Beyond Single Models: Mitigating Multimodal Hallucinations via Adaptive Token Ensemble Decoding

Anonymous ACL submission

Abstract

Large Vision-Language Models (LVLMs) have recently achieved impressive results in multimodal tasks such as image captioning and visual question answering. However, they remain prone to object hallucination-generating descriptions of nonexistent or misidentified objects. Prior work have partially mitigated this via auxiliary training objectives or external modules, but often lacks scalability, adaptability, or model independence. To address these limitations, we propose Adaptive Token **Ensemble Decoding (ATED)**, a training-free, token-level ensemble framework that mitigates hallucination by aggregating predictions from multiple LVLMs during inference. ATED dynamically computes uncertainty-based weights for each model, reflecting their reliability at each decoding step. It also integrates diverse decoding paths to improve contextual grounding and semantic consistency. Experiments on standard hallucination detection benchmarks demonstrate that ATED significantly outperforms state-of-the-art methods, reducing hallucination without compromising fluency or relevance. Our findings highlight the benefits of adaptive ensembling and point to a promising direction for improving LVLM robustness in high-stakes vision-language applications. Code is available at https:// anonymous.4open.science/r/ATED.

1 Introduction

011

017

In recent years, large language models (LLMs) have made significant breakthroughs in natural language processing (Touvron et al., 2023; Chiang et al., 2023; Achiam et al., 2023; Bai et al., 2023a) and have been increasingly extended to visionlanguage tasks, giving rise to large vision-language models (LVLMs) (Ye et al., 2023; Liu et al., 2023; Li et al., 2023a, 2024; Chen et al., 2024c,d; Bai et al., 2023b). These models have demonstrated strong capabilities in both understanding (Zhang



Figure 1: Comparison of image-description generation results from various LVLMs and our proposed *ATED* method. <u>Red</u> text indicates hallucination. <u>Green</u> text represents hallucination mitigating from *ATED*.

et al., 2025; Lai et al., 2024) and generating (Geng et al., 2023) multimodal content.

However, LVLMs often suffer from the problem of *object hallucination*, where the model generates details or objects that do not exist in the image (Li et al., 2023d; Wang et al., 2023; Gunjal et al., 2024; Liu et al., 2024b), significantly limiting their reliability in high-stakes applications, such as autonomous driving, medical image analysis, and remote sensing, where factual correctness and visual grounding are critical.

Early research on mitigating hallucinations primarily focused on enhancing data quality and training paradigms. Specifically, diverse instructiontuning datasets and multi-task training approaches were introduced to reduce the models' tendency to hallucinate during generation (Li et al., 2023c; Liu et al., 2024a). Other methods adopted posthoc strategies by implementing output-checking mechanisms to detect and correct hallucinated con-

077

084

101

102

103

104

105

106

107

108

109

110

111

112

113

063

tent (Yin et al., 2024; Zhou et al., 2024).

More recently, training-free approaches have emerged (Leng et al., 2023; Wang et al., 2024; Huang et al., 2024). However, many of these methods rely on additional annotations, large-scale finetuning, or complex inference, incurring substantial human and computational costs. Furthermore, single-model strategies are inherently limited by the knowledge scope of the underlying model, restricting generalization and adaptability. As shown in Figure 1, existing LVLMs exhibit notable levels of hallucination in image captioning tasks, highlighting the need for a more robust solution.

Ensemble learning (Polikar, 2012), which leverages the collective intelligence of multiple models, has proven highly effective at reducing errors and enhancing robustness in traditional classification and regression tasks (Mienye and Sun, 2022). More recently, it has been successfully extended to text generation tasks-particularly in LLMs-to improve output accuracy and mitigate issues such as hallucination (Jiang et al., 2023; Wan et al., 2024). Inspired by these advances, this paper presents a novel framework that integrates ensemble learning with diversified decoding strategies within the autoregressive generation process, harnessing the complementary strengths of different LVLMs. By aggregating outputs from multiple models or decoding paths, our approach not only improves the factual consistency and coherence of generated content, but also reduces the disparity in hallucination tendencies across models-ultimately enhancing the generalization and adaptability of LVLMs in multimodal tasks.

To this end, we introduce **Adaptive Token Ensemble Decoding (ATED)**—the first fine-grained, token-level ensemble strategy for multimodal LVLMs. ATED is a training-free framework that performs parallel inference across multiple LVLMs and aggregates their output logits at the token level. We estimate each model's hallucination tendency via output uncertainty, and use a greedy optimization algorithm to derive adaptive importance weights by minimizing overall uncertainty. Furthermore, ATED incorporates diverse decoding paths, substantially improving the factual accuracy and reliability of generated outputs while maintaining strong adaptability across diverse scenarios.

Our main contributions are as follows:

• We propose **ATED**, a training-free multimodal ensemble decoding method that mitigates hallucinations via fine-grained tokenlevel fusion.

114

115

116

117

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

- We introduce an uncertainty-minimization weighting mechanism that dynamically assigns weights based on model confidence, improving the reliability of ensemble decoding.
- Extensive experiments show that ATED consistently outperforms existing methods across multiple multimodal benchmarks, achieving superior accuracy and robustness.

2 Related Work

Hallucination in LVLMs Hallucination was initially observed in LLMs, referring to generated content that deviates from factual knowledge or user intent (Jing et al., 2024; Liu et al., 2024b). Large vision-language models (LVLMs) (Bai et al., 2025; Zhang et al., 2023), which extend LLMs with visual inputs, also exhibit hallucinations-typically manifesting as mismatches between generated text and visual content. Existing studies categorize hallucinations in LVLMs into three main types: *object* hallucination (Biten et al., 2021; Li et al., 2023d; Rohrbach et al., 2019), attribute hallucination, and relationship hallucination (Wu et al., 2024; Zhou et al., 2024). Object hallucination refers to fabricated or omitted objects; attribute hallucination involves incorrect properties such as color or size; relationship hallucination describes inaccurate relations among objects. These errors may arise from visual misinterpretation, flawed reasoning, or overreliance on language priors.

Hallucination Mitigation in LVLMs To address hallucination in LLMs and LVLMs, researchers have proposed a range of solutions, including improved instruction tuning (Jiang et al., 2024; Liu et al., 2024a; Yu et al., 2024a; Yue et al., 2024), reinforcement learning with human or AI feedback (Gunjal et al., 2024; Kim et al., 2024; Li et al., 2023b; Sun et al., 2023; Yu et al., 2024b,c), retrieval augmentation, and structural model enhancements (Zhai et al., 2024). More recently, several training-free decoding strategies have been developed to suppress hallucination in LVLMs. For example, conservative decoding methods applied to both original and perturbed inputs (Chen et al., 2024b; Favero et al., 2024; Huo et al., 2025; Leng et al., 2023; Wang et al., 2024; Woo et al., 2024) aim to reduce overreliance on language priors. Techniques such as input distortion-applied

to either visual content or instructions—amplify
hallucinations to better identify and suppress them
through contrastive decoding. Token-level pruning and related approaches (Favero et al., 2024;
Woo et al., 2024) also manipulate visual inputs to
mitigate hallucinations.

While existing methods mitigate hallucinations to some extent from various angles, they often suffer from limitations in scalability, generalizability, and practical deployment. This paper, therefore, focuses on reducing hallucinations in LVLMs without requiring additional training or complex reasoning pipelines, with an emphasis on adaptability to real-world applications.

3 Methodology

169

170

171

172

173

174

176

177

178

179

180

181

182

183

185

187

189

190

191

193

194

195

198

199

200

206 207

210

3.1 Preliminaries of LVLMs Generation

The generation mechanism of LVLMs can be deconstructed into three core modules: Vision Language Input, Model Forward Propagation, and Next Token Decoding.

Vision Language Input. LVLMs take both visual and textual inputs. Typically, raw images are processed by a vision encoder (e.g., a pretrained visual backbone), and the resulting features are projected into the input space of the language decoder via a cross-modal interface. These visual features are represented as visual tokens $V = \{v_1, v_2, ..., v_n\}$, where *n* is the number of visual tokens. Similarly, the textual input is tokenized into text tokens $T = \{t_1, t_2, ..., t_m\}$ using a tokenizer, where *m* is the number of text tokens. The visual and text tokens are then concatenated to form the final input sequence, denoted as $\{x_i, i \in [0, n + m - 1]\}$.

Model Forward Propagation. A large language model (LLM) parameterized by ϕ , such as Vicuna (Chiang et al., 2023), generates responses by conditioning on both the text and visual context. Following the autoregressive generation paradigm, LVLMs predict the probability of the next token x_t at time step t based on the previously generated tokens, the input text, and the visual features, over the vocabulary set V. This process can be formally expressed as:

$$p(x_t \mid v, t, x_{< t}) = \operatorname{softmax}(\operatorname{logit}_{\phi}(x_t \mid v, t, x_{< t})),$$
$$x_t \in v \quad (1)$$

where x_t denotes the token at time step t, $x_{< t}$ represents the sequence of previously generated tokens.

Next Token Decoding. Based on the predicted probabilities $p(x_t|v, t, x_{\leq t})$, various decoding strategies—such as greedy decoding, beam search, and contrastive decoding (e.g., VCD)-can be applied to generate output. While these strategies can marginally reduce hallucinations, they are typically restricted to single-model outputs, cannot leverage external knowledge, and fail to fully exploit complementary strengths across different models. As a result, they remain prone to errors, especially in open-domain scenarios. In contrast, our method adaptively fuses the token-level logits from multiple LVLMs that share the same vocabulary, immediately after the forward pass. By leveraging the diverse capabilities of different models, our approach more effectively mitigates hallucinations in both general-purpose and task-specific settings.

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

250

254

255

256

257

3.2 Adapative Token Ensemble Decoding

To leverage the complementary strengths of diverse LVLMs and enhance general task performance while reducing hallucinations, we propose a tokenlevel ensemble decoding approach. Specifically, we introduce **Adaptive Token Ensemble Decoding (ATED)**, a training-free, uncertainty-guided fusion method that dynamically integrates multiple LVLMs at inference time. The overall framework is illustrated in Figure 2.

Given two or more LVLMs $\{M_i\}$ at test time, ATED fuses their output logits using adaptive importance weights $\{\lambda_1, \ldots, \lambda_M\}$, where each weight λ_i reflects how well model M_i interprets the visual and textual inputs. At each decoding step t, each model M_i takes the visual token sequence $V = \{v_1, v_2, \ldots, v_n\}$ and the text token history $T = \{t_1, t_2, \ldots, t_{t-1}\}$ to generate a logit score $p_i = p_i(x_t \mid V, T)$ over the vocabulary. Assuming all models share the same vocabulary, the logits from the *i*-th model are computed as:

$$p(x_t|v, t, x_{< t})^{\operatorname{orig}, i} = LVLMs_i(X_{< t})\big|_t, \quad (2)$$

where i = 1, ..., M indexes the models, and $p \in R^v$ with v vocabulary size. ATED combines output logits to obtain the final probabilities as:

$$p(x_t|v, t, x_{< t}) = \left(\sum_{i=1}^{M} \underbrace{\lambda_i}_{unknown} p_i\right), \quad (3)$$

where p_i is the decoding logits of model M_i . We assume that the weights $\lambda_i \in [0, 1]$ are normalized, i.e., $\sum_{i=1}^{M} \lambda_i = 1$. To further enhance robustness, we additionally introduce multiple perturbed



Figure 2: Overall pipeline of Adapative Token Ensemble Decoding (ATED). Given a set of candidate LVLMs, a system instruction (inst), the original image, and its perturbed variants, the system produces multiple output streams. We ensemble the logits from each LVLM using uncertainty-guided weights and employ greedy uncertainty optimization to generate the final output. The entire process is dynamically repeated at each time step t of token generation, ensuring high reliability and consistency of the generated results.

variants of the original image *I* and employ contrastive decoding to effectively mitigate the backbone model's hallucinations. Formally, contrastive decoding can be expressed as:

261

265

269

270

274

275

279

281

$$p_{i,t} = \operatorname{softmax} \left[(1+\alpha) \operatorname{logit}_{\phi} (x_t \mid v, t, x_{< t}) - \alpha \operatorname{logit}_{\phi} (x_t \mid v_d, t, x_{< t}) \right],$$
(4)

where d and α indicate distortion operation and vision contrastive hyperparameter, respectively.

Inspired by Chen et al. (2024a); Qiu et al. (2025); Dey et al. (2025), we propose utilizing the entropy of probability distributions based on visual features as an uncertainty metric when LVLMs generate the next textual token. This metric reflects the model's token prediction confidence and associated weights λ_i under current multimodal inputs. The metric aligns with the training objectives of causal language modeling. By analyzing prediction entropy under visual conditions, we can evaluate the model's depth of understanding of specific visual content, the quality of vision-language alignment, and the distributional differences between the visual input and the model's training data (Gonen et al., 2024). This cross-modal uncertainty analysis provides a novel perspective for assessing the generalization capabilities of multimodal large models in open-world scenarios.

3.3 Uncertainty-Guided Weight

Uncertainty Minimization. Given a tokenized input X at time t, its uncertainty score is defined

as the distribution entropy of X at that time, with the following formula as:

$$H_{i,t}(X) = -\sum_{x_t \in \mathcal{V}} P_{i,t} \log P_{i,t}, \qquad (5)$$

287

290

291

292

293

295

296

297

300

301

302

303

304

305

306

307

308

309

310

311

312

313

where $P_{i,t}$ represents the normalized probability for the *i*-th token corresponding to vocabulary V, conditioned on the preceding tokens $x_{< t}$ according to model M_i .

We formulate the assignment of importance weights across M models as an optimization problem that requires no training or labeled data, and can be solved directly during next token prediction. Formally, our optimization framework is defined as follows:

$$\lambda_1^*, \dots, \lambda_i^* = \arg \min_{\lambda_1, \dots, \lambda_i} - \sum p(x_i | x_{< i}) \log p(x_i | x_{< i}),$$

$$p(x_i|x_{< i}) = \operatorname{softmax} \sum_{i=0}^{M} \lambda_i p_i(x_t|v, t, x_{< t}),$$
(6)

where the weights λ_i are inversely proportional to each model's normalized uncertainty score—i.e., models with lower uncertainty are assigned higher weights. All weights are constrained such that $\sum_{i=1}^{M} \lambda_i = 1$ and $\lambda_i \in [0, 1]$.

Uncertainty Greedy Optimization. To address the uncertainty minimization problem proposed in Equation 6, we introduce an efficient greedy optimization algorithm that incrementally ensembles LVLMs. Specifically, we first compute the uncertainty of each LVLM's next-token prediction using Equation 5, and then sort the LVLMs M_1, \ldots, M_i

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

358

based on their uncertainties scores. Let the sortedmodels be denoted as:

316

317

319

321

323

324

325

327

328

334

341

344

347

$$[M_1^*, M_2^*, \dots, M_i^*] = \operatorname{argsort}(H_{1,t}, \dots, H_{i,t}),$$
(7)

where $[M_1^*, M_2^*, \ldots, M_i^*]$ are ordered by the lowest to highest uncertainty scores, and set the weight of the top-ranked model to $\lambda_1^* = 1$ and $\lambda_{i>1}^* =$ $0, i = \{2, \ldots, M\}.$

We then iteratively consider incorporating the next-ranked model M_{i+1}^* into the current ensemble. A grid search is performed over interpolation weights λ_i with the step values between the two models, and the fused logits are defined as:

$$p_{\lambda}(x_{i} \mid x_{< i}) = \lambda p^{(i)}(x_{i} \mid x_{< i}) + (1 - \lambda) p^{(i+1)}(x_{i} \mid x_{< i}).$$
(8)

Finally, we compute the uncertainty score under each interpolation ratio with Equation 6, and select the weight that minimizes the uncertainty score. The same procedure is repeated to iterate through all LVLM for the next step and iteratively update the ensemble logits and weight allocation accordingly. We can perform early stopping when we find $\lambda(i) = 1$, i.e., the effect of the current LVLM is zero. The final generation is performed using the output probabilities from the adapative uncertainty-guided ensemble.

4 Experiments

4.1 Experimental Settings

4.1.1 Datasets

We evaluate our model on three datasets as listed below. More details are shown in the Appendix C.

POPE (Probability of Object Presence Estimation). Li et al. (2023d) is a benchmark dataset for evaluating object hallucination in LVLMs. POPE integrates the MSCOCO (Lin et al., 2015), A - OKVQA (Schwenk et al., 2022), and GQA (Hudson and Manning, 2019) datasets to form 27,000 query - answer pairs for evaluation. Performance is quantified using standard metrics, including *accuracy, precision, recall,* and *F1 score.*

352CHAIR(The Caption Hallucination Assessment353with Image Relevance). Rohrbach et al. (2019)354is a dataset for evaluating object hallucination in355image captioning. CHAIR has two main variants:356 $CHAIR_i$ and $CHAIR_s$, focusing on instance and357sentence levels, respectively.

MME (Multimodal Large Language Model Evaluation). Fu et al. (2024) assesses the performance of LVLMs in terms of two core capabilities: perception and cognition. In our evaluation, we focus on four representative sub-tasks: object existence, counting, position, and color. Model performance is measured using the *accuracy*+ metric.

4.1.2 Models

We integrate our proposed method with four popular LVLMs: InstructBLIP (Dai et al., 2023), MiniGPT-4 (Zhu et al., 2023), LLaVA-1.5 (Liu et al., 2024c), and LLaVA-Next (Liu et al., 2024d). All the LVLMs used have a language model size of 7 billion parameters (7B). InstructBLIP and MiniGPT-4 utilize a Q-former(Li et al., 2023a), which represents an image using only 32 tokens, effectively bridging the visual and textual modalities. LLaVA-1.5 and LLaVA-NeXT employ a linear projection layer to align features from the two modalities. All LVLMs adopt pretrained vision encoders such as the CLIP vision encoder(Radford et al., 2021), along with pre-trained large language models (LLMs) as language decoders, such as LLaMA(Touvron et al., 2023) or Vicuna v1.1(Chiang et al., 2023). Complete experimental details are provided in the Appendix. A.1

4.1.3 Baselines

For the object hallucination evaluation, we employ several widely used decoding strategies, such as multinomial sampling (default), greedy decoding, and four state-of-the-art training-free decoding methods. Greedy decoding selects tokens step by step by always choosing the one with the highest probability from the language model's logits. Based on greedy decoding, beam search maintains a set of beams to explore a wider range of candidates and eventually selects the best one among them. OPERA(Huang et al., 2024) is an improved method based on beam search, it alleviates hallucination by penalizing specific patterns of knowledge aggregation. VCD(Leng et al., 2023) reduces hallucination by decoding with noisy images in a contrastive manner. ICD(Wang et al., 2024) mitigates hallucination by designing negative prompts to interfere with the visual inputs during contrastive decoding. SID (Huo et al., 2025) mitigateas hallucinations by introspectively filtering low-relevance visual signals during generation. For all the baselines, we use the default hyperparameters provided by their original source code to ensure a fair com-

Model	Method	Random		Popular		Adversarial	
		Accuracy	F1 Score	Accuracy	F1 Score	Accuracy	F1 Score
	default	83.86	82.68	80.82	79.54	76.42	76.61
	OPERA	88.85	88.67	82.77	83.40	79.16	80.93
LLaVA-1.5	VCD	87.2	87.17	83.08	83.07	77.70	79.14
	ICD	83.15	83.91	83.15	83.91	79.13	80.41
	SID	89.46	89.62	85.13	85.94	83.24	82.21
	default	81.44	81.21	79.06	79.12	76.29	76.99
	OPERA	84.57	83.74	78.24	79.15	74.59	76.33
InstructBLIP	VCD	84.91	84.08	81.89	81.46	79.97	79.90
	ICD	81.12	82.25	81.12	82.25	76.82	78.99
	SID	87.23	86.90	81.16	82.57	78.51	81.26
	default	65.65	66.45	59.61	62.54	58.35	62.22
	OPERA	79.91	77.6	73.78	72.23	71.76	70.64
MiniGPT4	VCD	67.79	68.54	62.42	65.24	60.17	63.94
	ICD	71.89	75.63	64.58	75.33	61.77	67.61
	SID	75.20	76.12	68.94	72.93	66.57	69.40
	default	84.83	81.78	81.00	79.72	76.01	75.83
	OPERA	88.41	87.33	82.69	83.48	79.22	79.40
LLaVA-NeXT	VCD	86.01	85.20	81.90	82.23	78.00	79.12
	ICD	82.14	82.09	81.95	81.87	79.24	78.89
	SID	89.54	89.67	85.24	85.67	82.43	81.51
	ATED*	88.74	87.82	83.62	84.82	78.86	81.21
Ensemble	ATED&	89.21	89.39	85.32	85.66	81.51	82.32
	ATED#	89.83	89.35	86.71	85.97	82.96	82.78

Table 1: **POPE evaluation results on different decoding strategies.** Results are from the papers or re-implemented based on official codes. *Note:** denotes ours ensemble method without vision contrastive decoding, & denotes ensemble with LLaVA-1.5 and InstructBLIP, # denotes ensemble with LLaVA-1.5, InstructBLIP and LLaVA-NeXT.

parison. We posit that our method, being LVLMagnostic, can be easily integrated into various offthe-shelf LVLMs that share the same vocabulary.

4.2 Experimental Results

408

409

410

411

Results on POPE. We begin with the most 412 widely adopted benchmark for evaluating object 413 hallucination. Table 1 reports the average perfor-414 mance across three evaluation settings-random, 415 popular, and adversarial-on various datasets, 416 where *Default* refers to the unmodified backbone 417 model. Our evaluation of ATED includes three 418 configurations: two distinct LVLM ensemble vari-419 ants (ATED[&] and ATED[#]), as well as a ver-420 sion (ATED^{*}) based on the ATED[#] that ex-421 cludes vision-contrastive decoding (ATED*). For 422 clarity, we highlight the best baseline results for 423 424 each backbone in bold. Compared to each respective backbone, ATED achieves improvements of 425 4.20%-6.29% in Accuracy and 6.29%-6.97% in 426 F1-score. Furthermore, on both LLaVA-1.5 and 427 InstructBLIP, ATED consistently surpasses state-of-428

the-art methods ICD and VCD, attaining additional gains ranging from 0.89% to 5.10% in *Accuracy* and 0.80% to 2.94% in F1-*score*, thereby effectively mitigating hallucination issues.

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

Results on CHAIR. Beyond the binary "yes" or "no" evaluations on the POPE benchmark, we further validate the effectiveness of TADE in openended image captioning using the CHAIR metric. Specifically, we randomly sample 500 images from the validation split of the MSCOCO dataset and query various LVLMs with the prompt, "Please describe this image in detail." As shown in Table, when setting the maximum new token length to 64, our proposed ATED method significantly outperforms all baseline decoding approaches on the CHAIR metric, achieving improvements of 81.57% and 1.23% over the strongest baseline, respectively. Notably, when increasing the generation length to 512 tokens, ATED still attains the best performance on the CHAIR metric, with an improvement of approximately 30.0%. More detailed result are

Туре	METHOD	LLaVA-1.5		InstructBLIP CHAIR _S CHAIR _I		MiniGPT4 CHAIR _S CHAIR _I		LLaVA-NeXT CHAIR _S CHAIR _I	
	Default	24.8	8.9	30.3	13.9	19.8	8.5	24.26	8.51
	OPERA	21.8	8.2	28.4	9.7	22.6	8.2	21.33	7.73
Single	VCD	23.6	8.6	30.0	11.2	22.0	10.6	23.27	8.34
	ICD	21.0	8.7	21.8	8.2	20.0	8.7	20.59	8.54
	SID	20.7	8.4	20.7	8.4	23.1	10.7	19.37	7.83
Ensemble	Ours		AT	ED ^{&}			AT	E D [#]	
		CHA	IR_S	CHA	AIR _I	CHA	IR_S	CHA	IR _I
		15.3		10.9		11.4		8.1	

Table 2: **CHAIR evaluation results on different decoding strategies**. Results are from the papers or re-implemented based on official codes. *Note:* & denotes ensemble with LLaVA-1.5 and InstructBLIP, # denotes ensemble with LLaVA-1.5, InstructBLIP and LLaVA-NeXT.



Figure 3: MME evaluation on hallucination subset with decoding strategies in LLaVA-1.5(up) and Instruct-BLIP(down).

provided in Appendix E.

450

Results on MME. We extend the evaluation to 451 include hallucinations at the object attribute level. 452 We further conducted a systematic and comprehen-453 sive evaluation of the ATED method on the MME 454 hallucination subset, which includes object-level 455 tasks (existence identification and quantity judg-456 ment) and attribute-level tasks (location identifica-457 tion and color classification). As shown in Figure 458 4, ATED achieves the highest performance on loca-459 tion questions, and attains nearly perfect accuracy 460 461 on existence-related questions. Our ATED method significantly outperforms the default LVLMs and 462 other baseline methods across all four tasks (with 463 total score improvements of at least +61.7 and 464 +54.2 respectively). 465

Setting	Method	Accuracy	F1 Score	
	Uniform	87.03	86.65	
Dandam	Confidence-based	88.13	86.76	
Kandom	Ours (w/o UGO)	89.30	87.78	
	Ours	88.97	88.29	
	Uniform	84.57	84.37	
Descales	Confidence-based	86.27	84.79	
Popular	Ours (w/o UGO)	87.10	85.98	
	Ours	87.57	86.58	
	Uniform	81.07	81.65	
A J	Confidence-based	84.77	83.40	
Adversariai	Ours (w/o UGO)	85.17	84.25	
	Ours	85.37	84.64	

Table 3: Average results on the POPE COCO benchmark comparing various weights fusion method.

In addition, we provide several qualitative cases that proves ATED strong ability on mitigating hallucinations. These cases uses the instructions with "Please describe this image in detail.", and details are provided in Appendix E.3. 466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

4.3 Ablation Studies

Adapative Uncertainty-Guided Weight. To further validate the effectiveness of the adaptive uncertainty-guided weighting strategy, we conduct an extensive comparative analysis of various weighting approaches on the POPE benchmark. Specifically, we consider uniform weighting, confidence-based weighting, and uncertaintyguided weighting without the uncertainty greedy optimization (UGO) module. As presented in Table 3, our proposed ATED method delivers average improvements of 3.66% in Accuracy and 2.71% in F1 Score over the uniform weighting baseline. Moreover, when the UGO module is removed, the performance of the model ensemble deteriorates to varying extents, indicating that the lack of uncertainty-aware optimization impairs the

effectiveness of the ensemble strategy. These re-488 sults clearly highlight the crucial role of adaptive 489 uncertainty-guided weighting-particularly when 490 enhanced by greedy optimization—in maximizing 491 the performance gains of multimodal model ensem-492 bles. Overall, our findings provide strong empiri-493 cal evidence that the proposed adaptive weighting 494 strategy is fundamental for robust and effective 495 multimodal integration. 496

Impact of Vision Perturbations. We further in-497 vestigate the impact of visual perturbations on hal-498 lucination reduction in LVLM ensemble decoding 499 across different tasks. Specifically, we conduct systematic experiments on the POPE-MSCOCO and MME benchmarks to evaluate the performance 502 503 of dynamic model ensembles under various conditions, including the absence of visual perturbations (Ours(0)) and different levels of perturbation intensity. Experimental results in Table 4 demonstrate that, without adaptation to visual perturbations, 507 508 the performance of multimodal ensemble reasoning significantly degrades on both the POPE and 509 MME datasets—for example, Accuracy decreases 510 by 1.4%, F1-score drops by 2.7%, and Accuracy+ 511 decreases by 20. These findings further highlight 512 that introducing multi-path contrastive decoding 513 under visual perturbations can effectively mitigate 514 hallucinations and enhance reasoning performance. 515

> Additionally, we conduct an ablation study on the amplification factor between the output distributions of original and perturbed visual inputs, denoted as the hyperparameter α in Equation 4, to further validate the effect of visual contrast on model performance. Detailed experimental results are provided in Appendix E.2.

516

517

518

519

520

521

524

525

526

527

531

533

534

537

4.4 Performance Comparison of LVLMs Ensemble Strategy.

To investigate the performance of different LVLMs ensemble strategies across various tasks, we conducted ensemble experiments on the MME and POPE benchmarks usin LLaVA-1.5, InstructBLIP, and LLaVA-NeXT. The results are presented in Table 5. Our experiments reveal that when the performance gap between models is large (e.g., Instruct-BLIP and LLaVA-1.5 exhibit over a 10% difference on the MME benchmark), simple uniform(U) of token probabilities across models fails to improve results and may even degrade overall performance due to noise introduced by lower-performing models. Conversely, when the performance gap is small

Noise	РО	MME	
	Accuracy	F1 Score	Accuracy+
Ours(0)	85.90	84.19	616.67
Ours(200)	87.04	85.86	636.67
Ours(500)	86.73	85.44	608.33
Ours(700)	86.88	85.49	591.67
Ours(999)	87.17	86.59	576.67

Table 4: Evaluation results on POPE and MME with varying noise levels.

Model	РО	MME	
	Accuracy	F1 Score	Accuracy+
InstructBLIP	78.93	79.11	1385.87
LLaVA-1.5	80.37	79.61	1715.40
+ InstructBLIP(U)	83.90	85.14	1437.84
+ InstructBLIP	85.35	85.79	1718.18
+ LLaVA-NeXT	86.55	86.13	1788.09

Table 5: Ensemble performance on POPE and MME benchmarks

538

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

(for instance, LLaVA-1.5's F1 score on POPE exceeds that of LLaVA-NEXT by only about 5%), probability averaging can yield improvements over individual models. ATED overcomes these limitations. Unlike uniform methods, ATED employs an adaptive weighting strategy guided by uncertainty, effectively overcoming the aforementioned limitations and achieving stable performance improvements. This gives ATED greater robustness and broader applicability.

5 Conclusion

In this paper, we propose ATED, the first trainingfree multimodal ensemble decoding method that effectively mitigates hallucinations across diverse multimodal tasks. During inference, ATED performs parallel processing with multiple LVLMs and adaptively fuses token-level logits, enabling finer-grained semantic control and more consistent generation. By introducing an adaptive uncertaintyguided weighting mechanism, ATED dynamically adjusts model importance via uncertainty minimization, enhancing reliability of ensemble inference. Moreover, ATED supports diverse decoding paths, further improving the factual consistency and robustness of generated content. Extensive experiments show that ATED consistently outperforms previous methods on multiple benchmarks, achieving notable gains in both accuracy and robustness.

Limitations

Despite ATED's strong performance in mitigating hallucinations and enhancing robustness, several limitations remain. (i) ATED cannot address all 569 types of hallucinations in LVLMs, which is understandable since the method requires neither addi-571 572 tional training nor modifications to model architectures-factors that may constrain its effectiveness 573 in certain scenarios. (ii) The uncertainty minimization framework relies on estimating output logits over the vocabulary, but the accuracy of this es-576 timation may fluctuate across different tasks and 577 models, potentially affecting the reliability of the ensemble. (iii) Furthermore, our framework depends on the estimation of logits over a shared vocabulary among LVLMs; integrating models with 581 582 differing vocabularies remains an underexplored area for future research. 583

References

584

588

589

593

594

595

598

599

607

610

611

612

613

614

615

616

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- 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. *Preprint*, arXiv:2406.12718.
- Rie Kubota Ando and Tong Zhang. 2005. A framework for learning predictive structures from multiple tasks and unlabeled data. *Journal of Machine Learning Research*, 6:1817–1853.
- Galen Andrew and Jianfeng Gao. 2007. Scalable training of L1-regularized log-linear models. In *Proceedings of the 24th International Conference on Machine Learning*, pages 33–40.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, and 1 others. 2023a. Qwen technical report. *arXiv preprint arXiv:2309.16609*.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023b. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 1(2):3.
- Zechen Bai, Pichao Wang, Tianjun Xiao, Tong He, Zongbo Han, Zheng Zhang, and Mike Zheng Shou.
 2025. Hallucination of multimodal large language models: A survey. *Preprint*, arXiv:2404.18930.

Ali Furkan Biten, Lluis Gomez, and Dimosthenis Karatzas. 2021. Let there be a clock on the beach: Reducing object hallucination in image captioning. *Preprint*, arXiv:2110.01705.

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

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

- Shiqi Chen, Miao Xiong, Junteng Liu, Zhengxuan Wu, Teng Xiao, Siyang Gao, and Junxian He. 2024a. Incontext sharpness as alerts: An inner representation perspective for hallucination mitigation. *Preprint*, arXiv:2403.01548.
- Zhaorun Chen, Zhuokai Zhao, Hongyin Luo, Huaxiu Yao, Bo Li, and Jiawei Zhou. 2024b. Halc: Object hallucination reduction via adaptive focal-contrast decoding. *Preprint*, arXiv:2403.00425.
- Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, and 1 others. 2024c. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *Science China Information Sciences*, 67(12):220101.
- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, and 1 others. 2024d. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 24185–24198.
- Wei-Lin Chiang, Zhuohan Li, Ziqing Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, and 1 others. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. See https://vicuna. lmsys. org (accessed 14 April 2023), 2(3):6.
- Yeongjae Cho, Keonwoo Kim, Taebaek Hwang, and Sungzoon Cho. 2025. Do you keep an eye on what i ask? mitigating multimodal hallucination via attention-guided ensemble decoding. In *Proceedings* of the 2025 International Conference on Learning Representations (ICLR).
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. 2023. Instructblip: Towards general-purpose visionlanguage models with instruction tuning. *Preprint*, arXiv:2305.06500.
- Prasenjit Dey, Srujana Merugu, and Sivaramakrishnan Kaveri. 2025. Uncertainty-aware fusion: An ensemble framework for mitigating hallucinations in large language models. *arXiv preprint arXiv:2503.05757*.
- Thomas G Dietterich. 2000. ensemble methods in machine learning. *International workshop on multiple classifer systems*.
- Alessandro Favero, Luca Zancato, Matthew Trager, Siddharth Choudhary, Pramuditha Perera, Alessandro Achille, Ashwin Swaminathan, and Stefano Soatto. 2024. Multi-modal hallucination control by visual information grounding. *Preprint*, arXiv:2403.14003.

781

728

729

Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, Yunsheng Wu, and Rongrong Ji. 2024. Mme: A comprehensive evaluation benchmark for multimodal large language models. *Preprint*, arXiv:2306.13394.

674

675

685

686

691

698

703

704

707

710

711

712

713

714

715

716

718

719

720

721

722

723

724

725

727

- Zigang Geng, Binxin Yang, Tiankai Hang, Chen Li, Shuyang Gu, Ting Zhang, Jianmin Bao, Zheng Zhang, Han Hu, Dong Chen, and Baining Guo. 2023. Instructdiffusion: A generalist modeling interface for vision tasks. *Preprint*, arXiv:2309.03895.
- Hila Gonen, Srini Iyer, Terra Blevins, Noah A. Smith, and Luke Zettlemoyer. 2024. Demystifying prompts in language models via perplexity estimation. *Preprint*, arXiv:2212.04037.
- Anisha Gunjal, Jihan Yin, and Erhan Bas. 2024. Detecting and preventing hallucinations in large vision language models. *Preprint*, arXiv:2308.06394.
- Qidong Huang, Xiaoyi Dong, Pan Zhang, Bin Wang, Conghui He, Jiaqi Wang, Dahua Lin, Weiming Zhang, and Nenghai Yu. 2024. Opera: Alleviating hallucination in multi-modal large language models via over-trust penalty and retrospection-allocation. *Preprint*, arXiv:2311.17911.
- Drew A. Hudson and Christopher D. Manning. 2019. Gqa: A new dataset for real-world visual reasoning and compositional question answering. *Preprint*, arXiv:1902.09506.
- Fushuo Huo, Wenchao Xu, Zhong Zhang, Haozhao Wang, Zhicheng Chen, and Peilin Zhao. 2025.
 Self-introspective decoding: Alleviating hallucinations for large vision-language models. *Preprint*, arXiv:2408.02032.
- Chaoya Jiang, Haiyang Xu, Mengfan Dong, Jiaxing Chen, Wei Ye, Ming Yan, Qinghao Ye, Ji Zhang, Fei Huang, and Shikun Zhang. 2024. Hallucination augmented contrastive learning for multimodal large language model. *Preprint*, arXiv:2312.06968.
- Dongfu Jiang, Xiang Ren, and Bill Yuchen Lin. 2023. Llm-blender: Ensembling large language models with pairwise ranking and generative fusion. *Preprint*, arXiv:2306.02561.
- Liqiang Jing, Ruosen Li, Yunmo Chen, and Xinya Du. 2024. Faithscore: Fine-grained evaluations of hallucinations in large vision-language models. *Preprint*, arXiv:2311.01477.
- Minchan Kim, Minyeong Kim, Junik Bae, Suhwan Choi, Sungkyung Kim, and Buru Chang. 2024. Esreal: Exploiting semantic reconstruction to mitigate hallucinations in vision-language models. *Preprint*, arXiv:2403.16167.
- 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 Fei-Fei Li. 2017. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International Journal of Computer Vision*.

- Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. 2024. Lisa: Reasoning segmentation via large language model. *Preprint*, arXiv:2308.00692.
- 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. *Preprint*, arXiv:2311.16922.
- Bo Li, Kaichen Zhang, Hao Zhang, Dong Guo, Renrui Zhang, Feng Li, Yuanhan Zhang, Ziwei Liu, and Chunyuan Li. 2024. Llava-next: Stronger llms supercharge multimodal capabilities in the wild.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023a. Blip-2: Bootstrapping language-image pretraining with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR.
- Lei Li, Zhihui Xie, Mukai Li, Shunian Chen, Peiyi Wang, Liang Chen, Yazheng Yang, Benyou Wang, and Lingpeng Kong. 2023b. Silkie: Preference distillation for large visual language models. *Preprint*, arXiv: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. 2023c. M³it: A large-scale dataset towards multimodal multilingual instruction tuning. *Preprint*, arXiv:2306.04387.
- Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. 2023d. Evaluating object hallucination in large vision-language models. *Preprint*, arXiv:2305.10355.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. 2015. Microsoft coco: Common objects in context. *Preprint*, arXiv:1405.0312.
- 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. *Preprint*, arXiv:2306.14565.
- Hanchao Liu, Wenyuan Xue, Yifei Chen, Dapeng Chen, Xiutian Zhao, Ke Wang, Liping Hou, Rongjun Li, and Wei Peng. 2024b. A survey on hallucination in large vision-language models. *Preprint*, arXiv:2402.00253.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024c. Improved baselines with visual instruction tuning. *Preprint*, arXiv:2310.03744.

- 782 783 790 796 797 808 810 811 812 813 814 815 816 817 818 819 820 821 822 828 831

- 825

- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024d. Llavanext: Improved reasoning, ocr, and world knowledge.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. Advances in neural information processing systems, 36:34892-34916.
- Ibomoiye Domor Mienye and Yanxia Sun. 2022. A survey of ensemble learning: Concepts, algorithms, applications, and prospects. IEEE Access, 10:99129-99149.
- David Opitz and Richard Maclin. 1999. Popular ensemble methods: An empirical study. Journal of Artificial Intelligence Research.
- Robi Polikar. 2012. Ensemble learning. Ensemble machine learning: Methods and applications, pages 1 - 34.
- Zexuan Qiu, Zijing Ou, Bin Wu, Jingjing Li, Aiwei Liu, and Irwin King. 2025. Entropy-based decoding for retrieval-augmented large language models. Preprint, arXiv:2406.17519.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning transferable visual models from natural language supervision. Preprint, arXiv:2103.00020.
- Mohammad Sadegh Rasooli and Joel R. Tetreault. 2015. Yara parser: A fast and accurate dependency parser. Computing Research Repository, arXiv:1503.06733. Version 2.
- Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. 2019. Object hallucination in image captioning. Preprint, arXiv:1809.02156.
- Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. 2022. A-okvga: A benchmark for visual question answering using world knowledge. Preprint, arXiv:2206.01718.
- Rico Sennrich, Alexandra Birch, and Barry Haddow. 2016. Edinburgh neural machine translation system for wmt 16. In Proceedings of the First Conference on Machine Translation.
- 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. 2023. Aligning large multimodal models with factually augmented rlhf. Preprint, arXiv:2309.14525.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open

and efficient foundation language models. Preprint, arXiv:2302.13971.

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

- Fanqi Wan, Xinting Huang, Deng Cai, Xiaojun Quan, Wei Bi, and Shuming Shi. 2024. Knowledge fusion of large language models. Preprint, arXiv:2401.10491.
- Junyang Wang, Yiyang Zhou, Guohai Xu, Pengcheng Shi, Chenlin Zhao, Haiyang Xu, Qinghao Ye, Ming Yan, Ji Zhang, Jihua Zhu, Jitao Sang, and Haoyu Tang. 2023. Evaluation and analysis of hallucination in large vision-language models. Preprint, arXiv:2308.15126.
- Xintong Wang, Jingheng Pan, Liang Ding, and Chris Biemann. 2024. Mitigating hallucinations in large vision-language models with instruction contrastive decoding. Preprint, arXiv:2403.18715.
- Sangmin Woo, Donguk Kim, Jaehyuk Jang, Yubin Choi, and Changick Kim. 2024. Don't miss the forest for the trees: Attentional vision calibration for large vision language models. Preprint, arXiv:2405.17820.
- Mingrui Wu, Jiayi Ji, Oucheng Huang, Jiale Li, Yuhang Wu, Xiaoshuai Sun, and Rongrong Ji. 2024. Evaluating and analyzing relationship hallucinations in large vision-language models. Preprint, arXiv:2406.16449.
- Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, and 1 others. 2023. mplug-owl: Modularization empowers large language models with multimodality. arXiv preprint arXiv:2304.14178.
- Shukang Yin, Chaoyou Fu, Sirui Zhao, Tong Xu, Hao Wang, Dianbo Sui, Yunhang Shen, Ke Li, Xing Sun, and Enhong Chen. 2024. Woodpecker: hallucination correction for multimodal large language models. Science China Information Sciences, 67(12).
- Qifan Yu, Juncheng Li, Longhui Wei, Liang Pang, Wentao Ye, Bosheng Qin, Siliang Tang, Qi Tian, and Yueting Zhuang. 2024a. Hallucidoctor: Mitigating hallucinatory toxicity in visual instruction data. Preprint, arXiv:2311.13614.
- 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. 2024b. Rlhf-v: Towards trustworthy mllms via behavior alignment from fine-grained correctional human feedback. Preprint, arXiv:2312.00849.
- Tianyu Yu, Haoye Zhang, Qiming Li, Qixin Xu, Yuan Yao, Da Chen, Xiaoman Lu, Ganqu Cui, Yunkai Dang, Taiwen He, Xiaocheng Feng, Jun Song, Bo Zheng, Zhiyuan Liu, Tat-Seng Chua, and Maosong Sun. 2024c. Rlaif-v: Open-source ai feedback leads to super gpt-4v trustworthiness. Preprint, arXiv:2405.17220.

Zihao Yue, Liang Zhang, and Qin Jin. 2024. Less is more: Mitigating multimodal hallucination from an eos decision perspective. *Preprint*, arXiv:2402.14545.

893

896

900

901

902

903

904

905

906

907

908

909

910

911

912

913 914

915

916

917

918

919

920

921

922

923 924

925

930

931

933

935

937

938

941

- Bohan Zhai, Shijia Yang, Chenfeng Xu, Sheng Shen, Kurt Keutzer, Chunyuan Li, and Manling Li. 2024. Halle-control: Controlling object hallucination in large multimodal models. *Preprint*, arXiv:2310.01779.
- Muru Zhang, Ofir Press, William Merrill, Alisa Liu, and Noah A. Smith. 2023. How language model hallucinations can snowball. *Preprint*, arXiv:2305.13534.
- Shilong Zhang, Peize Sun, Shoufa Chen, Min Xiao, Wenqi Shao, Wenwei Zhang, Yu Liu, Kai Chen, and Ping Luo. 2025. Gpt4roi: Instruction tuning large language model on region-of-interest. *Preprint*, arXiv:2307.03601.
 - Zhiyuan Zhao, Bin Wang, Linke Ouyang, Xiaoyi Dong, Jiaqi Wang, and Conghui He. 2023. Beyond hallucinations: Enhancing lvlms through hallucinationaware direct preference optimization. *Preprint*, arXiv:2311.16839.
- Yiyang Zhou, Chenhang Cui, Jaehong Yoon, Linjun Zhang, Zhun Deng, Chelsea Finn, Mohit Bansal, and Huaxiu Yao. 2024. Analyzing and mitigating object hallucination in large vision-language models. *Preprint*, arXiv:2310.00754.
- Deyao Zhu, Xiang Shen, Xiang Li, and Mohamed Elhoseiny. 2023. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *Preprint*, arXiv:2304.10592.

A More Backgrounds

A.1 Architecture of Vision-Language Models

Large Vision-Language Models (LVLMs) integrate pretrained image encoders and large-scale language models to support tasks such as image captioning and visual question answering. Typically, a frozen vision encoder such as CLIP extracts dense image embeddings (Radford et al., 2021), which are projected into the language space and fed into a decoder-only LLM like LLaMA. Architectures such as BLIP-2 (Li et al., 2023a) and InstructBLIP (Dai et al., 2023) adopt this two-tower design and align the modalities using lightweight adapters or learned commands.

LLaVA-1.5 is a refinement over the original LLaVA model, featuring a simplified architecture and improved training pipeline. It employs CLIP-ViT-L as the vision encoder and a Vicuna-based decoder-only language model. The two modalities are connected via a trainable MLP projection layer, which maps visual tokens into the language embedding space. Trained with visual instruction tuning on synthetic datasets, it achieves strong results across various benchmarks. (Liu et al., 2023)

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

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

InstructBLIP builds on BLIP-2 by introducing an instruction-aware query transformer that conditions the vision encoder's output on task-specific prompts. It integrates a pretrained Vision Transformer (ViT) with a Q-Former, feeding the encoded visual queries to a language model such as Flan-T5 or Vicuna. It is trained using instruction tuning on a collection of 26 datasets. (Dai et al., 2023)

MiniGPT-4 aims to replicate the capabilities of GPT-4-based vision-language systems using opensource components. It integrates a frozen ViTbased vision encoder with Vicuna, connected via a lightweight linear projection layer. Training involves two stages: pre-alignment on image-text pairs followed by fine-tunning on high quality image descriptions. Its minimal parameter count enables efficient multimodal alignment with strong performance in image captioning. (Zhu et al., 2023)

LLaVA-NEXT is an enhanced version of LLaVA-1.5, optimized for higher visual reasoning fidelity. It retains the MLP projection structure but augments training with improved instructionfollowing datasets and higher-resolution visual inputs. It achieves better performance in OCR, compositional reasoning, and world knowledge benchmakrs. (Liu et al., 2024d)

A.2 Ensemble Learning in NLP

Ensemble learning has long been a reliable strategy in machine learning to improve robustness, reduce overfitting, and enhance generalization. By combining the predicitions of multiple models or decision rules, ensembles can correct individual biases and reduce the variance of outputs (Dietterich, 2000), Classical ensemble methods include bagging, boosting, and stacking, all of which have demonstrated strong performance in classification tasks such as sentiment analysis, topic classification, and named entity recognition (Opitz and Maclin, 1999).

In the domain of Natural Language Processing (NLP), ensemble methods have been applied extensively in both structured prediction and generation tasks. For example, ensemble decoding, which involves averaging or voting across multiple language models. has been shown to improve fluency and factuality in neural machine translation (Sennrich et al., 2016). Recent work has also explored

993

994

- 9
- 999 1000
- 100
- 1002 1003
- 1004
- 10

1006

1007

1008

1009

1011

1012

1013

1014

1015

1016

1017

1019

1020

1021

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035 1036

1037

1038

1040

ensemble inference for large language models, aggregating outputs per-token-level logits from multiple sources to improve consistency and reduce hallucination

In multimodal learning, ensemble approaches are gaining traction as a decoding-level intervention. Rather than relying on a single model's output, ensembles constructed across different model variants or decoding stratigies can better capture complementary evidence, making them suitable for suppressing the hallucination in vision-language tasks.

B Details about Baseline

To mitigate hallucination without retraining, a variety of decoding-time techniques have been proposed:

- ICD: (Wang et al., 2024) Instruction-Contrastice Decoding leverages instructionlevel perturbations to reduce hallucination in multimodal large language models. It operates by introducing minimal semantic alterations to the input prompt, such as inserting irrelevant phrases or modifiying the question structure, and then comparing the model's output distributions under both the original and perturbed instructions. Tokens exhibiting instability across these variants are identified as potentially hallucinated and are downweighted during generation.
 - SID (Self-Introspective Decoding): (Huo et al., 2025) Self-Introspective Decoding mitigateas hallucinations by introspectively filtering low-relevance visual signals during generation. It evaluates the contextual alignment of visual tokens with both the preceding textual context and the decoding history, retaining only those with strong semantic relevance. By pruning distractive or semantically weak visual features early in decoding, SID improves grounding accuracy, particularly in complex or visually dense scenarios.
 - VCD (Visual-Contrastive Decoding): (Leng et al., 2023) Visual-Contrastive Decoding aims to improve visual consistency by introducing small-scale perturbations to the visual input and contrasting the model's responses. The approach applies controlled distortions, such as Gaussian blur, occlusion, or token

masking, to the image embeddings and mea-
sures output divergence. Tokens highly sensi-
tive to such perturbations are treated as visu-
ally fragile and are penalized during decoding.1041
1042

1045

1046

1047

1048

1049

1050

1051

1052

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1067

1068

1069

1070

1071

1072

- OPERA (Overtrust Penalty with Retrospective Adjustment): (Huang et al., 2024) Overtrust Penalty with Retrospective Adjustiment introduces a two-stage mechanism to address hallucination in multimodal generation: overtrust penalty and retrospective adjustment. During decoding, it applies a regularization term to suppress overconfident token predictions that exhibit weak visual grounding. After generation, a retrospective evaluation is performed to re-rank or adjust outputs based on their semantic agreement with the image.
- Ensemble Decoding (ED): (Cho et al., 2025) Ensemble Decoding combines multiple generation pathways to improve robustness and reduce hallucination. It operates by aggregating outputs from a set of models or decoding configurations, such as different random seed, visual crops, or temperature settings, and fusing them through majority voting, logit averaging, or response re-ranking. This ensemble process helps to mitigate the influence of unstable or outlier predicitions by emphasizing consensus across multiple decoders

All these methods operate without modifying model parameters, offering flexible, training-free solutions for enhancing visual faithfulness during inference.

C Evaluation Metric Details

The Polling-based Object Probing Evaluation 1074 (POPE) benchmark is a systematic framework 1075 designed to assess object hallucination in Large 1076 Vision-Language Models (LVLMs) during image 1077 description tasks. POPE employs a binary question-1078 answering format, using prompts such as "Does the 1079 image contain ____?" to evaluate a model's ability 1080 to accurately determine the presence or absence of 1081 specific objects within images. To construct nega-1082 tive samples-instances where the object is absent 1083 from the image-POPE utilizes three distinct strate-1084 gies: random sampling involves selecting objects 1085 that do not appear in the image at random; pop-1086 ular sampling selects absent objects from a pool 1087 of frequently occurring objects across the dataset; 1088

adversarial sampling prioritizes objects that com-1089 monly co-occur with present objects but are absent 1090 in the current image. The benchmark integrates 1091 three datasets: MSCOCO, A-OKVQA, and GQA. 1092 From each dataset, 500 images are selected, and 1093 six questions are generated per image, resulting in 1094 a total of 27,000 query-answer pairs for evaluation. 1095 Performance is measured using standard metrics, 1096 including accuracy, precision, recall, and F1 score, 1097 with higher values indicating a model's superior 1098 capability in mitigating hallucinations such as fabricated objects and erroneous descriptions. 1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

The Caption Hallucination Assessment with Image Relevance (CHAIR) metric is a specialized evaluation framework designed to quantify object hallucination in image captioning models. CHAIR assesses the alignment between generated captions and the actual visual content by comparing the objects mentioned in the captions against groundtruth annotations from datasets like MSCOCO.

$$C_{S} = \frac{|\{\text{hallucinated objects}\}|}{|\{\text{all mentioned objects}\}|}$$
(9)
$$C_{I} = \frac{|\{\text{captions w/ hallucinated objects}\}|}{|\{\text{all captions}\}|}$$
(10)

The metric comprises two variants: CHAIRi (instance-level) and CHAIRs (sentence-level). CHAIRi calculates the proportion of hallucinated object mentions relative to all object mentions in the generated captions, while CHAIRs measures the fraction of sentences that contain at least one hallucinated object. Lower values in both metrics indicate better performance in mitigating object hallucinations.

The Multimodal Model Evaluation (MME) 1120 benchmark offers a comprehensive framework for 1121 assessing Large Vision-Language Models (LVLMs) 1122 across a spectrum of tasks, encompassing both 1123 perceptual and cognitive dimensions. Specifically, 1124 MME comprises ten perception-oriented subtasks 1125 and four cognition-focused ones, facilitating a 1126 holistic evaluation of LVLM capabilities .In the 1127 context of object-level hallucination evaluation, 1128 MME includes dedicated subsets targeting the "ex-1129 istence" and "count" tasks. The "existence" task 1130 assesses a model's ability to accurately identify the 1131 presence or absence of specific objects within an 1132 1133 image, while the "count" task evaluates the model's proficiency in determining the correct number of 1134 instances of a given object. These tasks is quan-1135 tified using a combined metric of accuracy and 1136 accuracy+. Accuracy measures the proportion of 1137

correct predictions, while accuracy+ accounts for 1138 near-correct responses. 1139

1140

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

D Implementation Details

In all experimental settings, the hyper-parameter 1141 α is fixed at 1. For visual perturbations in the 1142 model ensemble, we adopt a noise-injection strat-1143 egy, setting the noise steps T to 200 for MME, 1144 500 for LLaVA-Bench, and 999 for POPE. For 1145 OPERA, VCD, and SID, we use the default set-1146 tings as specified in their original papers. Greedy 1147 decoding is used for comparison methods, while for 1148 open-ended generation tasks (such as CHAIR and 1149 LLaVA-Bench), we employ sampling with Top-p 1150 = 1. All experiments are conducted on Nvidia A40 1151 GPUs. 1152

E More Detailed Comparison

E.1 More Results on on CHAIR

The hyperparameter max new tokens, which controls the maximum length of generated responses, plays a critical role in CHAIR-based evaluation. In the main text, we report results using a setting of max new tokens = 64. Additional results under a relaxed constraint of max new tokens = 512are provided in Table 6. As Table illustrates, the generation length limit has a substantial impact on LVLM performance under the CHAIR metric. when the token budget is increased from 64 to 512, our method consistently outperforms all baselines on the metric CHAIR_S, highlighting its robustness and adaptability under varying generation lengths. Furthermore, our model produces responses with an average length of 107.4 tokens as shown in Table 7, indicating that the observed reduction in object hallucinations is achieved without compromising the richness of the generated descriptions.

E.2 More Results on on MME

ATED is designed to integrate the expertise of mul-1174 tiple models, thereby bridging the hallucination gap 1175 that exists among different LVLMs during infer-1176 ence. To further investigate whether our approach 1177 not only preserves but also potentially enhances the 1178 fundamental perception and reasoning capabilities 1179 of LVLMs across a broader range of multimodal 1180 tasks, we also analyze the comprehensive perfor-1181 mance on the MME benchmark, which consists of 1182 14 sub-tasks for evaluating perception and recog-1183 nition. As shown in Table 8, our method $(Ours_{\#})$ 1184 significantly outperforms all baseline approaches 1185

Туре	METHOD	LLaVA-1.5 CHAIR _S CHAIR _I		InstructBLIP CHAIR _S CHAIR _I		MiniGPT4 CHAIR _S CHAIR _I		LLaVA-NeXT CHAIR _S CHAIR _I	
	Default	51.3	16.8	55.6	24.2	33.6	19.4	42.6	14.1
	OPERA	46.4	13.0	47.1	12.4	26.4	10.7	39.4	11.8
Single	VCD	51.7	15.6	51.0	16.7	30.4	14.2	41.1	12.9
	ICD	47.4	13.9	46.3	15.3	32.6	13.1	42.1	12.6
	SID	44.2	12.2	42.3	12.4	28.5	11.7	40.8	13.0
	Ours		ATI	ED ^{&}		ATED [#]			
Ensemble		CHAIR _S		CHAIRI		$CHAIR_S$		CHAIRI	
		-		-		34.0		17.1	

Table 6: **CHAIR evaluation results on different decoding strategies**. Results are from the papers or re-implemented based on official codes. *Note:* & denotes ensemble with InstructBLIP, # denotes ensemble with InstructBLIP and LLaVA-NeXT.

Method	Length
Default	100.6
OPERA	98.6
VCD	100.4
ICD	106.3
Ours [#]	107.4

Table 7: Comparison of CHAIR performance across different methods in terms of output length on LLaVA-1.5.

Method	Accuracy+
Default	1715.40
OPERA	1773.52
VCD	1756.02
ICD	1749.43
SID	1770.43
Ours [#]	1788.09

Table 8: Comparison of total accuracy+ cacross different methods on LLaVA-1.5.

based on the LLaVA-1.5 backbone, surpassing both the original LVLMs and the best-performing baselines by a substantial margin (+18.34). These results indicate that our approach not only effectively manages hallucination during inference but also improves the accuracy of the underlying LVLMs on fundamental tasks.

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195 1196

1197

1198

1199

1200

Table 9 presents the quantitative evaluation results of the model under different α values on object-level metrics (Existence, Count), attributelevel metrics (Position, Color), and the overall accuracy (Total Accuracy+). As α increases from 0.5 to 1.0, all metrics demonstrate varying degrees of improvement, with Color showing the most substantial gain—from 140 to 155. These improvements are reflected in the Total Accuracy+, which rises from 595.00 to 636.67 as α increases. Moreover, we observe that attribute-level metrics are more sensitive to changes in the intensity of visioncontrastive regularization compared to object-level metrics, resulting in greater improvements. This finding indicates that appropriately tuning the α parameter not only enhances the model's ability to confirm object information during adaptive ensemble inference but also significantly improves its capability to capture fine-grained attribute details. As a result, the overall prediction accuracy and robustness are further strengthened. 1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

E.3 Qualitative Analysis

To further evaluate whether ATED effectively mitigates hallucinations beyond quantitative metrics in open-ended generation tasks, we conducted a qualitative analysis on the MSCOCO dataset, using several decoding strategies as baselines. The LVLMs are prompt with "Please describe this image in detail", with the maximum token limit set to 150. As illustrated in Figure 8 and Figure 9, baseline methods including the default decoding, OPERA, and VCD often produce hallucinated content (highlighted in red). In contrast, ATED dynamically selects and weights token-level outputs from multiple models at each decoding step, guided by a greedy uncertainty-minimization strategy. This enables the model to better adapt to contextual environments and significantly improves the credibility and robustness of the generated content.

In addition, we perform GPT-assisted evaluation on the LLaVA-Bench benchmark (Liu et al., 2023). Following evaluation protocol proposed by (Yin et al., 2024; An et al., 2025), the model is presented with an image and two candidate descriptions, structured according to the prompt format

α	Object-	Level	Attribute	e-Level	Total Accuracy+	
u	Existence	Count	Position	Color	10001 Hecurucy I	
0.5	180.00	143.33	131.67	140	595.00	
0.7	180.00	143.33	136.67	140	600.00	
1.0	185.00	158.33	138.33	155	636.67	

Table 9: Quantitative results on Object-level (Existence, Count), Attribute-level (Position, Color), and Total Accuracy+ for using various noise steps.

shown in Figure 5. The GPT-40 API is employed to evaluate the generated responses in terms of factual accuracy (Accuracy) and descriptive richness (Detailedness).

Furthermore, we conducted an additional evaluation based on GPT-4, following the methodology outlined in (Zhao et al., 2023). Specifically, we randomly sampled 200 images from the Visual 1245 1246 Genome (VG-100K) dataset (Krishna et al., 2017) and assessed model performance by comparing the generated descriptions with the region descriptions 1248 associated with each image. This comparison al-1249 lows for effective identification of hallucinated content based on semantic inconsistencies. We comprehensively analyzed five key metrics: sentences per image (SPI), words per image (WPI), hallucinated sentence ratio (HSR), hallucinated word ratio (HWR), and mean hallucination ratio (MHR). No-1255 tably, higher SPI and WPI values, as well as lower 1256 HSR, HWR, and MHR, indicate better model performance. In the radar charts, a larger area reflects 1258 superior performance. Multiple models and decoding strategies were included as baselines for comparison. The detailed results are presented in Figure 4. As shown, the proposed ATED method substantially reduces hallucination and effectively 1263 suppresses misleading content during generation.



Figure 4: GPT-4 assisted hallucination evaluation.

1238

1239 1240

1241

1242

1243

1244

1247

1250 1251

1252

1253

1254

1257

1259

1260

1261

1262

GPT-40 Prompt

You are an AI designed to evaluate and score the performance of two AI assistants in describing a given image. Your primary focus is on the accuracy and detailedness of their descriptions. You will assess the accuracy by checking for hallucinations-any part of the description that is inconsistent with the image content. For detailedness, you will consider how rich the response is in necessary details, excluding any hallucinated parts. You will provide scores on a scale from 1 to 10 for each assistant separately, based on these criteria. After scoring, you will offer an explanation for your evaluation, ensuring it is free from bias and not influenced by the order of presentation of the responses. Input format:

[Assistant 1] {} [End of Assistant 1] [Assistant 2] {} [End of Assistant 2] Output format: Accuracy: Scores of the two answers: Detailedness: Scores of the two answers:

Figure 5: GPT-40 prompt format.



GPT-4o Evaluation: LLaVA-1.5 (OPERA) Accuracy: 4/10 Detailedness: 5/10

GPT-4o Evaluation: Ours Accuracy: 9/10 Detailedness: 7/10



Instr: Please describe this image in detail Max New Token = 150



The image depicts a large, <u>white skyscraper</u> towering over a city with a clear blue sky in the background. The skyscraper is located in the center of the city, surrounded by other buildings of various heights. The city is bustling with activity, as evidenced by the <u>presence of numerous cars and trucks on the streets</u> <u>below the skyscraper</u>.



The image you've shared is a close-up of a tall, slender tower. The tower is predominantly white, with a few black accents that add a touch of contrast. The tower is situated in a city, as evidenced by the buildings and structures that surround it.

Accuracy: Assistant 1: This description contains significant hallucinations. The structure in the image is the Space Needle in Seattle, which is not a "large, white skyscraper" and does not resemble one. It is a single, iconic observation tower with a distinct design. The claim of "numerous cars and trucks on the streets" is also inaccurate, as there are no visible streets bustling with traffic in the image. Assistant 2: This response is mostly accurate. The description identifies the subject as a "tall, slender tower" that is "predominantly white with a few black accents," which matches the appearance of the Space Needle. It also matches the setting by mentioning the city with surrounding buildings.

Detailedness: Assistant 1: While this response attempts to provide a more elaborate description of the setting, most of the details are fabricated and therefore detract from the evaluation. It does not mention defining features of the tower, such as its observation deck, tapered legs, or unique aesthetic. Assistant 2: This description is more concise but relatively to the point. The mention of black accents and the general shape of the tower improve its detailedness. However, it omits several crucial features, such as the observation deck, the flag, and the surrounding greenery in view.

Figure 6: Qualitative cases on LLaVA-Bench. The hallucinated content is highlighted in red.



GPT-4o Evaluation: LLaVA-1.5 (VCD) Accuracy: 4/10 Detailedness: 6/10

GPT-4o Evaluation: Ours Accuracy: 9/10 Detailedness: 7/10 A cute little dog is laying down on a wooden floor, resting LLaVA-1.5 vcD and appears to be sleeping on a blue-colored surface.

> The image features a small dog with a predominantly white coat, lying on its side on a wooden surface. The dog appears to be resting or sleeping, with its head comfortably resting on the wooden surface. The background is out of focus, but it seems to be an indoor setting.

Accuracy: Assistant 1 incorrectly described the dog as "white and brown" and stated that it was "sleeping on a blue-colored surface" which is inaccurate. The dog's coat is light tan with slight white markings, and it is lying directly on a blue wooden floor, not "a bluecolored surface" like a blanket. Additionally, Assistant 1 claimed the dog was "resting its head on its paws" which is not the case—its head is resting on the floor. Assistant 2, on the other hand, provided an accurate description of the dog and its environment, correctly identifying that the dog is lying on its side on the blue wooden floor and avoiding any inaccurate claims.

Detailedness: Assistant 1 included some descriptive details, such as identifying the dog's coloration (albeit inaccurately) and the mention of the "blue-colored surface". However, these inclusions were not meaningful due to inaccuracies. Assistant 2 provided sufficient and accurate details without adding unnecessary or hallucinated elements, such as describing the dog, its position, and the surrounding background accurately. Assistant 2's response is slightly more detailed due to a richer contextual hint about the background being "out of focus" and the scene being "an indoor setting."

Figure 7: Qualitative cases on LLaVA-Bench. The hallucinated content is highlighted in red.

Ours

LLaVA-1.5 (VCD): The image features a woman lying on a bed, with her head resting on her hand. She <u>appears</u> to be looking at the camera, possibly posing for a picture. The bed is covered with a striped blanket, adding a cozy touch to the scene.



Instr: Please describe this image in details

Max new token = 64

LLaVA-1.5 (OPERA): The image features a beautiful young woman <u>sitting on a couch</u>, wearing a white shirt. She appears to be posing for the camera, possibly for a magazine or advertisement. The couch she is sitting on has a striped pillow, adding a touch of color to the scene.

InstructBLIP (OPERA): The image features a young woman lying on a bed, wearing a white shirt and a red lipstick. She is positioned in the middle of the bed, with her hair cascading down her back. There are two pillows visible in the scene, one on the left side of the bed and

Ours: The image features a beautiful young woman with long, dark hair. She is lying down on a bed, dressed in a white shirt. The bed she is on is covered with a white loomforter.

Figure 8: Qualitative cases on MSCOCO. The hallucinated content is highlighted in red.



Instr: Please describe this image in details Max new token = 64 *LLaVA-1.5 (VCD):* The image features a woman sitting at a dining table, holding a large pizza in her hands. She is smiling and appears to be enjoying the moment. The pizza is <u>placed in the center of the table</u>, and the woman is positioned on the left side of the table."

LLaVA-1.5 (OPERA): The image features a woman sitting at a dining table, holding up a large pizza with a smile on her face. The pizza is placed in front of her, occupying a significant portion of the table. <u>There are</u> several other people in the scene, some sitting and others standing around the dining

InstructBLIP (OPERA): The image features a woman isitting at a dining table, holding a large pizza in her hands. She is smiling and appears to be eager to share the pizza with others. There are several chairs placed around the itable, suggesting that there are other people present.

Ours: The image features a woman sitting in front of a dining table. She is holding a large pizza in her hands, which is placed on the table. The woman appears to be in the process of serving the pizza, as she is holding it up.

Figure 9: Qualitative cases on MSCOCO. The hallucinated content is highlighted in red.