

CorrectFlow: On-the-Spot Correction for Multimodal Reasoning with Multi-Agent Collaboration

Anonymous CVPR submission

Paper ID 2259

Abstract

001

002

003

004

005

006

007

008

009

010

011

012

013

014

015

016

017

018

019

020

021

022

023

024

025

026

027

028

029

030

031

032

033

034

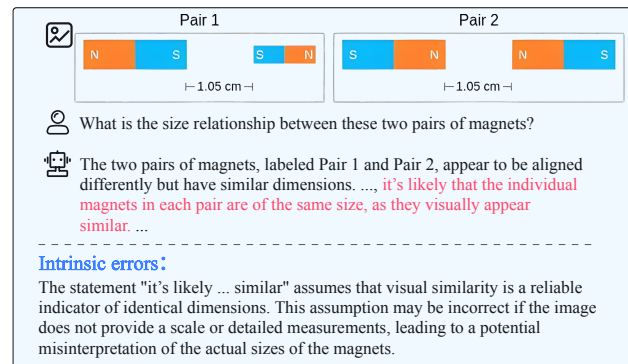
035

036

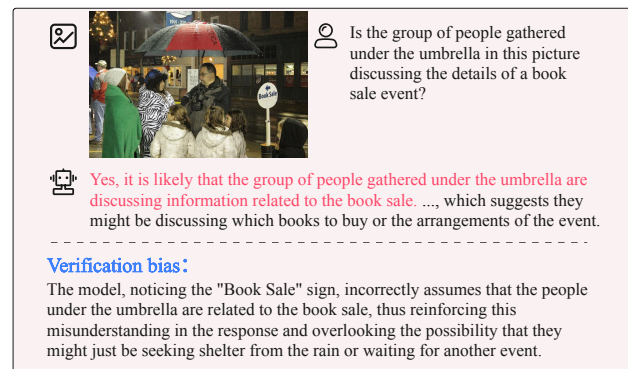
Multimodal Large Language Models (MLLMs) have shown great potential in addressing complex reasoning tasks. However, their progress is often hindered by misleading or ambiguous internal knowledge resulting from training biases. Although Chain-of-Thought (CoT) reasoning and its variants have proven effective in enhancing reasoning task performance, they often fail to correct errors in intermediate steps. Similarly, judge-based methods, while useful for validating reasoning steps, frequently struggle to identify and rectify specific mistakes. To tackle these challenges, we propose **CorrectFlow** framework, a novel approach consisting of two key agents. In **CorrectFlow**, one agent extracts knowledge from visual and textual modalities to mitigate internal knowledge errors in MLLMs, while the other serves as a multi-level agent responsible for intermediate reasoning and solution generation. The multi-level agent serves as the core of our approach and implements three key strategies: self-suspect mechanism, active-abandonment mechanism, and advanced agent takeover mechanism. Within this framework, when lower-level agents identify potential reasoning errors, they positively abandon the current thought process and delegate it to higher-level agents based on the task's complexity. By incorporating this real-time correction mechanism into a multimodal reasoning framework, **CorrectFlow** significantly enhances the accuracy and reliability of collaborative agents, particularly in complex reasoning scenarios. Comprehensive qualitative and quantitative experiments on widely used benchmarks demonstrate that **CorrectFlow** surpasses existing baseline methods, underscoring its effectiveness in improving model performance and addressing both MLLM limitations.

1. Introduction

Recent advancements in Multimodal Large Language Models (MLLMs) [1, 6, 17, 20, 29–31, 37, 48] have signifi-



(a) Example for Intrinsic Errors



(b) Example for Verification Bias

Figure 1. Examples to illustrate two limitations in MLLMs: (a) **intrinsic errors** and (b) **verification bias**.

icantly propelled the fields of perception, such as object detection [42], segmentation [15, 32], and video understanding [7, 18, 41], driving the development of specialized MLLMs tailored for these downstream tasks. However, despite their success in perception, MLLMs still face substantial challenges when it comes to complex reasoning tasks, especially those scenarios involving complex and long-horizon problem-solving.

To tackle such tasks, methods like Chain-of-Thought (CoT) reasoning [36] and its variants [10, 12, 22, 28, 44, 47]

037

038

039

040

041

042

043

044

045

046

047 have emerged. These methods break down complex reason-
048 ing processes into smaller steps and use strategies such as
049 self-correction and self-criticism to evaluate the accuracy
050 of reasoning paths [36]. While these techniques can help
051 mitigate logical errors during inference, two critical issues
052 remain largely unaddressed: (1) **intrinsic errors** that arise
053 from the MLLMs themselves, and (2) **verification bias**,
054 which stems from the limited capability of MLLMs to ac-
055 tively correct errors and instead only verify the reasoning
056 steps.

057 In Fig. 1, we illustrate both of the aforementioned errors.
058 From Fig. 1 (a), it is evident that intrinsic bias can mani-
059 fest as a form of hallucination, which is difficult to mitigate
060 in the absence of external knowledge about the problem-
061 solving object. When solving reasoning problems, these
062 intrinsic errors could easily mislead the model into an er-
063 roneous reasoning path right from the first step, ultimately
064 leading to failure in subsequent steps. Meanwhile, we also
065 show the verification bias in Fig. 1 (b). This verification bias
066 reveals another characteristic of MLLMs: they can detect
067 or suspect reasoning errors, but they do not actively correct
068 them; instead, they only verify the validity of their reason-
069 ing steps.

070 To overcome these challenges, we introduce **Correct-**
071 **Flow**, a novel framework that leverages multi-agent collab-
072 oration to effectively address both intrinsic and verification
073 biases. CorrectFlow features a two-agent system: the first
074 agent, a *knowledge extractor*, gathers objective knowledge
075 from both image and text data to provide relevant back-
076 ground information about the object being reasoned about.
077 This knowledge is essential, as MLLMs can become con-
078 fused, especially when dealing with visually similar ob-
079 jects or insufficient background context. By leveraging the
080 knowledge extractor, CorrectFlow mitigates internal errors
081 during the initial reasoning steps.

082 However, reasoning tasks often demand more than just
083 knowledge extraction; they require continuous validation
084 and correction. This is where CorrectFlow’s *multi-level*
085 *agent system* comes into play. The system introduces three
086 key strategies to address verification biases: (1) Confidence
087 Check, (2) Path Pruning, and (3) Expert Intervention. These
088 strategies ensure that the MLLM performs self-evaluation,
089 expands reasoning paths, and receives real-time corrections
090 when necessary. In CorrectFlow, the highest-level agent
091 evaluates the intermediate reasoning path from the root to
092 the current step, classifying it as accurate, erroneous, or un-
093 certain. Unlike previous Chain-of-Thought (CoT) methods
094 and their variants, CorrectFlow introduces a novel mech-
095 anism: when a lower-level agent experiences self-doubt or
096 detects potential errors in the reasoning path, a higher-level
097 agent takes over, redirecting the reasoning process. This
098 approach ensures more robust and reliable reasoning. In
099 summary, our contributions are as follows:

- We present CorrectFlow, a novel multi-agent collabora- 100
tion framework designed to overcome the limitations of a 101
single MLLM in mitigating intrinsic errors and verifica- 102
tion biases. By separating knowledge extraction from reason- 103
ing validation, CorrectFlow enhances robustness and 104
minimizes internal reasoning errors. 105
- Three core strategies are proposed: (1) Confidence Check 106
which stimulates the MLLM’s ability to self-evaluate; (2) 107
Path Pruning to facilitate the expansion of thought paths; 108
and (3) Expert Intervention for providing real-time cor- 109
rection for reasoning paths. These strategies work collec- 110
tively to ensure reliable and accurate reasoning. 111
- CorrectFlow pioneers a dynamic escalation mechanism 112
that enables lower-level agents to transfer control to 113
higher-level agents upon identifying potential reasoning 114
errors, thereby surpassing traditional passive validation 115
methods. This active intervention leads to more refined 116
and robust reasoning outcomes. 117
- We extensively evaluate CorrectFlow on public bench- 118
mark datasets, demonstrating its superior performance 119
compared to existing baseline methods, effectively ad- 120
dressing both intrinsic biases and verification limitations 121
in multimodal reasoning tasks. 122

2. Related Work 123

Multi-modal Large Language Model. Since the advent 124
of large language models (LLMs), their remarkable suc- 125
cess across numerous language-based applications has in- 126
spired the development of multimodal large language mod- 127
els (MLLMs). These models aim to bridge the gap between 128
vision and language modalities, enabling richer understand- 129
ing and reasoning across both domains. In early research, 130
MLLMs are regarded as a special way to extend the capa- 131
bilities of LLMs to handle diverse tasks and modalities, by 132
connecting specialized vision models. These models mainly 133
include MiniGPT [3, 48], VisualChatGPT [37], Hugging- 134
GPT [30], LMDrive [29], and MM-REACT [38], which 135
integrate LLMs with vision models to facilitate complex in- 136
teractions between visual and textual information. Recently, 137
the focus of MLLMs has shifted towards aligning visual 138
and language representations more effectively. This has 139
been accomplished through extensive training on datasets 140
consisting of image-caption pairs or image-question dia- 141
logues. Two main effective approaches have been pro- 142
posed. The first approach, LLaVA [20], trains an MLP 143
projector to map image tokens to a representation space 144
aligned with pre-trained LLMs, fostering effective modal- 145
ity integration. The second approach, BLIP-2 [17], uti- 146
lizes a query transformer (Q-Former) to learn image embed- 147
dings by employing learnable queries after extracting im- 148
age features. Besides the model architecture, a two-stage 149
training strategy has been explored and become a popu- 150
lar approach for MLLMs [1, 6, 31, 48]. In the first stage, 151

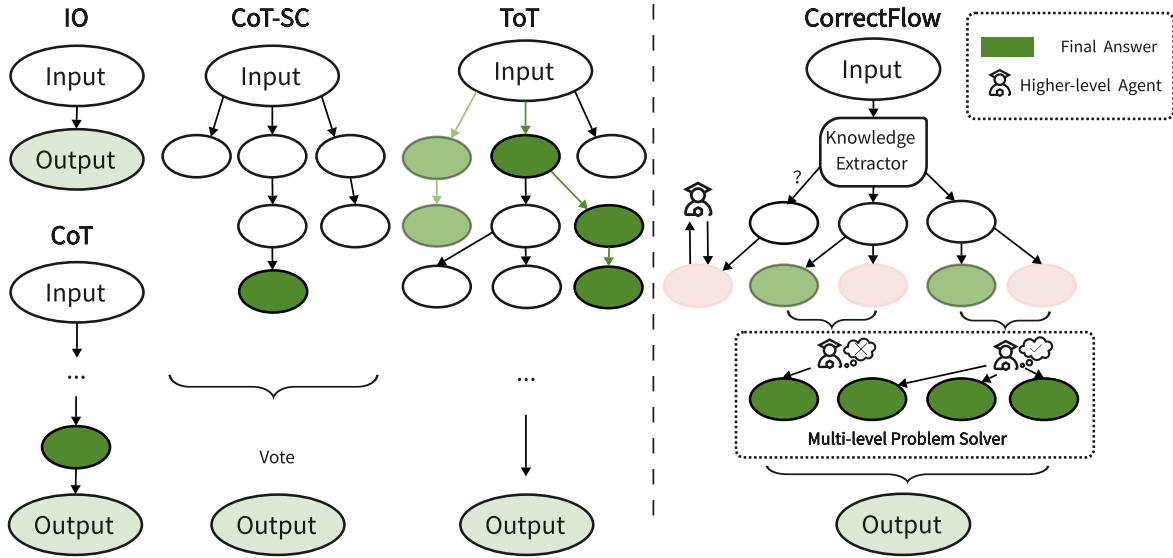


Figure 2. CorrectFlow is a multi-agent collaboration framework that enhances the robustness of MLLMs by separating knowledge extraction from reasoning validation, thereby mitigating intrinsic errors and verification biases. It introduces three core strategies—Confidence Check, Path Pruning, and Expert Intervention—that collectively ensure reliable and accurate reasoning. Compared to CoT and its variants, CorrectFlow features a dynamic escalation mechanism, allowing lower-level agents to transfer control to higher-level agents upon detecting reasoning errors, resulting in more refined and robust outputs.

152 the models undergo pre-training using large-scale image-
 153 caption datasets, laying the foundation for cross-modal un-
 154 derstanding. The second stage focuses on refining align-
 155 ment between modalities using question-answering triplets
 156 to ensure nuanced understanding and precise reasoning.
 157 With the model structure and training strategies, MLLMs
 158 have achieved promising performance in various perception
 159 tasks, including fine-grained localization [15, 32], such as
 160 object detection [42], video understanding [7, 18, 41], and
 161 image generation [13, 27]. Although MLLMs have shown
 162 promising results in perception tasks, they still face signif-
 163 icant challenges in reasoning tasks, which stem not only
 164 from limitations in their perception capabilities but also
 165 from biases inherent in the models themselves, leading to
 166 misunderstandings.

167 **CoT Reasoning in LLMs and MLLMs.** Recent stud-
 168 ies have proven using Chain-of-Thought (CoT) reason-
 169 ing to improve problem-solving skills. CoT prompts encour-
 170 age LLMs to express intermediate reasoning steps, which
 171 considerably enhances their reasoning ability. Studies such
 172 as [36] and [14] have demonstrated that simple prompting
 173 techniques or a few detailed examples can significantly en-
 174 hance the reasoning performance of LLMs in both zero-
 175 shot and few-shot scenarios. The type methods mainly cur-
 176 rent research focuses on optimizing these methods through
 177 a more refined selection of examples based on factors like
 178 similarity, diversity, and complexity [10, 22, 28, 44], while
 179 also incorporating structured approaches, including pro-

180 gramming [5], problem decomposition [12, 47], and ratio-
 181 nale calibration [33].

182 Similar to LLMs, Chain-of-Thought (CoT) prompt-
 183 ing has also shown significant effectiveness in enhanc-
 184 ing the performance of multimodal large language models
 185 (MLLMs). For example, [45] leveraged visual inputs to
 186 generate relevant rationales, thereby improving the model’s
 187 reasoning capabilities. [46] approached the problem by
 188 breaking down questions into sub-questions and utilizing
 189 answers from a visual question answering (VQA) model
 190 to develop rationales. In addition, [40] directed the model
 191 to solve complex questions involving multiple image inputs
 192 by assessing similarities and differences across the images.
 193 Moreover, [24] was a pioneer in using LLMs to generate
 194 scene graphs, subsequently using these models to formulate
 195 answers.

196 3. CorrectFlow

197 Here we present **CorrectFlow**, an innovative zero-shot
 198 prompting approach that leverages a multi-agent framework
 199 to enhance the reasoning capabilities of Multimodal Large
 200 Language Models (MLLMs) in tackling complex tasks.
 201 CorrectFlow enables zero-shot learning solely through
 202 prompts, bypassing the need for annotated data for fine-
 203 tuning. The core idea is to coordinate multiple MLLM
 204 agents to dynamically correct erroneous reasoning paths,
 205 thereby broadening and deepening the reasoning process.

206 Briefly reviewing MLLMs, these models utilize a pre-

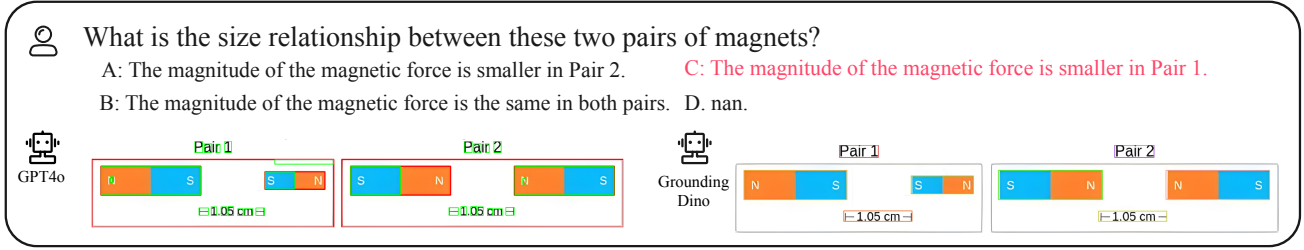


Figure 3. Open-world detectors often struggle to accurately identify relevant objects in images.

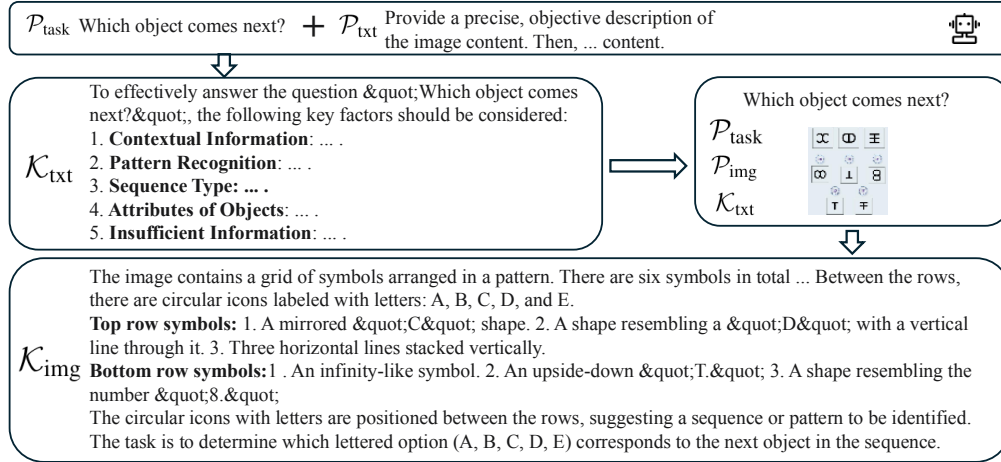


Figure 4. The pipeline of our knowledge extractor.

207 trained vision encoder $\phi_{\mathbf{w}}(\cdot)$, parameterized by \mathbf{w} , to convert
 208 an image \mathbf{I} into an embedding, and a language encoder
 209 $\psi_{\mathbf{o}}$, parameterized by \mathbf{o} , to encode the task prompt $\mathcal{P}_{\text{task}}$
 210 (e.g., a question or caption request). These embeddings are
 211 then fed into a pretrained language model f_{θ} , parameterized
 212 by θ , to generate a response \mathcal{R} :

$$\mathcal{R} = f_{\theta}(\phi_{\mathbf{w}}(\mathbf{I}), \psi_{\mathbf{o}}(\mathcal{P}_{\text{task}})).$$

214 As the vision encoder $\phi_{\mathbf{w}}(\cdot)$ has been pretrained on paired
 215 image-text descriptions, it maps visual input into a shared
 216 text embedding space, enabling the LLM f_{θ} to reason seam-
 217 lessly within a unified embedding space.

218 3.1. Knowledge Extractor

219 As shown in Fig. 1 (upper), MLLMs often display inter-
 220 nal errors when reasoning through complex and incomplete
 221 tasks. These errors likely stem from training biases that favor
 222 certain reasoning paths, leading MLLMs to overlook
 223 parts of the user’s problem and produce incomplete or incor-
 224 rect conclusions. Additionally, MLLMs commonly struggle
 225 to accurately detect and relate all relevant objects in an im-
 226 age, further complicating reasoning.

227 To address this, a straightforward solution would be
 228 to extract all objects relevant to the user’s query to help
 229 MLLMs understand object relationships in context. How-
 230 ever, in practice, open-world detectors frequently fail to

231 identify these objects accurately as illustrated in Fig. 3.
 232 Even when successfully detected, linking these objects to
 233 the user’s question presents a major challenge, as incorrect
 234 initial relationships between objects established during the
 235 initial reasoning step can lead to cascading errors in reason-
 236 ing.

237 To mitigate these challenges, we introduce a dedicated
 238 “knowledge extractor” agent, designed to provide MLLMs
 239 with objective contextual knowledge drawn from the image
 240 and task prompt. Fig. 4 illustrates this process. This agent
 241 gathers supplementary information to guide the reasoning
 242 process, enabling a more comprehensive understanding of
 243 the user’s question. The process begins with extracting key
 244 information from both the task prompt $\mathcal{P}_{\text{task}}$ and image \mathbf{I} .
 245 Given $\mathcal{P}_{\text{task}}$ and a textual key point generation prompt \mathcal{P}_{txt} ,
 246 we first derive textual key points \mathcal{K}_{txt} directly from $\mathcal{P}_{\text{task}}$.
 247 Next, we use \mathcal{K}_{txt} and an image key point generation prompt
 248 \mathcal{P}_{img} to identify objective facts \mathcal{K}_{img} within the image \mathbf{I} , thus
 249 isolating each modality to avoid cross-modal interference:

$$\begin{aligned} \mathcal{K}_{\text{txt}} &= f_{\theta}(\psi_{\mathbf{o}}(\text{cat}(\mathcal{P}_{\text{task}}, \mathcal{P}_{\text{txt}}))), \\ \mathcal{K}_{\text{img}} &= f_{\theta}(\phi_{\mathbf{w}}(\mathbf{I}), \psi_{\mathbf{o}}(\text{cat}(\mathcal{P}_{\text{img}}, \mathcal{K}_{\text{txt}}, \mathcal{P}_{\text{task}}))), \end{aligned} \quad (1)$$

251 where the operation `cat` denotes concatenating the inputs.

252 This knowledge extractor agent systematically identifies
 253 the query’s objects and extracts objective attributes related

254	to them from both the image and text. This enriched context	306
255	allows the MLLM to holistically interpret the problem, fo-	307
256	cus on the full image rather than isolated objects or par-	308
257	tial conclusions. By integrating this knowledge, the MLLM	309
258	can better align its reasoning steps with the intended so-	310
259	lution, significantly reducing the chance of errors in initial	311
260	reasoning steps.	312
261	3.2. Multi-level Problem Solver	313
262	At the core of our CorrectFlow framework is a novel multi-	314
263	level problem-solving approach that enhances the robust-	315
264	ness and reliability of reasoning in complex tasks. Inspired	316
265	by automated curriculum learning, we introduce a hierar-	317
266	chical structure consisting of multiple low-level problem-	318
267	-solving agents and a high-level problem-solving agent. In	319
268	this framework, the low-level agents act as “students,” while	320
269	the high-level agent functions as a “teacher,” guiding the	321
270	reasoning process in real time. This relationship allows	322
271	students to leverage the teacher’s higher expertise, ensur-	323
272	ing more accurate and efficient reasoning. The motiva-	324
273	tion behind this setup is intuitive: when students encounter	325
274	uncertainty, they can request help from the teacher to re-	326
275	solve potential errors. This teacher-student dynamic mirrors	327
276	real-world learning, where a teacher helps correct misun-	328
277	derstandings, fostering more accurate decision-making and	329
278	deeper understanding.	330
279	In our multi-level problem-solving framework, we adopt	331
280	an O1-inspired approach, combining self-correction and	332
281	cross-validation. At each reasoning step, multiple indepen-	333
282	dent paths are generated, allowing the system to evaluate	334
283	the state from different perspectives or knowledge sources.	335
284	This multi-perspective validation enables a thorough re-	336
285	assessment, comparing outcomes across paths. When sig-	337
286	nificant discrepancies are detected, the system analyzes and	338
287	adjusts the reasoning to correct biases or errors. While the-	339
288	oretically effective, MLLMs often fail to identify the root	340
289	causes of mistakes, allowing errors to propagate along the	341
290	reasoning path.	342
291	To overcome this limitation, we introduce three piv-	343
292	otal strategies at each step of reasoning: (1) Confidence	344
293	Check, (2) Path Pruning, and (3) Expert Intervention. These	345
294	strategies work together to ensure robustness, accuracy, and	346
295	adaptability in the reasoning process.	347
296	Confidence Check. It enables the reasoning agent to self-	348
297	assess the validity of each reasoning step. When an agent	349
298	detects potential flaws or inconsistencies, it generates a	350
299	“self-suspect” signal. This signal prompts further investi-	351
300	gation or assistance from a higher-level agent. This process	352
301	is inspired by iterative questioning, where doubts lead to	353
302	deeper scrutiny, ensuring more reliable conclusions.	354
303	Path Pruning. It discards the unreliable path before errors	355
304	can propagate when an agent is uncertain about the correct-	
305	ness of a reasoning path. This ensures the system only pro-	
	gresses along valid reasoning paths, preventing the system	
	from getting stuck or moving forward with flawed reason-	
	ing.	
	Expert Intervention. When a self-suspect signal or error	
	is detected, control is escalated to the high-level agent, the	
	“teacher,” which performs a more thorough analysis of the	
	reasoning process. The high-level agent evaluates the rea-	
	soning path and decides on the next action:	
	• Correct: If the reasoning path is validated, no further ac-	
	tion is needed.	
	• Wrong: If a logical error is identified, teacher agent cor-	
	rects the reasoning path based on prior steps.	
	• Uncertain: If teacher agent is unsure, “Path Pruning” is	
	triggered to discard the uncertain path.	
	In Appendix, we provide a detailed figure to summarize the	
	reasoning steps in our multi-level problem-solving agents.	
	These mechanisms create a robust feedback loop that	
	continuously refines the reasoning process, minimizing er-	
	ror propagation. The dynamic interplay between low-level	
	agents (students) and the high-level agent (teacher) creates	
	an adaptable, self-correcting system that is both efficient	
	and reliable. This makes it ideal for complex problem-	
	-solving tasks where precision is critical. Our multi-level	
	problem-solving approach bridges the gap between theoret-	
	ical advancements and real-world applications, empowering	
	agents to tackle challenging tasks with confidence and ac-	
	curacy.	
	4. Experiment	
	4.1. Implementation Details	
	GPT-4O. The architectural and pretraining details of GPT-	
	4O [26] are not publicly available. Nevertheless, we use	
	GPT-4O as the MLLM backbone due to its state-of-the-art	
	language reasoning capabilities. This allows us to evalu-	
	ate the performance of our proposed method on an LMM	
	with advanced reasoning skills, providing insights into its	
	effectiveness in solving complex multi-step problems. In	
	addition to using GPT-4O, we also conducted experiments	
	with other (MLLMs) as our base models. Detailed results	
	of these additional experiments are provided in the supple-	
	mentary materials.	
	4.2. Multimodal Reasoning Benchmarks	
	The implementation of CorrectFlow has undergone rig-	
	orous evaluation using several benchmark datasets, in-	
	cluding MME, MathVista [23], BLINK [9], MMStar [4],	
	CCBench [21], and RealWorldQA [43]. These benchmarks	
	are specifically designed to assess the multimodal percep-	
	tion and reasoning capabilities of large multimodal lan-	
	guage models (LMMs). Both MME and MathVista feature	
	different splits that evaluate general visual perception and	
	reasoning. For instance, MME includes perception tasks	

Table 1. Performance Comparison on Mathvista dataset.

Method	Overall	SCI	TQA	NUM	ARI	VQA	GEO	ALG	GPS	MWP	LOG	FQA	STA
LLaVA-OneVision-72B (SI)	66.9	64.8	63.3	51.4	61.2	54.2	75.3	70.8	77.4	77.4	21.6	62.1	71.4
InternVL2-Llama3-76B	65.6	63.1	66.5	41.7	62.0	49.2	66.1	65.8	67.8	75.8	32.4	67.3	76.7
Ovis1.5-Gemma2-9B	65.6	64.8	60.1	50.7	66.3	54.7	62.8	58.7	63.5	87.1	13.5	62.8	74.1
InternVL2-40B	64.0	60.7	63.3	41.7	64.9	58.7	56.5	56.9	57.2	71.0	21.6	68.4	76.7
NVLM-D-72B	63.9	66.4	68.4	40.3	50.7	41.9	76.2	73.0	78.4	65.1	16.2	63.9	71.8
InternLM-XComposer2.5	63.7	55.7	55.7	43.8	64.0	53.1	62.8	56.9	62.0	83.3	16.2	63.2	73.4
Ovis1.5-Llama3-8B	63.0	63.1	60.8	48.6	65.7	57.5	61.1	57.3	61.1	79.6	21.6	58.0	67.1
POINTS-Qwen2.5-7B	63.0	61.5	61.4	49.3	61.5	58.1	70.7	66.5	72.6	71.0	13.5	54.3	63.5
POINTS-Yi-1.5-9B	63.0	61.5	59.5	46.5	61.5	55.9	70.3	66.2	72.1	73.1	10.8	55.8	65.8
LLaVA-OneVision-7B	62.3	65.6	60.8	45.1	57.5	47.5	68.6	64.1	70.2	76.9	16.2	56.9	66.1
Claude3.5-Sonnet	61.6	75.4	74.1	31.2	53.5	45.8	58.6	61.6	57.7	59.1	35.1	69.5	77.7
RBDash-v1.2-72B	61.6	59.8	65.8	40.3	53.3	41.9	68.2	66.9	69.7	74.2	24.3	57.2	69.1
Qwen2-VL-7B	61.4	66.4	63.3	41.0	58.9	57.0	51.0	51.6	51.0	66.1	27.0	68.0	73.8
GPT-4o (0806, high)	62.7	71.3	75.3	42.4	56.9	48.0	65.7	68.3	65.9	68.3	32.4	58.7	69.1
CorrectFlow (GPT-4o _(0806, high))	67.0	72.1	73.4	49.3	63.2	55.9	69.5	70.1	70.2	72.0	37.8	64.7	74.4

Table 2. Performance Comparison of reasoning task on MME dataset.

Method	Overall	Code Reasoning	Numerical Calculation	Text Translation	Commonsense Reasoning
Qwen-VL-Max-0809	723.9	177.5	170.0	200.0	176.4
InternVL2-Llama3-76B	658.6	152.5	185.0	162.5	158.6
NVLM-D-72B	655.7	160.0	162.5	162.5	170.7
GPT-4o (0513, low)	719.3	182.5	170.0	192.5	174.3
LLaVA-OneVision-72B	583.9	145.0	177.5	100.0	161.4
GPT-4o (0513, high)	696.1	177.5	147.5	192.5	178.6
Qwen-VL-Plus-0809	633.9	157.5	125.0	200.0	151.4
InternVL2-40B	572.1	137.5	117.5	170.0	147.1
JT-VL-Chat	608.2	145.0	170.0	132.5	160.7
Qwen-VL-Max	576.1	132.5	107.5	192.5	143.6
CongRong	521.8	97.5	100.0	185.0	139.3
MiniCPM-V-2.6	597.9	155.0	117.5	177.5	147.9
GPT-4o (0806, high)	696.4	185.0	147.5	192.5	171.4
CorrectFlow (GPT-4o _(0806, high))	766.1	185.0	200.0	188.57	192.5

Table 3. Performance Comparison of perception task on MME dataset.

Method	Overall	OCR	Artwork	Color	Count	Existence	Landmark	Position	Posters	Scene
Qwen-VL-Max-0809	1585.5	177.5	156.2	190.0	170.0	200.0	183.5	155.0	189.1	164.2
InternVL2-Llama3-76B	1572.4	147.5	173.2	180.0	180.0	195.0	179.8	173.3	188.4	164.2
NVLM-D-72B	1586.2	185.0	141.8	190.0	170.0	200.0	179.5	168.3	187.1	164.5
GPT-4o (0513, low)	1562.7	192.5	144.0	180.0	190.0	195.0	175.5	145.0	192.2	148.5
LLaVA-OneVision-72B	1570.5	162.5	153.2	185.0	170.0	200.0	178.8	178.3	183.7	159.0
GPT-4o (0513, high)	1546.2	192.5	145.2	185.0	185.0	185.0	182.0	133.3	191.2	147.0
Qwen-VL-Plus-0809	1513.3	155.0	150.0	180.0	158.3	180.0	185.0	160.0	182.0	163.0
InternVL2-40B	1565.0	162.5	170.0	188.3	180.0	190.0	180.2	153.3	189.5	151.2
JT-VL-Chat	1535.0	117.5	161.5	185.0	170.0	195.0	185.0	173.3	184.7	163.0
Qwen-VL-Max	1528.6	177.5	150.2	168.3	160.0	190.0	191.0	140.0	187.8	163.8
CongRong	1576.8	177.5	151.0	176.7	175.0	195.0	187.2	168.3	171.1	175.0
MiniCPM-V-2.6	1519.4	192.5	149.0	168.3	160.0	195.0	177.5	146.7	177.9	152.5
GPT-4o (0806, high)	1550.3	200.0	139.5	178.3	190.0	195.0	189.2	113.3	193.5	151.5
CorrectFlow (GPT-4o _(0806, high))	1540.7	192.5	148.3	190	180	195.0	148.3	153.3	190.4	143.0

356 that assess an LMM’s ability to identify instances and under-
 357 stand instance attributes, as well as higher-order reason-
 358 ing tasks such as scene understanding and instance inter-
 359 action. MathVista, on the other hand, contains complex
 360 mathematical problems, often requiring extensive inference
 361 steps. We evaluate our method on MME, excluding the in-

stance identification task, and on the entirety of MathVista.
 Additionally, we use the reasoning sets of BLINK, MM-
 Star, CCBench, and RealWorldQA to further evaluate our
 approach, focusing on the LMMs’ ability to provide de-
 tailed, long-form answers to visual questions.

362
 363
 364
 365
 366

Table 4. Performance Comparison on MMBenc benchmark.

Method	Code Reasoning	Numerical Calculation	Text Translation	Commonsense Reasoning
CoT	185.0	192.5	177.5	176.4
CoT-SC	192.5	192.5	185.0	180.7
SoT	177.5	132.5	177.5	172.1
ToT	192.5	200	177.5	172.5
CorrectFlow	185	200	188.57	192.5

Table 5. Performance Comparison on RealWorldQA benchmark.

Method	Overall
Qwen2-VL-72B	76.7
GPT-4o (0513, high)	75.4
Qwen-VL-Max-0809	74.2
LLaVA-OneVision-72B	73.9
LLaVA-OneVision-72B (SI)	73.7
Molmo-72B	73.7
InternVL2-Llama3-76B	72.7
Ovis1.6-Gemma2-9B	70.7
Qwen-VL-Plus-0809	70.1
InternVL2-40B	70.1
LLaVA-OneVision-7B	69.9
NVLM-D-72B	69.9
OmChat-v2.0-13B	69.8
Step-1.5V	69.7
LLaVA-OneVision-7B (SI)	69.5
GPT-4o (0806, high)	76.5
CorrectFlow (GPT-4o _(0806, high))	77.3

367

4.3. Baseline

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

386

387

In our experiments, we compared our proposed CorrectFlow methodology with two prompting baselines. The first baseline aimed to evaluate the added value of our method to pretrained LMMs [2, 8, 11, 16, 19] by applying the model to the benchmark without any prompt engineering as shown in Table 1, Table 2, and Table 3. The second baseline utilized a zero-shot (ZS) Chain-of-Thought (CoT) prompting method to assess the benefits of CorrectFlow compared to a state-of-the-art (SOTA) CoT [35] approach. The ZS-CoT method involves two main steps: (i) given the input question and text, the reasoning prompt "Let's think step-by-step." is appended after the question to guide the model in generating reasoning for an answer, and (ii) since the answer is implicitly embedded in the generated reasoning, the second step involves passing the image, question, generated reasoning, and an answer extraction phrase to produce the response in the desired format. We also compared CorrectFlow to recent SOTA multimodal CoT prompting methods, including COT-SC [34], SoT [25], and ToT [39], on the reasoning split of the MME benchmark, as summarized in Table 4.

388

4.4. Result

389

390

391

392

393

394

395

396

397

398

399

Results are presented in Table 1, Table 2, and Table 3. One notable advantage of our method is its significant improvement in performance on several multimodal reasoning tasks, including the complex MathVista benchmark. We demonstrate that applying CorrectFlow to GPT-4O outperforms the base models across various benchmarks, highlighting the effectiveness of our approach. Figure 5 provides specific examples where CorrectFlow enhances performance over the baselines, as well as instances where it still encounters challenges. Additional results can be found in the supplementary materials.

400

401

402

403

404

405

406

407

408

Multimodal Reasoning Tasks. CorrectFlow outperforms the baselines in the reasoning test category across RealWorldQA, BLINK, MMStar, and CCBench in Table 5, Table 6, Table 7 and Table 8. From these datasets, we observe that CorrectFlow significantly improves performance on complex tasks involving mathematical reasoning and logical understanding, particularly those requiring extended reasoning steps. Notably, our method also proves effective on Chinese datasets. These results provide strong evidence

Table 6. Performance Comparison on Blink benchmark.

Method	Multi-view Reasoning	Spatial Relation
Qwen-VL-Max-0809	40.6	88.1
Gemini-1.5-Pro	53.4	79.7
Phi-3.5-Vision	48.1	69.2
Gemini-1.5-Flash	57.1	77.6
InternVL2-26B	42.9	84.6
Yi-Vision	48.1	82.5
MiniCPM-V-2.6	55.6	81.1
LLaVA-OneVision-7B	54.1	80.4
LLaVA-Next-Interleave-7B	44.4	71.3
GPT-4o (0806, high)	45.1	82.5
CorrectFlow (GPT-4o _(0806, high))	47.4	83.2

Table 7. Performance Comparison on MMstar benchmark.

Method	Logical Reasoning	Math
Qwen-VL-Max-0809	72.4	76.0
Qwen2-VL-72B	72.4	72.8
InternVL2-Llama3-76B	72.4	75.2
LLaVA-OneVision-72B	68.8	74.4
LLaVA-OneVision-72B (SI)	67.2	72.0
Step-1.5V	68.4	64.4
InternVL2-40B	69.2	70.0
JT-VL-Chat-V3.0	69.6	76.8
GPT-4o (0513, high)	72.0	66.4
NVLM-D-72B	68.8	70.8
Molmo-72B	65.2	60.8
GPT-4o (0806, high)	72.0	67.6
CorrectFlow (GPT-4o _(0806, high))	73.6	72.4

that our approach enhances LMMs' long inference capabilities in general multimodal reasoning tasks.

We further conducted a comparative evaluation against Chain-of-Thought (CoT) and its variants, including CoT-SC, SoT, and ToT, within the reasoning category of the

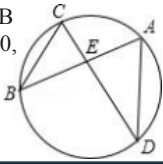
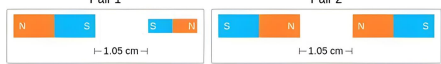

MathVista	MMBench
<p>Q: As shown in the figure, the two chords AB and CD in the circle intersect at E, $\angle D = 35.0$, $\angle AEC = 105.0$, then $\angle C = ()$</p> <p>CorrectFlow: 70° ✓</p> <p>GPT4o: 35° ✗</p> 	<p>Q: What is the size relationship between these two pairs of magnets?</p>  <p>CorrectFlow: The magnitude of the magnetic force is smaller in Pair 1. ✓</p> <p>GPT4o: The individual magnets in each pair are of the same size. ✗</p>
<p>Q: Is the number of tiny gray bicycles that are on the left side of the brown metal sedan greater than the number of things that are to the left of the tiny green bicycle?</p> <p>CorrectFlow: No ✗ Answer: Yes</p> 	<p>Q: The image shows a python code. Is the output of the code '11'?</p> <pre>a = [1, 2, 6, 9] b = max(a) + a[1] print(b)</pre> <p>CorrectFlow: No ✗</p> <p>Answer: Yes</p>

Figure 5. Comparative visualization of our method’s performance on the MathVista and MMBench datasets.

Table 8. Performance Comparison on CCbench benchmark.

Method	Sketch Reasoning
InternVL2-Llama3-76B	91.1
InternVL2-8B	91.1
InternVL2-1B	86.7
BlueLM-V-3B	91.1
Step-1.5V	91.1
Qwen-VL-Max-0809	88.9
MMAIaya2	91.1
Qwen2-VL-72B	86.7
CongRong	91.1
GPT-4o (0806, high)	88.9
CorrectFlow (GPT-4o _(0806, high))	92.2

Table 9. Effects of Knowledge Extractor (KE) and Multi-level Problem Solver (MPS) on the MME benchmark.

Method	Code Reasoning	Numerical Calculation	Text Translation	Commonsense Reasoning
w/o-KE	185.0	200.0	188.5	176.4
w/o-MPS	185.0	192.5	185.0	180.7
Our	185.0	200.0	188.57	192.5

MMBench benchmark. Table 4 provides a detailed performance comparison across various reasoning tasks, such as code reasoning, numerical calculation, text translation, and commonsense reasoning. This analysis demonstrates the advantages of CorrectFlow in handling complex reasoning tasks, which can be attributed to its real-time correction mechanism during agent collaboration.

Multimodal Perception Tasks. Table 3 presents the experimental results for perception tasks. From the table, we observe that CorrectFlow has minimal impact on perception task performance.

4.5. Ablation Study

We conducted a comprehensive ablation study on reasoning tasks in MME-Bench using our GPT-4o-CorrectFlow model. The study highlights the effectiveness of our knowledge extractor (KE) and multi-level problem solver (MPS),

as presented in Table 9. Without KE, the model’s performance dropped significantly in the commonsense category due to internal errors. Similarly, without MPS, relying solely on methods like COT-SC, the performance in several categories deteriorated, attributed to the lack of a thorough consideration of the reasoning path. For more ablation results, please refer to the supplementary materials.

4.6. Visualization Analysis

Figure 5 presents sample outputs from our method. On the left, we highlight successful cases of CorrectFlow, demonstrating its effectiveness in accurately handling complex reasoning tasks through agent collaboration. On the right, we display failure cases, offering insights into the current limitations and potential areas for improvement of our approach. For additional qualitative visualizations and a detailed analysis, please refer to the supplementary materials.

5. Conclusion

Our CorrectFlow offers a robust solution for addressing the intrinsic limitations of single MLLMs in multimodal reasoning tasks. By introducing a two-agent framework that separates knowledge extraction from reasoning validation, CorrectFlow significantly enhances accuracy and reliability. The implementation of core strategies such as the self-suspect mechanism, active abandonment, and advanced agent takeover enables dynamic intervention and escalation, effectively reducing reasoning errors and overcoming verification biases. Our extensive evaluations on several public benchmark datasets show that CorrectFlow outperforms existing methods, paving a road in the pursuit of dependable multimodal reasoning systems.

Limitation. CorrectFlow has extra computational overhead due to error correction, may impacting efficiency in rapid-response system. Future work will optimize these costly processes.

464

References

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

- [1] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023. 1, 2
- [2] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *CoRR*, abs/2308.12966, 2023. 7
- [3] Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechun Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhoseiny. Minigt-v2: large language model as a unified interface for vision-language multi-task learning. *arXiv preprint arXiv:2310.09478*, 2023. 2
- [4] 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 preprint arXiv:2403.20330*, 2024. 5
- [5] Wenhui Chen, Xueguang Ma, Xinyi Wang, and William W Cohen. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *arXiv preprint arXiv:2211.12588*, 2022. 3
- [6] Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, et al. Pali: A jointly-scaled multilingual language-image model. In *The Eleventh International Conference on Learning Representations*, 2023. 1, 2
- [7] Yukang Chen, Shengju Qian, Haotian Tang, Xin Lai, Zhijian Liu, Song Han, and Jiaya Jia. Longlora: Efficient fine-tuning of long-context large language models. In *The Twelfth International Conference on Learning Representations*, 2024. 1, 3
- [8] Zhe Chen, Jiannan Wu, Wenhui Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. *CoRR*, abs/2312.14238, 2023. 7
- [9] 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. In *Computer Vision - ECCV 2024 - 18th European Conference, Milan, Italy, September 29-October 4, 2024, Proceedings, Part XXIII*, pages 148–166. Springer, 2024. 5
- [10] Yao Fu, Hao Peng, Ashish Sabharwal, Peter Clark, and Tushar Khot. Complexity-based prompting for multi-step reasoning. In *The Eleventh International Conference on Learning Representations*, 2023. 1, 3
- [11] Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang, Yuxiang Huang, Weilin Zhao, Xinrong Zhang, Zhen Leng Thai, Kai Zhang, Chongyi Wang, Yuan Yao, Chenyang Zhao, Jie Zhou, Jie Cai, Zhongwu Zhai, Ning Ding, Chao Jia, Guoyang Zeng, Dahai Li, Zhiyuan Liu, and Maosong Sun. Minicpm: Unveiling the potential of small language models with scalable training strategies. *CoRR*, abs/2404.06395, 2024. 7
- [12] Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. Decomposed prompting: A modular approach for solving complex tasks. In *The Eleventh International Conference on Learning Representations*, 2023. 1, 3
- [13] Jing Yu Koh, Daniel Fried, and Russ R Salakhutdinov. Generating images with multimodal language models. *Advances in Neural Information Processing Systems*, 36, 2024. 3
- [14] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213, 2022. 3
- [15] Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Reasoning segmentation via large language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9579–9589, 2024. 1, 3
- [16] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer. *CoRR*, abs/2408.03326, 2024. 7
- [17] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023. 1, 2
- [18] Yanwei Li, Chengyao Wang, and Jiaya Jia. Llama-vid: An image is worth 2 tokens in large language models. In *European Conference on Computer Vision*, pages 323–340. Springer, 2025. 1, 3
- [19] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. 7
- [20] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024. 1, 2
- [21] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, Kai Chen, and Dahua Lin. Mmbench: Is your multi-modal model an all-around player? In *Computer Vision - ECCV 2024 - 18th European Conference, Milan, Italy, September 29-October 4, 2024, Proceedings, Part VI*, pages 216–233. Springer, 2024. 5
- [22] 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 *The Eleventh International Conference on Learning Representations*, 2023. 1, 3
- [23] 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. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. 5
- [24] Chancharik Mitra, Brandon Huang, Trevor Darrell, and Roei Herzig. Compositional chain-of-thought prompting for large multimodal models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14420–14431, 2024. 3
- [25] Xuefei Ning, Zinan Lin, Zixuan Zhou, Zifu Wang, Huazhong Yang, and Yu Wang. Skeleton-of-thought: Prompting llms for efficient parallel generation. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. 7
- [26] OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023. 5
- [27] Shengju Qian, Huiwen Chang, Yuanzhen Li, Zizhao Zhang, Jiaya Jia, and Han Zhang. Strait: Non-autoregressive generation with stratified image transformer. *arXiv preprint arXiv:2303.00750*, 2023. 3
- [28] Ohad Rubín, Jonathan Herzig, and Jonathan Berant. Learning to retrieve prompts for in-context learning. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2022. 1, 3
- [29] Hao Shao, Yuxuan Hu, Letian Wang, Guanglu Song, Steven L Waslander, Yu Liu, and Hongsheng Li. Lmdrive: Closed-loop end-to-end driving with large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15120–15130, 2024. 1, 2
- [30] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugginggpt: Solving ai tasks with chatgpt and its friends in hugging face. *Advances in Neural Information Processing Systems*, 36, 2024. 2
- [31] Weihang Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. Cogvlm: Visual expert for pretrained language models. *arXiv preprint arXiv:2311.03079*, 2023. 1, 2
- [32] Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, et al. Visionllm: Large language model is also an open-ended decoder for vision-centric tasks. *Advances in Neural Information Processing Systems*, 36, 2024. 1, 3
- [33] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*, 2023. 3
- [34] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. 7
- [35] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022. 7
- [36] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022. 1, 2, 3
- [37] Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. Visual chatgpt: Talking, drawing and editing with visual foundation models. *arXiv preprint arXiv:2303.04671*, 2023. 1, 2
- [38] Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. Mm-react: Prompting chatgpt for multimodal reasoning and action. *arXiv preprint arXiv:2303.11381*, 2023. 2
- [39] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. 7
- [40] Daoan Zhang, Junming Yang, Hanjia Lyu, Zijian Jin, Yuan Yao, Mingkai Chen, and Jiebo Luo. Cocot: Contrastive chain-of-thought prompting for large multimodal models with multiple image inputs. *arXiv preprint arXiv:2401.02582*, 2024. 3
- [41] Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*, 2023. 1, 3
- [42] Shilong Zhang, Peize Sun, Shoufa Chen, Min Xiao, Wenqi Shao, Wenwei Zhang, Yu Liu, Kai Chen, and Ping Luo. Gpt4roi: Instruction tuning large language model on region-of-interest. *arXiv preprint arXiv:2307.03601*, 2023. 1, 3
- [43] Yi-Fan Zhang, Huanyu Zhang, Haochen Tian, Chaoyou Fu, Shuangqing Zhang, Junfei Wu, Feng Li, Kun Wang, Qingsong Wen, Zhang Zhang, Liang Wang, Rong Jin, and Tieniu Tan. Mme-realworld: Could your multimodal LLM challenge high-resolution real-world scenarios that are difficult for humans? *CoRR*, abs/2408.13257, 2024. 5
- [44] Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in large language models. In *The Eleventh International Conference on Learning Representations*, 2023. 1, 3
- [45] Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. Multimodal chain-of-thought reasoning in language models. *arXiv preprint arXiv:2302.00923*, 2023. 3
- [46] Ge Zheng, Bin Yang, Jiajin Tang, Hong-Yu Zhou, and Sibe Yang. Ddcot: Duty-distinct chain-of-thought prompting for multimodal reasoning in language models. *Advances in Neural Information Processing Systems*, 36:5168–5191, 2023. 3

- 693 [47] Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan
694 Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier
695 Bousquet, Quoc Le, et al. Least-to-most prompting en-
696 ables complex reasoning in large language models. In *The*
697 *Eleventh International Conference on Learning Representa-*
698 *tions*, 2023. 1, 3
- 699 [48] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mo-
700 hamed Elhoseiny. Minigt-4: Enhancing vision-language
701 understanding with advanced large language models. In *The*
702 *Twelfth International Conference on Learning Representa-*
703 *tions*, 2024. 1, 2