fine-grained and explainable factuality evaluation for multimodal summarization. *arXiv preprint arXiv:2402.11414*.
  - Yiyang Zhou, Chenhang Cui, Rafael Rafailov, Chelsea Finn, and Huaxiu Yao. 2024. Aligning modalities in vision large language models via preference finetuning. *CoRR*, abs/2402.11411.
  - Yiyang Zhou, Chenhang Cui, Jaehong Yoon, Linjun Zhang, Zhun Deng, Chelsea Finn, Mohit Bansal, and Huaxiu Yao. 2023. Analyzing and mitigating object hallucination in large vision-language models. *CoRR*, abs/2310.00754.
  - Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *CoRR*, abs/2304.10592.

# A Hallucinations in Different Components

We show the potential hallucinations of each component of LVLMs, and the corresponding mitigation methods in Table 8101110121013

1010

1014

1015

1019

1024

1025

# B More Experiments on Hallucination Benchmark

We further add experiments on another hallucina-<br/>tion benchmark, Amber. The experimental results1016of Table 9 show the effectiveness of our method.1018

## C Case Study

We showed some hallucinated examples in Figure10203. We can see that the hallucination caused by1021CLIP can be further input to the LVLM, causing1022the hallucination in the LVLM.1023



Question: Is there a bed in the image?

LLaVA output: There is a bed in the image.

Text 1: a bird is in the image. CLIP probability: 0.1848

Text 2: a bed is in the image. CLIP probability: 0.8154

Figure 3: The illustration of the hallucinated case for CLIP and LLaVA.

# D Comparison with the Existing Hallucination Mitigation Method

To verify the effectiveness or our methods, we further add more baselines on POPE, as shown in 1027

Component	Hallucination?	Mitigation
Vision Backbone	$\checkmark$	w-ECLIP & w-FineIns
Projector	$\checkmark$	Int. Align. 🔥 & Int. Align. 🕸 & Sep. Ctrs. Align.
LLM	×	N/A

Table 8: Illustration of potential hallucinations in the components of LVLMs, and the corresponding mitigation methods

1028Table 10. From this table, our methods show com-<br/>petitive performance with the best baseline (i.e.,<br/>10301030Less is more). This further demonstrates the effec-<br/>tiveness of our method.

# E Experiment on General Benchmark

To verify the impact of the proposed method on general capabilities, we further conduct experiments on the general benchmark LLaVA-Bench (Liu et al., 2023e). The results of Table 11 show the effectiveness of our method.

#### F Ablation Study

1032

1033

1034

1035

1036

1037

1038

1039In this section, we conduct ablation experiments to1040assess the contribution of each component in the1041loss function by individually removing the weights1042 $\lambda_1$  and  $\lambda_2$ . The results are shown in the Table 12.1043These results demonstrate that both components1044play meaningful roles in enhancing model perfor-1045mance.

Dataset	LLaVA-7B	w-ECLIP	w-FineIns	Int. Align. 🔥	Int. Align. 🏶	Sep. Ctrs. Align.
Existence	83	93	92	88	91	87
Attribute	64	81	81	75	78	76
Relation	65	69	70	57	62	59
All	71	73	81	73	77	74

Table 9: Performance on the Amber dataset across different model variants. Bold indicates best scores per row.

Method	F1 Score
DoLa (Chuang et al., 2024)	80.2
ITT (Li et al., 2023b)	83.7
VCD (Leng et al., 2023)	83.2
AGLA (An et al., 2025)	84.6
OPERA (Huang et al., 2024a)	85.2
DOPRA (Wei and Zhang, 2024)	85.6
HALC (Chen et al., 2024c)	83.9
FastV (Chen et al., 2024b)	81.3
Less is more (Yue et al., 2024)	86.0
CCA-LLAVA (Xing et al., 2024)	85.5
LRV (Liu et al., 2024a)	80.0
Amber (Wang et al., 2023)	81.6
EAH (Zhang et al., 2024a)	85.7
w-ECLIP	85.9
w-FineIns	85.5
Int. Align. 🔥	85.5
Int. Align. 🏶	85.5
Sep. Ctrs. Align	86.0

Table 10: POPE F1 scores for baselines and proposed methods. Bold indicates the highest score.

Model	Conv	Detail	Complex	Full
LLaVA-7B	92	75	75	81
w-ECLIP	93	84	87	88
w-FineIns	94	86	86	89
Int. Align. 🔥	95	87	83	89
Int. Align. 🍀	93	84	82	86
Sep. Ctrs. Align	99	85	87	90

Table 11: Model performance comparison on different categories and the full set on LLaVA-Bench.

Method	POPE					
	Random		Popular		Adversarial	
	Acc	F1	Acc	F1	Acc	F1
w-ECLIP	87.80	86.87	87.30	86.04	85.87	84.70
$\lambda_1 = 0$	87.50	86.38	86.93	85.84	85.62	83.97
$\lambda_2 = 0$	87.52	86.47	86.79	85.88	85.47	84.11

Table 12: Ablation study on the impact of loss function components  $\lambda_1$  and  $\lambda_2$  across different POPE test subsets.