

# Corvid<sup>🐦</sup>: Improving Multimodal Large Language Models Towards Chain-of-Thought Reasoning

## Supplementary Material

### A. Supplementary Training Details

In Tab. 10, we provide detailed hyperparameter settings and time costs for the three-stage training of Corvid. Unless otherwise specified, all training and inference are conducted using 8 NVIDIA A800 (80G) GPUs by default.

Configuration	Stage 1	Stage 2	Stage 3
Batch size	256	256	128
Peak learning rate	1e-3	2e-5	2e-6
Learning rate schedule	Cosine	Cosine	Cosine
Learning rate warm-up ratio	0.03	0.03	0.03
Weight decay	0	0	0
Epoch	1	1	3
Optimizer	AdamW	AdamW	AdamW
Float precision	bfloat16	bfloat16	bfloat16
Deepspeed configuration	zero2	zero3	zero3
Training modules	GateMixer	GateMixer, LLM	LLM
Data Size	1M	1M	320K
Training hours	~11	~28	~20

Table 10. Training hyperparameter setting.

### B. Additional Experiment Results

#### B.1. Comparison with o1-Like MLLMs

In ??, we compare our models against o1-like MLLMs on various benchmarks, including MMStar [2], MMB [12], MMVet [26], MathVista (MathV) [15], AI2D [9], and Halusion [7], using their benchmark metrics computed with official implementations. Here, Corvid-o1<sup>†</sup>, LLaVA-o1 [23], and LlamaV-o1 [20] utilize the same baseline Llama-3.2-11B-Vision-Instruct [17]. Results in the table showcase that Corvid-o1-8B surpasses existing o1-like MLLMs on multiple benchmarks, particularly outperforming llamaV-o1 and Mulberry-o1-7B [24] on MathVista by 10.5 and 14.5 points, respectively. Additionally, Corvid-o1 achieves the best overall performance across all benchmarks. These results highlight the effectiveness of Corvid-o1, establishing it as a competitive MLLM that exceeds existing o1-like MLLMs with similar parameter sizes.

#### B.2. Additional Evaluation on VRC-Bench

We additionally evaluate our model on VRC-Bench [20], which is specifically designed for multimodal step-by-step reasoning tasks. The results in Tab. 11 show that Corvid-o1 achieves leading accuracy in final answers but exhibits limited performance on reasoning steps. This is because, compared to Llava-CoT and LlamaV-o1, Corvid-o1’s reasoning traces do not strictly adhere to the annotated multi-step struc-

ture in VRC-Bench. It tends to generate more streamlined and simplified reasoning processes rather than following the predefined step-by-step format.

Model	Llama-3.2 Vision [17]	Mulberry [24]	LLaVA-o1 [23]	LlamaV-o1 [20]	Corvid-o1 (Ours)
Final Answer	48.40	51.90	54.09	56.49	61.90
Steps	58.37	63.86	66.21	68.93	63.93

Table 11. Comparison with o1-like MLLMs on VRC-Bench.

#### B.3. Influence of $\alpha$ on Self-Verification

In the proposed self-verification strategy,  $\alpha$  is a weighting factor used to trade-off the cross-modal representation similarity  $\mathcal{S}$  and the model confidence  $\mathcal{C}$  for the final decision  $\mathcal{SC}$ . This relationship is formally expressed as:

$$\mathcal{SC} = (1 - \alpha)\mathcal{S} + \alpha\mathcal{C}.$$

To analyze the influence of  $\alpha$ , we conduct ablation studies by varying  $\alpha$  from 0.0 to 1.0 with a step size of 0.1. Figure 5 illustrates the relationship between model performance and the weighting factor  $\alpha$  in our self-verification strategy. As  $\alpha$  increases from 0.0 to 0.7, accuracy rises significantly from 48.6 to a peak of 55.6, demonstrating the advantage of incorporating model confidence into the final answer selection. Beyond  $\alpha = 0.7$ , performance gradually declines, suggesting that overweighting confidence relative to cross-modal similarity degrades effectiveness. The optimal value ( $\alpha = 0.7$ ) indicates that while both components contribute meaningfully to verification performance, a configuration that slightly prioritizes confidence yields superior results.

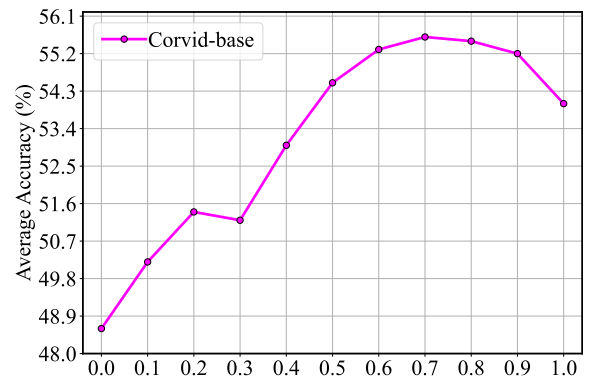



Figure 5. Influence of  $\alpha$  on our self-verification strategy.



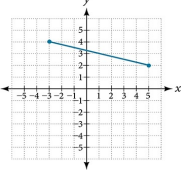
As shown in the diagram, it is known that the circumference of a circle is 6.28cm, A is the center of the circle, and BC are two points on the circle. If the radius of the circle is 1cm, what is the area of triangle ABC in  $\text{cm}^2$ ?

A.  $1 \text{ cm}^2$  B.  $1.5 \text{ cm}^2$  C.  $0.5 \text{ cm}^2$  D.  $6.2 \text{ cm}^2$   
E. No correct answer

**Corvid Response**

<Thought process>:  
S1: Since the circumference of the circle is 6.28cm, we have  $2\pi r = 6.28\text{cm}$ . Therefore,  $r = 1\text{cm}$ .  
S2: Since A is the center of the circle, and BC are two points on the circle.  
S3: Thus, the area of triangle ABC is  $0.5 \times r \times r = 0.5 \times 1 \times 1 = 0.5\text{cm}^2$ .  
S4: Therefore, the answer is C.

<Answer>: C



Find the distance between the two endpoints using the distance formula. The two end points of the line are (-3, 4) and (5, 2), respectively. Round to three decimal places.

How many degrees south of east is the new course that the plane is ordered to fly?

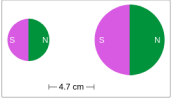
**Corvid Response**

<Thought process>:  
S1: The distance between two points is found using the distance formula:  
$$\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$
  
S2: Substitute the coordinates of the two points into the formula:  $\sqrt{(5 - (-3))^2 + (2 - 4)^2}$   
S3: Evaluate the expression:  $\sqrt{8^2 + 2^2} = \sqrt{64 + 4} = \sqrt{68}$ .  
S4: Evaluate to three decimal places: 8.246. So the answer is 8.246.

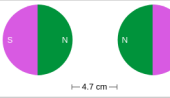
<Answer>: 8.246

(a) The CoT responses generated by Corvid in mathematical reasoning.

The images below show two pairs of magnets. The magnets in different pairs do not affect each other. All the magnets shown are made of the same material, but some of them are different sizes. Think about the magnetic force between the magnets in each pair. Which of the following statements is true?



Pair 1



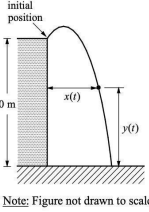
Pair 2

A. The magnitude of the magnetic force is greater in Pair 1.  
B. The magnitude of the magnetic force is greater in Pair 2.  
C. The magnitude of the magnetic force is the same in both pairs.

**Corvid Response**

<Thought process>:  
S1: The size of magnets influences the magnitude of the magnetic force.  
S2: Consider magnets that have the same shape and material.  
S3: When comparing magnets of different sizes, the larger the magnets, the greater the magnitude of the magnetic force between them.  
S4: In Pair 1, Magnet A is the same size, while in Pair 2, Magnet B is larger than in Pair 1.  
S5: Consequently, the magnitude of the magnetic force is greater in Pair 2 compared to Pair 1.  
S6: Therefore, the correct choice is option B.

<Answer>: B



An object is thrown upward into the air 10 meters above the ground. The figure shows the initial position of the object and the position at a later time. At time  $t$  seconds after the object is thrown upward the horizontal distance from the initial position is given by  $x(t)$  meters, and the vertical distance from the ground is given by  $y(t)$  meters, where  $dx/dt = 1.4$  and  $dy/dt = 4.2 - 9.8t$ , for  $t \geq 0$ . Find the angle  $\theta$ ,  $0 < \theta < \pi/2$ , between the path of the object and the ground at the instance the object hit the ground.


Note: Figure not drawn to scale. A. 0.524 B. 1.047 C. 1.475 D. 1.570

**Corvid Response**

<Thought process>:  
S1: To find the angle  $\theta$ , we can do this by solving the equation  $y(t) = 0$ . Given the equation for  $y(t)$ , we have  $4.2 - 9.8t = 0$ . Solving for  $t$ , we get  $t = 4.2/9.8 = 0.4286$  seconds.  
S2: Now, we can find the horizontal distance  $x(t)$  at this time using the equation for  $x(t)$ :  $x(t) = 1.4t = 1.4(0.4286) = 0.5954$  meters.  
S3: To find the angle  $\theta$ , we can use the tangent function:  $\tan \theta = y(t)/x(t) = 4.2/0.5954 = 7.06$ .  
S4: Now, we can find the angle  $\theta$  by taking the inverse tangent of both sides:  $\theta = \arctan(7.06) = 1.475$  radians. S5: Therefore, the answer is option C.

<Answer>: C

(b) The CoT responses generated by Corvid in science problem-solving.




You are a driving assistant. The height of current vehicle is 4 meter. Based on current image, what is the best action to take?

A. Keep driving B. Change to other road  
C. Speed up D. Turn right E. Turn left

**Corvid Response**

<Thought process>:  
The image shows a road with a bridge and a sign. The height of the current vehicle is 4 meters. Based on the given information, the best action to take is to keep driving. The bridge is high enough for the vehicle to pass underneath, and the sign does not indicate any restrictions or hazards. Therefore, the answer is A.

<Answer>: A



You are a driving assistant. Based on current image, what should you do if you want to overtake the vehicle in front of you?

A. Change to the left lane and pass.  
B. Change to the right lane and pass.  
C. Keep following the car until the road is suitable for overtaking.

**Corvid Response**

<Thought process>:  
Upon observing the image, there is a no overtaking sign on the road ahead. This suggests that overtaking is prohibited in this area. Therefore, the correct and legal action to take would be option C, keep following the car until the road is suitable for overtaking.

<Answer>: C

(c) The CoT responses generated by Corvid in visual comprehension.

Figure 6. Qualitative result of Corvid in performing mathematical reasoning, science problem-solving, and visual comprehension.

## B.4. Inference Efficiency

To evaluate the efficiency and effectiveness of our inference-time scaling strategy, we compare the average inference time and accuracy of our Corvid, LLaVA-o1, and LlamaV-o1 on the MathVista benchmark comprising 1,000 test instances. LLaVA-o1 and LlamaV-o1 utilize stage-level beam search and sentence-level beam search, respectively. Following their optimal configurations, we set the beam size to 2 for

LLaVA-o1 and 4 for LlamaV-o1. Tab. 12 shows that Corvid achieves significantly lower inference latency while maintaining higher accuracy. Specifically, Corvid-o1 reduces the inference time per instance to 11.4 seconds, yielding a  $4.4\times$  speedup over LLaVA-o1 and a  $1.35\times$  speedup over LlamaV-o1, while also surpassing both models in accuracy. This improvement is primarily attributed to Corvid-o1's more streamlined and simplified intermediate reasoning process during inference, as its inference time grows linearly with

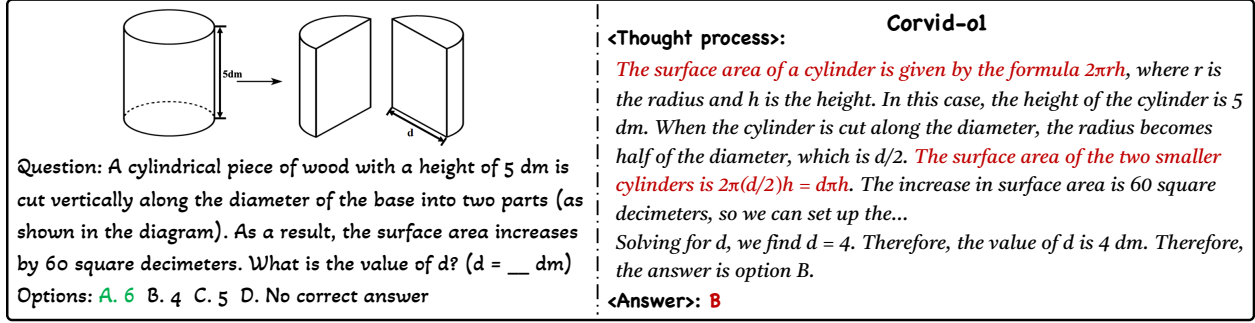


Figure 7. A failure case where Corvid-o1 fails to reason effectively due to insufficient domain-specific knowledge.

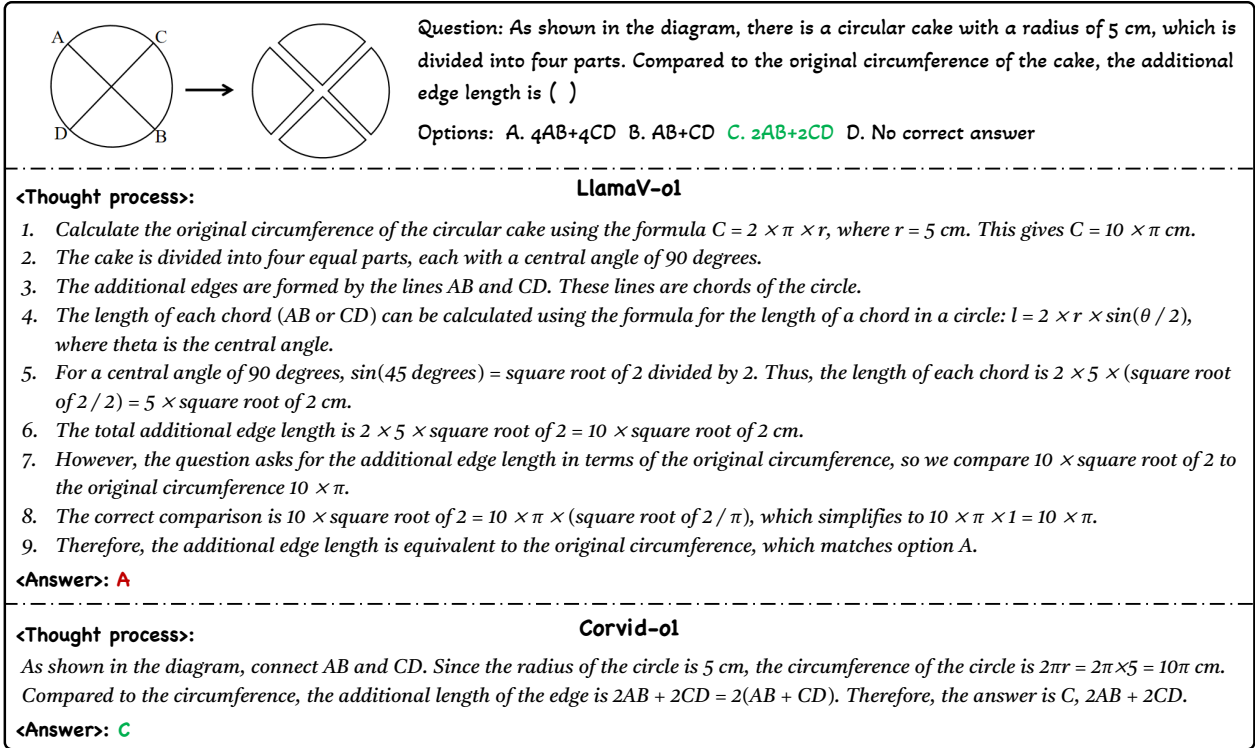


Figure 8. Qualitative comparison between LlamaV-o1 and Corvid-o1 in mathematical reasoning.

the number of generated tokens.

MLLMs	LLaVA-o1	LlamaV-o1	Corvid-o1 <sup>†</sup>	Corvid-o1-8B
Time (second)	50.6	15.4	11.4	11.3
Accuracy	56.1	54.4	61.5	72.0

Table 12. The average inference time per instance on MathVista, evaluated using a single NVIDIA A800 (80G) GPU.

## B.5. Qualitative Results

In Figure 6, we provide an intuitive understanding of Corvid’s CoT reasoning capabilities. As illustrated, when performing science and math reasoning, as well as visual comprehension, Corvid-o1 consistently generates faithful and detailed thought processes before arriving at an answer, enhancing the reliability and interpretability of its answer

and demonstrating exceptional CoT capabilities.

## B.6. Additional Failure Case

In addition to the case shown in ??, Figure 7 presents a typical failure case in mathematical reasoning, where Corvid-o1 fails to arrive at the correct answer due to insufficient domain-specific knowledge.

## B.7. Qualitative Comparison

Figures 8 to 10 visualize several qualitative comparisons between LlamaV-o1 and Corvid-o1 across tasks.

## C. MCoT-Instruct

In this section, we introduce MCoT-Instruct-287K, our high-quality multimodal CoT instruction-following dataset.

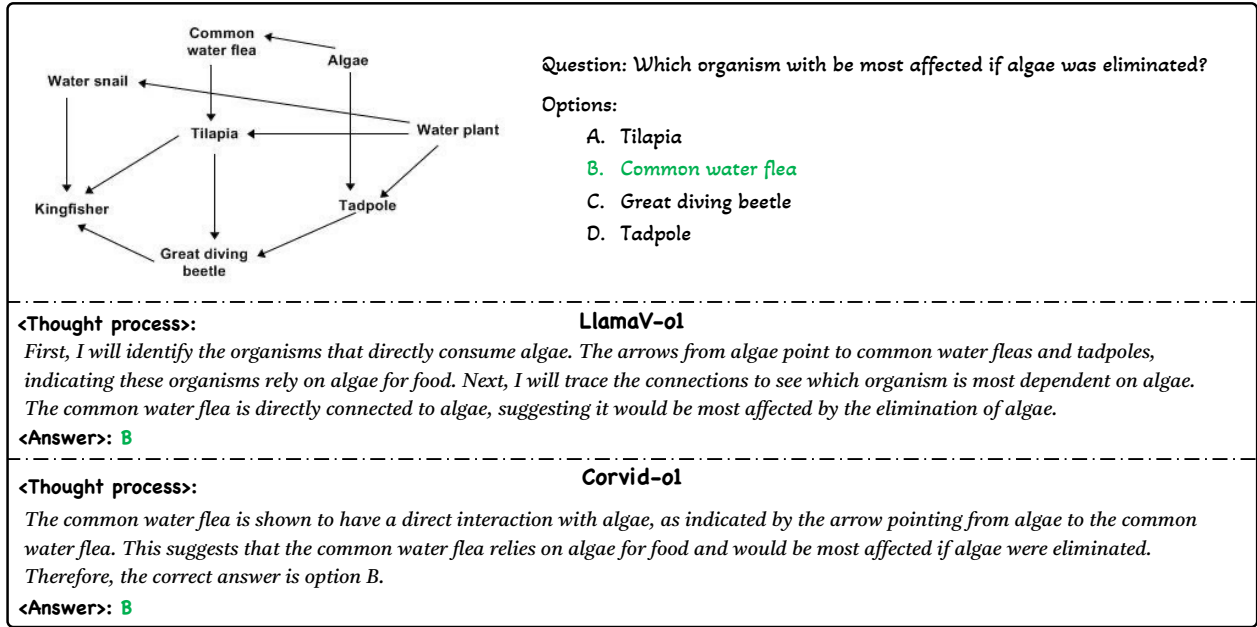


Figure 9. Qualitative comparison between LlamaV-o1 and Corvid-o1 in science problem-solving.

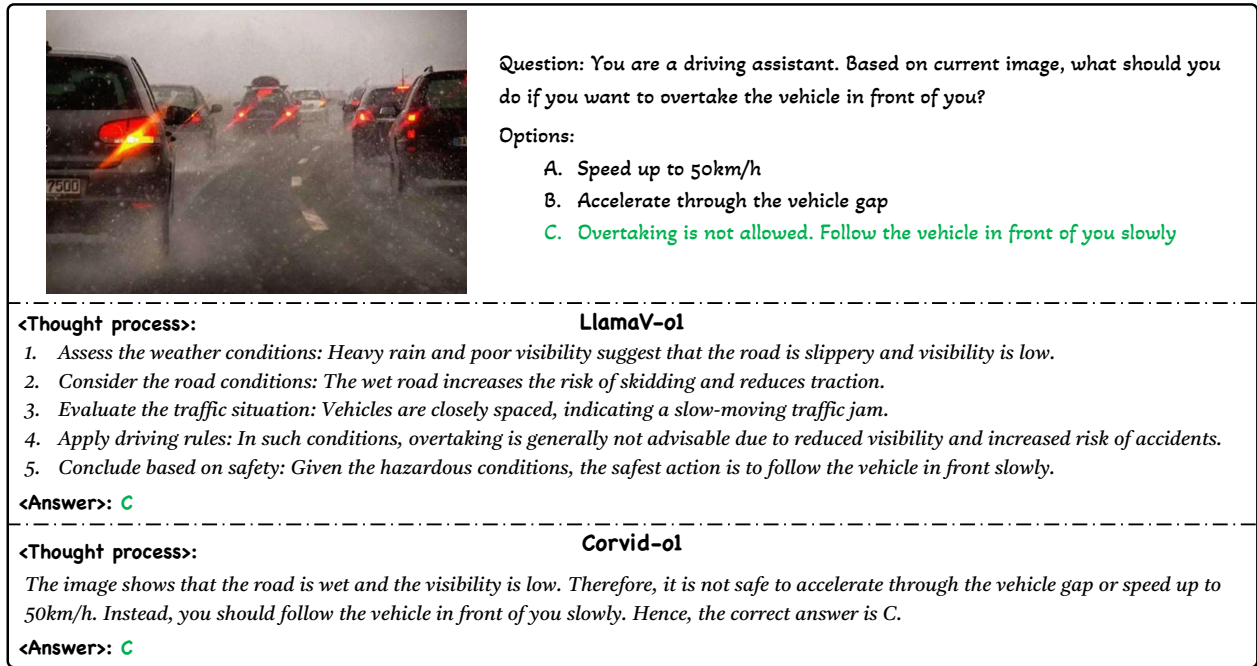


Figure 10. Qualitative comparison between LlamaV-o1 and Corvid-o1 in visual comprehension.

Specifically, we first describe its sources and then elaborate the process of improving the quality of raw CoTs.

### C.1. Source of Raw Data

As detailed in Tab. 13, we collect data from seven manually created reasoning datasets and three AI-assisted generated reasoning datasets, totaling 292K raw instances spanning diverse reasoning types and domains, to construct MCoT-Instruct. Although all datasets provide initial rationales that

serve as CoT responses, significant quality issues exist: AI-assisted generated CoTs may contain errors and duplications, while manually-created CoTs are usually brief and logically incoherent, rendering the raw data too noisy and unstandardized for effective CoT-enhancement training.

### C.2. Improving the Quality of Raw CoT

To improve the quality of raw CoTs, we separately refine and standardize the aforementioned manually created and AI-




Reasoning Type	Raw Dataset	Size
① General visual reasoning	GPT-VQA [30]	26K
② Knowledge-intensive visual reasoning	A-OKVQA [19]	18K
③ Visual Commonsense Reasoning	VCR [28]	84K
	M <sup>3</sup> CoT [3]	9K
④ Science Problem-Solving	SQA-IMG (train) [14]	8K
	ArxivQA [11]	54K
⑤ Geometric Reasoning	GeomVerse [8]	9K
	R-CoT [5]	53K
⑥ Numerical Reasoning	GeoQA [1]	7K
⑦ Mathematical reasoning	TabMWP [16]	24K

Table 13. **Raw data of MCoT-Instruct.** Here, GPT-VQA, R-CoT, and ArxivQA are the three AI-assisted generated datasets.

**Given Multimodal Input**

Which term matches the picture?

A. endotherms B. ectotherms



**Raw CoT:**  
 Endotherms regulate their temperature internally. Horses and other mammals are endotherms. On a hot day, horses can sweat to regulate their body temperature.

**Rewritten CoT:**  

Upon observing the image, we can find there is a horse.

Among the given options, endotherms are organisms that have the inherent ability to regulate their body temperature internally.

Indeed, horses and other mammals are examples of endotherms.

Therefore, the term that matches the picture is option A, endotherms.

<Answer>: A

Figure 11. **Comparison between raw and rewritten CoTs.**

generated reasoning datasets with GPT assistance through the following two steps:

- **CoT Rewriting.** As shown in Figure 12, we design a specialized prompt to instruct GPT-4o to refine these raw CoTs from manually-created datasets, enhancing their diversity and logical consistency. As demonstrated in Figure 11, the rewritten CoTs remain faithful and consistent with the given context while becoming more detailed, logically coherent, and standardized.
- **Quality Verification and Data Filtering.** To guarantee the quality of all rewritten CoTs and those from AI-assisted generated datasets, we employ GPT to evaluate free-text CoTs across three dimensions: *faithfulness*, *relevance*, and *completeness*. Inspired by the success of LLMs in automatic evaluation [4, 13], we design a base prompt, as shown in Figure 13, to instruct GPT-4o to assign an overall score (0 - 1) to each rewritten CoT and its corresponding raw CoT. The CoT with the higher score is selected as the high-quality CoT. After that, we filter out instances with an overall score below 0.6.

With these steps, we ultimately obtain 287K instance with high-quality CoT responses that are consolidated into

Benchmarks	Task Format	Metric	#Sample
MMStar [2]	multi-choice	Accuracy	1,500
MMMU [27]	multi-choice	Accuracy	900
SQA-IMG [14]	multi-choice	Accuracy	2,017
AI2D [9]	multi-choice	Accuracy	3,088
WeMath [18]	multi-choice	Accuracy	1,740
MathVista [15]	multi-choice&free-text	Accuracy	1,000
MathVerse [29]	multi-choice&free-text	Accuracy	3,940
MathVision [21]	multi-choice&free-text	Accuracy	3,040
DynaMath [31]	multi-choice&free-text	Accuracy	5,010
SEED-IMG [10]	multi-choice	Accuracy	14,232
MMT-Val [25]	multi-choice	Accuracy	31,325
RWQA [22]	multi-choice	Accuracy	1901
BLINK [6]	multi-choice	Accuracy	1,901
MMB [12]	multi-choice	Accuracy	6,666
MMVet [26]	free-text	GPT Score	218
Hallusion [7]	multi-choice	Accuracy	254

Table 14. **Summary of evaluation benchmarks.**

single-turn conversation instances of MCoT-Instruct. Notably, no testing or validation instances from any evaluation benchmark were included in this process. We used only the training split of ScienceQA for data curation, and Corvid’s in-domain performance was evaluated exclusively on its respective test set.

## D. Benchmark Details

Tab. 14 presents all evaluation benchmarks used in this work. The task formats of MathVista, MathVerse, MathVision, and DynaMath encompass both multiple-choice question answering and free-text generation, while MMVet formats tasks as free-text generation. All other benchmarks are limited to multiple-choice question answering. Each benchmark adopts accuracy as its primary metric, except for MMVet, which utilizes a GPT-based score. Notably, SQA-IMG includes human-annotated CoTs, serving as references for assessing the quality of model’s CoT responses in ??.

## References

- [1] Jiaqi Chen, Jianheng Tang, Jinghui Qin, Xiaodan Liang, Lingbo Liu, Eric Xing, and Liang Lin. GeoQA: A geometric question answering benchmark towards multimodal numerical reasoning. In *ACL*, pages 513–523, 2021. 19
- [2] Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, et al. Are we on the right way for evaluating large vision-language models? *arXiv:2403.20330*, 2024. 15, 19
- [3] Qiguang Chen, Libo Qin, Jin Zhang, Zhi Chen, Xiao Xu, and Wanxiang Che. M<sup>3</sup>CoT: A novel benchmark for multi-domain multi-step multi-modal chain-of-thought. *arXiv:2405.16473*, 2024. 19
- [4] Cheng-Han Chiang and Hung-yi Lee. A closer look into using large language models for automatic evaluation. In *EMNLP*, pages 8928–8942, 2023. 19
- [5] Linger Deng, Yuliang Liu, Bohan Li, Dongliang Luo, Liang

### System message

You are an AI assistant that can do text rewritten.

### Prompt

I want you to act as a Chain-of-Thought (CoT) Rewriter. Given a question with several options and its CoT response (i.e., the intermediate reasoning steps or rationales that lead to the correct answer), your objective is to rewrite the given CoT into a more standardized version.

#### The rewritten CoT must follow the following rules:

- 1) Keep the logic of reasoning-then-answering to ensure that the reasoning can be performed step by step.
- 2) Be faithful enough to ensure that the reasoning can accurately lead to the correct answer.
- 3) Be clear and concise, without factual errors or repeated content, and no key intermediate reasoning steps are omitted.
- 4) Do not mention or refer to the given CoT in your responses directly.

#### You can rewrite the given CoT using the following methods:

1. Improve existing reasoning steps or rationales to make the CoT more coherent and smooth.
2. Add more intermediate reasoning steps or rationales to make the CoT more specific and detailed.

Please make sure you have read and understood these instructions carefully.

Following is an exemplar:

#### ### Given question & options:

Which person is in charge?

- A. person3 [0.1231, 0.1601, 0.4296, 0.9848] is the leader
- B. person2 [0.5440, 0.1151, 0.7528, 0.9990] is in charge
- C. person3 [0.1231, 0.1601, 0.4296, 0.9848] seems to be the one in charge
- D. person3 [0.1231, 0.1601, 0.4296, 0.9848] has the higher rank

#### ### Given CoT:

person3 [0.1231, 0.1601, 0.4296, 0.9848] is walking ahead of the other two and speaking while person1 [0.3602, 0.0979, 0.5477, 0.9689] looks to be more of a worker.

#### ### Rewritten CoT:

The person who is in charge can be determined by observing their behavior and position. In this case, person3 [0.1231, 0.1601, 0.4296, 0.9848] is walking ahead of the other two individuals and is also the one who is speaking, indicating that they have a leadership role. On the other hand, person1 [0.3602, 0.0979, 0.5477, 0.9689] appears to be more of a worker. Based on these observations, it can be concluded that person3 seems to be the one in charge.

Hence, the correct answer is B.

Here is the example to be rewritten:

#### ### Given Question & Options:

{}

#### ### Given CoT:

{}

#### ### Rewritten CoT:

Now you can start to rewrite the given CoT.

Figure 12. **Prompt template of CoT rewriting.** Here, we use one instance from the VCR dataset [28] as the exemplar. For other source datasets, replace the given exemplar with data-specific instances.

- Wu, Chengquan Zhang, Pengyuan Lyu, Ziyang Zhang, Gang Zhang, Errui Ding, et al. R-cot: Reverse chain-of-thought problem generation for geometric reasoning in large multi-modal models. *arXiv:2410.17885*, 2024. 19
- [6] Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A Smith, Wei-Chiu Ma, and Ranjay Krishna. Blink: Multimodal large language models can see but not perceive. *arXiv:2404.12390*, 2024. 19
- [7] Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, et al. Hallusionbench: an advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models. In *CVPR*, pages 14375–14385, 2024. 15, 19
- [8] Mehran Kazemi, Hamidreza Alvari, Ankit Anand, Jialin Wu, Xi Chen, and Radu Soricut. Geomverse: A systematic evaluation of large models for geometric reasoning. *arXiv:2312.12241*, 2023. 19
- [9] Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. A diagram is worth a dozen images. In *ECCV*, pages 235–251, 2016. 15, 19
- [10] Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-bench: Benchmarking multimodal llms with generative comprehension. *arXiv:2307.16125*, 2023. 19
- [11] Lei Li, Yuqi Wang, Runxin Xu, Peiyi Wang, Xiachong Feng, Lingpeng Kong, and Qi Liu. Multimodal arxiv: A dataset for improving scientific comprehension of large vision-language

### System message

You are a helpful AI assistant that can evaluate the quality of free-text chain-of-thought (CoT) responses generated by a multimodal large language model (MLLM).

### Prompt

You will be provided with the input context to the MLLM (i.e., an image description, a question, and several options for the question), along with the corresponding CoT response generated by the MLLM. Your task is to evaluate the free-text CoT responses and give a final overall score (0 - 1) based on the following three perspectives:

- ❑ **Faithfulness** (0 - 1): it refers to how accurately the CoT response reflect the actual reasoning process of the MLLM. A faithful CoT response is one that genuinely represents the factors and logic the MLLM used to arrive at its answer. For example, if the MLLM generates an answer based on certain key points in the given context, a faithful CoT response would accurately describe how it picked those points and how they led to the answer. The focus of faithfulness is on the transparency and truthfulness of the explanation.
- ❑ **Relevance** (0 - 1): it measures how the CoT response aligns with and supports the answer generated by the MLLM. A consistent CoT response should logically justify the answer, demonstrating a clear and direct connection between the CoT response and the inferred answer. That is, a consistent CoT response should not only be aligned with the answer but also provide sufficient and convincing reasons for why the answer is valid.
- ❑ **Completeness** (0 - 1): it evaluates whether the CoT response provided by the MLLM encompasses all essential information and reasoning necessary to understand the MLLM's answer reasoning process. A complete CoT response should cover all critical aspects and steps of the MLLM's reasoning without omitting key details.

### Evaluation Steps:

1. Understand and analyze the provided image description, question, and options.
2. Read the MLLM's response and systematically assess the CoT response from the three perspectives of Faithfulness, Relevance, and Completeness.
3. Assign a final overall score (0 - 1) by averaging Faithfulness, Relevance, and Completeness.

*Please make sure you read and understand these instructions carefully.*

The sample to be scored:

### Image Description:

{ }

### Question & Options:

{ }

### CoT Response:

{ }

### Evaluation Form:

Answer by starting with "Scoring:" and then give the explanation of the score by "Explanation:"

- Overall:

Figure 13. Prompt template for CoT quality evaluation.

- models. *arXiv:2403.00231*, 2024. 19
- [12] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. MMBench: Is your multi-modal model an all-around player? *arXiv:2307.06281*, 2023. 15, 19
- [13] Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. G-eval: NLG evaluation using gpt-4 with better human alignment. In *EMNLP*, pages 2511–2522, 2023. 19
- [14] Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. In *NeurIPS*, pages 2507–2521, 2022. 19
- [15] Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv:2310.02255*, 2023. 15, 19
- [16] Pan Lu, Liang Qiu, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, Tanmay Rajpurohit, Peter Clark, and Ashwin Kalyan. Dynamic prompt learning via policy gradient for semi-structured mathematical reasoning. In *ICLR*, 2023. 19
- [17] Meta AI. Llama 3.2: Revolutionizing edge ai and vision with open, customizable models. <https://ai.meta.com/blog/llama-3-2-connect-2024-vision-edge-mobile-devices/>, 2024. 15
- [18] Runqi Qiao, Qiuna Tan, Guanting Dong, Minhui Wu, Chong Sun, Xiaoshuai Song, Zhuoma GongQue, Shanglin Lei, Zhe Wei, Miaoxuan Zhang, et al. We-math: Does your large multimodal model achieve human-like mathematical reasoning? *arXiv:2407.01284*, 2024. 19
- [19] Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. A-OKVQA: A benchmark for visual question answering using world knowledge. In *ECCV*, pages 146–162, 2022. 19
- [20] Omkar Thawakar, Dinura Dissanayake, Ketan More, Ritesh Thawkar, Ahmed Heakl, Noor Ahsan, Yuhao Li, Mo-

- hammed Zumri, Jean Lahoud, Rao Muhammad Anwer, et al. Llamav-o1: Rethinking step-by-step visual reasoning in llms. *arXiv:2501.06186*, 2025. 15
- [21] Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Mingjie Zhan, and Hongsheng Li. Measuring multimodal mathematical reasoning with math-vision dataset. *arXiv:2402.14804*, 2024. 19
- [22] X.AI. Grok-1.5 vision preview. <https://x.ai/blog/grok-1.5v>, 2024. 19
- [23] Guowei Xu, Peng Jin, Li Hao, Yibing Song, Lichao Sun, and Li Yuan. Llava-o1: Let vision language models reason step-by-step. *arXiv:2411.10440*, 2024. 15
- [24] Huanjin Yao, Jiaxing Huang, Wenhao Wu, Jingyi Zhang, Yibo Wang, Shunyu Liu, Yingjie Wang, Yuxin Song, Haocheng Feng, Li Shen, et al. Mulberry: Empowering mllm with o1-like reasoning and reflection via collective monte carlo tree search. *arXiv:2412.18319*, 2024. 15
- [25] Kaining Ying, Fanqing Meng, Jin Wang, Zhiqian Li, Han Lin, Yue Yang, Hao Zhang, Wenbo Zhang, Yuqi Lin, Shuo Liu, et al. Mmt-bench: A comprehensive multimodal benchmark for evaluating large vision-language models towards multitask agi. *arXiv:2404.16006*, 2024. 19
- [26] Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv:2308.02490*, 2023. 15, 19
- [27] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruochi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *CVPR*, pages 9556–9567, 2024. 19
- [28] Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. From recognition to cognition: Visual commonsense reasoning. In *CVPR*, pages 6720–6731, 2019. 19, 20
- [29] Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan Lu, Kai-Wei Chang, Peng Gao, et al. Mathverse: Does your multi-modal llm truly see the diagrams in visual math problems? *arXiv:2403.14624*, 2024. 19
- [30] Zhiyuan Zhao, Linke Ouyang, Bin Wang, Siyuan Huang, Pan Zhang, Xiaoyi Dong, Jiaqi Wang, and Conghui He. Mllm-dataengine: An iterative refinement approach for mllm. *arXiv:2308.13566*, 2023. 19
- [31] Chengke Zou, Xingang Guo, Rui Yang, Junyu Zhang, Bin Hu, and Huan Zhang. Dynamath: A dynamic visual benchmark for evaluating mathematical reasoning robustness of vision language models. *arXiv:2411.00836*, 2024. 19