MITIGATING MULTIMODAL HALLUCINATIONS VIA GRADIENT-BASED SELF-REFLECTION

Anonymous authors

Paper under double-blind review

A APPENDIX

A.1 FIRST ORDER TAYLOR EXPANSION

The first order Taylor expansion of \mathcal{F}_{θ^*} w.r.t \mathbf{t}^v , \mathbf{t}^p and, all the previous output tokens \mathbf{y}_i with i < m:

$$\begin{aligned} \mathcal{F}_{\theta^{*}}(\mathbf{t}^{v},\mathbf{t}^{p})_{m} \approx & \mathcal{F}_{\theta^{*}}(\mathbf{t}^{v(0)},\mathbf{t}^{p(0)})_{m} + \sum_{s=1}^{S} \frac{\partial(\mathcal{F}_{\theta^{*}})_{m}}{\partial t_{s}^{v}} \Big|_{\mathbf{t}^{v(0)}} (t_{s}^{v} - t_{s}^{v(0)}) \\ &+ \sum_{n=1}^{N} \frac{\partial(\mathcal{F}_{\theta^{*}})_{m}}{\partial t_{n}^{p}} \Big|_{\mathbf{t}^{p(0)}} (t_{n}^{p} - t_{n}^{p(0)}) + \sum_{i=1}^{m-1} \frac{\partial(\mathcal{F}_{\theta^{*}})_{m}}{\partial y_{i}} \Big|_{y_{i}^{(0)}} (y_{i} - y_{i}^{(0)}) \qquad (1) \\ &= \sum_{s=1}^{S} \mathbf{g}_{ms}^{v} \cdot t_{s}^{v} + \sum_{n=1}^{N} \mathbf{g}_{mn}^{p} \cdot t_{n}^{p} + \sum_{i=1}^{m-1} \mathbf{g}_{mi}^{y} \cdot y_{i} + Const, \end{aligned}$$

where $y_i^{(0)} := \mathcal{F}_{\theta^*}(\mathbf{t}^{(0)}, \mathbf{t}^{p(0)})_i$ is the *i*th output token given the input.

A.2 INTERPRETING CONTRASTIVE DECODING THROUGH KL DIVERGENCE

Kullback-Leibler (KL) divergence can be used to interpret contrastive decoding, For example, in GACD-VA, it measures the divergence between the reference distribution $p(y_{cm}|\mathcal{F}_{\theta^*}, \mathbf{t}^p, y_{< m})$ to the visual joint distribution $p(y_{cm}|\mathcal{F}_{\theta^*}, \mathbf{t}^v, \mathbf{t}^p, y_{< m})$.

$$D_{KL} = \sum_{c} p(y_{cm} | \mathcal{F}_{\theta^{\star}}, \mathbf{t}^{v}, \mathbf{t}^{p}, y_{

$$= \sum_{c} p(y_{cm} | \mathcal{F}_{\theta^{\star}}, \mathbf{t}^{v}, \mathbf{t}^{p}, y_{

$$= \sum_{c} p(y_{cm} | \mathcal{F}_{\theta^{\star}}, \mathbf{t}^{v}, \mathbf{t}^{p}, y_{

$$- \mathcal{F}_{\theta^{\star}}(\mathbf{t}^{p})_{m}[c] + \log(\sum \exp(\mathcal{F}_{\theta^{\star}}(\mathbf{t}^{v}, \mathbf{t}^{p})_{m} - \mathcal{F}_{\theta^{\star}}(\mathbf{t}^{p})_{m})[c] + Const),$$
(2)$$$$$$

where $p(y_{cm}|\mathcal{F}_{\theta^*}, \mathbf{t}^v, \mathbf{t}^p, y_{< m}) = \sigma(\mathcal{F}_{\theta^*}(\mathbf{t}^v, \mathbf{t}^p)_m), p(y_{cm}|\mathcal{F}_{\theta^*}, \mathbf{t}^p, y_{< m}) = \sigma(\mathcal{F}_{\theta^*}(\mathbf{t}^p)_m) \text{ and } c$ represents a class in the predefined vocabulary. The adjustment term $\mathcal{F}_{\theta^*}(\mathbf{t}^v, \mathbf{t}^p)_m - \mathcal{F}_{\theta^*}(\mathbf{t}^p)_m$ increases the KL divergence, thereby emphasizing the impact of visual tokens.

A.3 DIFFERENT SAMPLING STRATEGIES

Tab. 1 presents an ablation study on sampling strategies, showing that our method performs better with greedy sampling in both discriminative and generative VQA tasks. This reflects the precision of our logit adjustments, where the highest-confidence prediction is generally the most accurate.

	POPE MSCOCO Adversarial					LLaVA-QA90						
strategy	LLaVA-v1.5		InstructBLIP		mPLUG-Owl2		LLaVA-v1.5		InstructBLIP		mPLUG-Owl2	
I	Acc	F1	Acc	F1	Acc	F1	Acc	Det	Acc	Det	Acc	Det
random 8 greedy 8	32.3 33.5	81.1 82.1	82.2 82.5	81.8 82.1	83.2 84.2	82.9 83.7	5.79 6.20	4.74 5.13	5.98 6.28	4.64 4.77	6.07 6.69	5.72 6.28

Table 1: Ablation Study on Sampling Strategies.

A.4 ABLATION STUDY ON THE GACD PROCESS WITH 100% CONFIDENCE

We conducted an ablation study to evaluate the effect of processing GACD even with 100% confidence. Results in Tab. 2 show that applying GACD with 100% confidence yields only marginal improvements, indicating that the benefits of GACD are more pronounced when the model's confidence is lower.

Table 2: Ablation Study on 100% Confidence with mPLUG-Owl2

w/ 100% Confiden	onfidence POPE N	ISCOCO Adversarial	LLaVA-QA90		
	Acc	F1	Acc	Det	
X	84.2	83.7	6.69	6.28	
1	84.2	83.7	6.75	6.29	

A.5 MLLMs Architectures

Tab. 3 shows detailed information about the vision encoder and LLM components of the MLLM architectures used in our experiments.

Table 3: Details of the used MLLM architectures.

ILLMs	Vision encoder	LLM
LLaVA-v1.5	CLIP-L-336px	Vicuna-v1.5-7B
LLaVA-v1.6	CLIP-L-336px	Vicuna-v1.5-7B
InstructBLIP	BLIP-2	Vicuna-v1.1-7B
mPLUG-Owl2	CLIP-L	LLaMA-2-7B

A.6 OTHER RESULTS OF POPE

We report our experimental results on the POPE dataset, in addition to MSCOCO and adversarial settings, in Tab. 4. The results indicate that our method improves performance across all baseline MLLMs, with more significant gains observed in the adversarial setting. This is expected, as visual inputs provide crucial clues for identifying adversarial objects.

- A.7 INFLUENCE RATIO IN VQA

Fig. 1 illustrates the influence ratio across predicted tokens in VQA tasks, comparing baseline pre-dictions with those obtained after applying the GACD-VA addon. The analysis confirms that as more tokens are generated, the influence of both visual and prompt tokens diminishes, while the contribution of previous output tokens steadily increases, shifting the primary drivers of the model's predictions. It also shows that visual tokens generally have a lower influence compared to prompt tokens across MLLMs, particularly at the beginning. The application of GACD-VA rebalances these influences, boosting the visual token influence to align with the dominant components-prompt at the beginning and previous outputs later-thereby reducing the likelihood of hallucinations in the model's predictions.

Setting

Random

Popular

Random

Popular

Random

Popular

Adversarial X

Adversarial X

w/Ours

× ✓

X

1

× √

X

1

× √

X

1

108 109

110

- 111 112
- 113 114
- 115
- 116 117
- 118
- 119 120

121

122

123

124 125

Dataset

MSCOCO

A-OKVQA

GQA

126 127

128 129

130 131

132

133 134 135

136 137

1	3	8
1	3	9

140 141

142

143 144

145



148

149 150

151 152

153

154

155 156

157

158

159

160

161

\mathbf{a}
-
. /

Table 4: More Results on POPE ?.	
----------------------------------	--

F1 ↑

84.8

85.1

83.8

84.0

87.6

87.4

85.1

85.1

79.9

80.1

88.2

88.2

84.1

84.1

81.3

81.6

InstructBLIP

Acc \uparrow

87.1

87.9

84.2

85.0

88.5

88.8

81.9

82.3

74.8

75.3

87.2

87.2

78.6

78.8

75.9

76.1

 $F1\uparrow$

85.7

86.8

83.6

84.3

88.5

88.8

83.1

83.4

77.9

78.2

87.1

87.2

80.4

80.4

78.4

78.5

mPLUG-Owl2

Acc \uparrow

86.0

87.9

84.6

86.4

86.5

88.4

82.4

85.1

74.7

78.2

85.2

86.1

78.7

81.0

76.4

79.2

F1 ↑

84.4

87.1

83.2

85.7

85.7

88.1

82.2

85.3

76.9

79.9

84.0

85.0

78.5

80.5

76.8

79.1

LLaVA-v1.5

Acc \uparrow

86.5

86.8

85.5

85.6

88.0

88.1

85.5

85.5

79.1

79.5

88.9

88.9

84.1

84.2

80.8

81.1



Figure 1: Influence Ratio across Predicted Tokens in VQA: (left) Baseline predictions; (right) Predictions with GACD-VA. In general, visual tokens initially contribute less than prompt tokens. As more tokens are generated, the contribution of visual tokens and prompt tokens decreases, while the influence of previous output tokens increases. GACD-VA effectively boosts the influence of visual tokens to match the dominant components—prompt tokens at the start and previous output tokens toward the end—thereby mitigating hallucinations.