Critic-V: VLM Critics Help Catch VLM Errors in Multimodal Reasoning

Supplementary Material

7. Pseudo-code for Main Algorithms

Algorithm 1 Bug Insertion and Rule-based Reward for Preference Data Collection

- 1: Input: True answer A_{true} , Question-Image pair $(Q^{(i)}, I^{(i)})$
- 2: **Output:** Critique score score
- 3: Step 1: Generate a fake answer with inserted bugs
- 4: $\mathcal{A}_{\text{fake}} \leftarrow \mathcal{A}_{\text{true}}$
- 5: Randomly choose number of fake details $n \leftarrow$ random integer between 1 and 5
- 6: for each fake detail d_j from 1 to n do
- 7: Insert bug into $\mathcal{A}_{\text{fake}}$ by adding $d_i \in \mathcal{D}_{\text{fake}}$
- 8: end for
- 9: Fake answer generation complete.
- 10: Step 2: Extract details from true answer and fake answer
- 11: $\mathcal{D}_{true} \leftarrow Extract details from \mathcal{A}_{true}$
- 12: $\mathcal{D}_{fake} \leftarrow Extract details from \mathcal{A}_{fake}$
- 13: Step 3: Generate critique from a VLM
- 14: $\mathcal{D}_{detected} \leftarrow \text{Use VLM}$ to detect errors in \mathcal{A}_{fake}
- 15: Step 4: Calculate critique quality using Jaccard index
- 16: $\mathcal{D}_{\text{true}} \leftarrow \{d_1, d_2, \dots, d_m\}$
- 17: $\mathcal{D}_{\text{detected}} \leftarrow \{d'_1, d'_2, \dots, d'_n\}$
- 18: intersection $\leftarrow \mathcal{D}_{true} \cap \mathcal{D}_{detected}$
- 19: union $\leftarrow \mathcal{D}_{true} \cup \mathcal{D}_{detected}$
- 20: $Jaccard(\mathcal{D}_{true}, \mathcal{D}_{detected}) \leftarrow \frac{len(intersection)}{len(union)}$
- 21: Step 5: Calculate critique score based on rule-based reward
- $Jaccard(\mathcal{D}_{true}, \mathcal{D}_{detected}) + 0.1 \times$ 22: score GPT-based score
- 23: Return: score

Algorithm 2 Training Critic Model with DPO

- 1: Input: Dataset $\mathcal{D}_{cri} = \{(Q^{(i)}, I^{(i)}, C_w^{(i)}, C_l^{(i)})\}_{i=1}^N,$ base model π_{ref} , learning rate α , and hyperparameter β
- 2: Initialize critic model $\pi_{\theta} \leftarrow \pi_{\text{ref}}$
- 3: for each batch (Q^i, I^i, C^i_w, C^i_l) in \mathcal{D}_{cri} do
- Compute the critique logits for the preferred (C_w) $4 \cdot$ and disfavored (C_l) critiques
- Compute the DPO loss 5:
- Compute gradients of \mathcal{L}_{DPO} w.r.t. π_{θ} 6:
- Update π_{θ} using gradient descent 7:
- 8: end for
- 9: **Output:** Trained critic model π_{θ}

Alg	gorithm 3 Reasoner-Critic Framework				
1:	Input: Query Q, Input image I, Reasoner $\pi_{\theta^{reasoner}}$				
	Critic $\pi_{\theta^{critic}}$, Maximum iterations <i>max_iterations</i>				
2:	Output: Final response <i>R</i> _{final}				
3:	Initialize Preasoner (initial prompt for Reasoner)				
4:	response = $\pi_{\theta^{reasoner}}(P^{reasoner}, I)$				
	(generate initial response)				
5:	for iteration = 1 to max_iterations do				
6:	Critic evaluates response:				
7:	critique = $\pi_{\theta^{critic}}(\delta P^{reasoner} P^{reasoner},Q,R)$				
	(Critic generates critique)				
8:	If Critic determines that critique is satisfactory:				
9:	break (end loop if critique is satisfactory)				
10:	Else:				
11:	reasoner updates prompt: $P^{reasoner} \leftarrow P^{reasoner} +$				
	$\delta P^{reasoner}$				
12:	response = $\pi_{\theta^{reasoner}}(P^{reasoner}, I)$				
	(generate new response)				

- 13: end for
- 14: **Return:** $R_{final} = response$

8. Prompt Template

For multiple-choice questions (MCQ), the template of prompt is designed as follows,

Hint: {hints}
Question: {question}
Options: {options}
Please select the correct answer from the options
above.

As well as open-ended visual question-answering (VQA) tasks,

{question_text} Please try to answer the question with short words or phrases if possible.

We utilize the prompt above to help Reasoner generate explanations of their answers. Then, we let the Critic generate critiques on the answer with the prompt below:

Question
{question}
Answer
{result}
Task
Please provide a critique of the answer above. What
are the weaknesses of the answer?

After that, we use these weaknesses(or errors) from

Evaluation Criterion	Evaluation Criterion Question Description		
Coverege Applysic	Did the model identify all the hallucinations mentioned in the correct answer?	Vac / No	
Coverage Analysis	Are there any significant hallucinations that were missed?	Tes / INO	
Accuracy Assessment	Are the detected items genuine hallucinations (true positives)?	Vac / No	
Accuracy Assessment	Are there any false detections (false positives)?	1057110	
Precision of Description	How precise and clear are the model's descriptions of the detected hallucinations?	Vac / No	
recision of Description	Is the explanation specific enough to understand what exactly is wrong?	1037110	
Overall Effectiveness	How effective is this detection compared to an ideal detection?	Ves / No	
Overall Effectiveness	Does it provide practical value for hallucination identification?	1057 110	
Comprehensive Evaluation	Based on your analysis, please provide a brief explanation of your evaluation.	Text Input	
Final Score	Based on the scoring criteria, provide a final score from 0 to 10.	0-10	

Table 5. GPT-4o Evaluation for Erroneous Detection

Critic to let Reasoner correct their answers with the following prompt:

Reflection on former answer:	
{critics}	
{original_question}	

9. The GPT-40 Evaluation Rules

In this section, we provided a detailed description of the evaluation criteria for erroneous detected by the VLMs as shown in Table 5.

10. Hyperparameters of Critic Model's Training

We use the Qwen2-VL-7B as our base due to its strong performance across numerous vision-language tasks. For fullparameter fine-tuning, we apply DPO on a total of 29,012 samples from the critique-VQA dataset, where the α was set to 0.1. During data preprocessing, the sequence length is limited to 1024 tokens. The hyperparameters are configured as follows: the preference parameter β is set to 0.1, the batch size to 2, and the learning rate to 5e-5. We use a cosine learning rate decay strategy with a 1% warm-up period. Distributed training is conducted using DeepSpeed [44] for 5 epochs, with training performed on four A100 GPUs, requiring approximately 2.5 GPU hours per epoch.

11. Evaluation Hyperparameters for experiments.

In this section, we will list out the hyperparameters we choose for evaluation.

For the Qwen2-VL-7B and DeepSeek-VL-7B, we set the generation parameters as follows: max_new_tokens to

Table 6. Token consuming analysis of Critic-V across benchmarks.

Benchmark	Average of token count	standard deviation of token count	
MathVista 40.64		51.42	
MMBench	50.26	41.54	
MMStar	39.39	44.96	
MMT-Bench	84.43	86.93	
ScienceQA	84.64	18.11	
RealWorldQA	30.29	10.70	
SEED	41.50	28.58	
MathVerse	43.13	34.76	

1024, top_p to 0.001, top_k to 1, temperature to 0.01, and repetition_penalty to 1.0. For LLaVA-v1.5-7B, we apply a different set of parameters geared toward deterministic generation. Specifically, we set do_sample to False, temperature to 0, max_new_tokens to 512, top_p to None, num_beams to 1, and enabled use_cache to enhance efficiency. The η for controlling the update of Reasoner's prompt was set to 1.0, which indicates a full concatenation of old prompt with new critique.

12. Token Consumption

We explore the token consumption of Critic-V across different benchmarks as shown in Table 6.

13. Visualization of Training Process

In this section, we show the entire training process by several visual aids. You can find them in Figure 6, Figure 7 and Figure 8. We can obviously discover that our method convergence well experimentally.



Figure 6. Training loss vs. training steps.



Figure 7. Training accuracy vs. training steps.



Figure 8. Learning rate vs. training steps.

14. Our critique-VQA Dataset Example

In this section, we show three examples in Figure 9, Figure 10 and Figure 11 sampled from critique-VQA dataset.



Figure 9. A math example. Fake Answer indicates the answer is inserted some errors by GPT-40.



Figure 10. A real-world example of public market signage. Fake Answer indicates the answer is inserted some errors by GPT-40.



Figure 11. A driving car example. Fake Answer indicates the answer is inserted some errors by GPT-40.

15. Details of Training data and Benchmarks for Evaluation

In this section, we list some details of our training data and benchmarks for evaluation, as Table 7 and Table 8 shows.

Table 7. Details of training set. Number of tokens counted.

Part	Max length	Min length	Avg Length
Question	679	41	181.96
Chosen Critique	714	5	60.48
Reject Critique	1048	5	49.32

Table 2	8 T	Details	of	evaluation	henchmarks
1 auto o	0. L	Juans	U1	cvaluation	ocheminarks.

Benchmark	Description	#samples	
MathVista	Multimodal Math QA	1000(testmini)	
MMBench	Multimodal QA	4329	
MMStar	Multimodal QA	1500	
MMT-Bench	Multimodal QA	3127	
RealWorldQA	Multimodal QA	764	
ScienceQA	Multimodal/Text Scientific QA	4241	
SEED Multimodal QA		14233	
MathVerse	Multimodal Math QA	3940	



Figure 12. The prompt template used to inject errors in our dataset.



Figure 13. The prompt template used to detect errors in our dataset.

16. The Prompt for critique-VQA Generation

The prompt designed to inject errors into the original VQA dataset is presented in Figure 12, while the prompt used to detect errors is shown in Figure 13.