# Multimodal Mathematical Reasoning with Diverse Solving Perspective

**Anonymous ACL submission** 

#### Abstract

Recent progress in large-scale reinforcement learning (RL) has notably enhanced the reasoning capabilities of large language models 004 (LLMs), especially in mathematical domains. However, current multimodal LLMs (MLLMs) for mathematical reasoning often rely on oneto-one image-text pairs and single-solution supervision, overlooking the diversity of valid reasoning perspectives and internal reflections. In this work, we introduce MathV-DP, a novel dataset that captures multiple diverse solution trajectories for each image-question pair, fos-013 tering richer reasoning supervision. We further propose Qwen-VL-DP, a model built upon Qwen-VL, fine-tuned with supervised learning and enhanced via group relative policy opti-017 mization (GRPO), a rule-based RL approach that integrates correctness discrimination and diversity-aware reward functions. Our method emphasizes learning from varied reasoning perspectives and distinguishing between correct yet distinct solutions. Extensive experiments on the MathVista's minitest and Math-V benchmarks demonstrate that Qwen-VL-DP significantly outperforms prior base MLLMs in both accuracy and generative diversity, highlighting the importance of incorporating diverse perspectives and reflective reasoning in multimodal mathematical reasoning. We will make our data and model public available.

## 1 Introduction

037

041

Large language models (LLMs) have demonstrated remarkable abilities in reasoning tasks (Wei et al., 2022; Wang et al., 2023; Zhou et al., 2023). This has spurred significant interest in their application to solving math problems described in natural language (Luo et al., 2023; Yue et al., 2023b; Gou et al., 2023; Jiang et al., 2023). Meanwhile, a more challenging direction involves multimodal mathematical reasoning (Lu et al., 2023), where models must interpret various types of images and apply advanced logical skills to address mathematical questions with visual components. Open-source multimodal large language models (MLLMs), such as LLaVA (Liu et al., 2023) and Qwen-VL (Bai et al., 2023), have achieved strong results on visual question answering benchmarks (Guo et al., 2023). However, when it comes to intricate mathematical problems that require visual understanding, these models still lag behind close-source counterparts like GPT-4V and Gemini (OpenAI, b; Google). 042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

078

079

081

Humans frequently engage in intuitive chainof-thought (CoT) processes to address complex reasoning tasks (Ericsson and Simon, 1980). Recent research (Wei et al., 2022) has demonstrated that LLMs are capable of exhibiting similar CoT reasoning. By employing straightforward prompting strategies or fine-tuning methods (Wang et al., 2023; Hsieh et al., 2023), CoT can both boost the reasoning abilities of LLMs and increase transparency in their decision-making procedures. Notably, recent progresses, such as OpenAI o1 (OpenAI, c), have enabled LLMs to generate more elaborate internal CoT sequences. Despite these successes in natural language contexts, adapting CoT approaches for multimodal tasks remains fully unexplored. In contrast to the rich supply of textcentric CoT data used during language model training, there is a marked shortage of multimodal CoT datasets within predominantly text-based online resources (Dai et al., 2024). This scarcity constrains the development and reasoning capacity of MLLMs.

Recent advancements in large-scale reinforcement learning (RL) (Sutton et al., 1998) have significantly enhanced the reasoning capacity of LLMs especially within mathematical reasoning tasks. o1 (OpenAI, c) and DeepSeek-R1 (Guo et al., 2025) illustrate that extensive RL applied during posttraining can lead to substantial gains in complex reasoning performance, in some instances surpassing outcomes achieved via supervised fine-tuning



Figure 1: An multimodal mathematical reasoning example with alternative solutions that reaches the final answer. Existing open-source image instruction datasets containing limited solution per image-question, do not fully exploit diverse solution with reflection to enhance the multimodal mathematical reasoning capabilities of MLLMs.

(SFT) (Radford et al., 2019). There has been growing interest within the research community to adapt the rule-based RL used in DeepSeek-R1 to multimodal scenarios (Chen et al., 2025; Yang et al., 2025). These works just explore using final answer and thinking format of image instruction dataset as reward signal.

087

100

101

102

104

105

106

108

110

111

112

Furthermore, most existing MLLMs focus on pre-training and post-training by using one-to-one image-text data to improve the final answer accuracy on mathematical reasoning but neglect diverse perspective of internal thought. As shown in Figure 1, for an image-question pair, there are usually multiple reasonable inference solutions to reach the final correct answer. Constrained by limited thinking perspectives tend to derive wrong solution and answer. Existing open-source image instruction datasets for fine-tuning or reinforcement learning, containing limited solution per image-question, do not fully exploit diverse solution with reflection to enhance the multimodal mathematical reasoning capabilities of MLLMs.

To bridge the gap, we construct MathV-DP dataset involving a variety of solutions for imagequestion corresponding to a single thought solution, and train the model Qwen-VL-DP based on the Qwen-VL-7B (Bai et al., 2025; Wang et al., 2024c) through supervised fine-tuning and group relative policy optimization (GRPO) (Shao et al., 2024) as rule-based reinforcement learning. In addition, the discrimination of diverse correct solutions and the preference for different correct and incorrect solutions are introduced in the reward function. Experiments on MathVista's minitest (Lu et al., 2023) and Math-V (Wang et al., 2024a) show that learning the correctness and diversity from multiple solution perspectives significantly improves the accuracy and generation diversity of base MLLMs on multimodal mathematical reasoning.

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

153

154

155

156

157

158

159

160

161

# 2 Related Works

#### 2.1 Multimodal Reasoning

The progress of MLLMs has significantly advanced research in multimodal reasoning (Chen et al., 2024; You et al., 2023). A widely adopted strategy involves augmenting existing question-answer datasets in specialized domains to further finetune MLLMs. For answer enhancement, rationales have been either human-authored (Zhang et al., 2023b) or extracted from leading LLMs (Wang et al., 2024b; Lin et al., 2023a; Chen and Feng, 2023; Li et al., 2024). Furthermore, VPD (Hu et al., 2023) introduced a method for converting programmatic answer representations into natural language explanations. On the question side, DDCoT (Zheng et al., 2023) employed LLMs to decompose complex queries into simpler sub-questions. Math-LLaVA (Shi et al., 2024) explored raw visual information presented in images to construct more questions. To provide a more comprehensive assessment of MLLM multimodal reasoning, several benchmarks have emerged: MathVista (Lu et al., 2023), and Math-V (Wang et al., 2024a) address diverse mathematical reasoning tasks, while MMMU (Yue et al., 2023a) spans multiple disciplines. Despite these progresses, open-source MLLMs still exhibit substantial room for improvement in complex multimodal reasoning scenarios.

# 2.2 Reinforcement Learning

Reinforcement learning (RL) (Littman and Moore, 1996) represents a foundational paradigm within machine learning, wherein an agent interacts with its environment by executing actions, receiving corresponding feedback in the form of rewards, and iteratively updating its policy to optimize cumulative returns over time. Classical RL algorithms, such as Q-learning (Watkins and Dayan, 1992), have demonstrated broad applicability across domains including robotics, game playing, and autonomous systems. With the advent of LLMs (Brown et al., 2020; Radford et al., 2018), reinforcement learning from human feedback (RLHF) (Bai et al., 2022) has emerged as an essential strategy for model fine-tuning, utilizing human-annotated preference data. RLHF commonly incorporates optimization methods like proximal policy optimization (PPO) (Schulman et al., 2017) and direct preference optimization (DPO) (Rafailov et al., 2023), facilitating improved response alignment, coherence, and utility in generated outputs.

162

163

164

165

166

167

168

169

171

172

173

174

175

176

177

178

179

181

183

184

187

188

189

192

193

194

195

196

197

198

Recently, there has been a growing interest in leveraging RL to enhance the reasoning abilities of LLMs (Team et al., 2025; Guo et al., 2025; Shao et al., 2024; Luong et al., 2024), particularly within the scope of mathematical reasoning. The central approach involves designing reward functions or evaluative models that preferentially reinforce highquality reasoning steps and discourage inadequate reasoning, thereby steering the optimization process toward more organized and comprehensible reasoning patterns through RL techniques. For instance, ReST-MCTS (Zhang et al., 2024) utilizes a process reward model (PRM) to assess the correctness of individual reasoning steps within solution paths. Moreover, recent research indicates that even straightforward rule-based, outcome-level reward functions can serve as robust and informative signals during RL, as demonstrated by DeepSeek-R1 (Guo et al., 2025). DeepSeek-R1 incorporates group relative policy optimization (GRPO) (Shao et al., 2024) combined with outcome-based reward assessments, effectively advancing the reasoning proficiency of LLMs. In this work, we focus on further enhancing the reasoning capabilities of MLLMs through reinforcement learning.

# 3 Method

Our proposed method is composed of two compo-199 nents: (1) bootstrapping a substantial set of both positive and negative chain-of-thought (CoT) solu-201 tions with reflection for collected multimodal mathematical question-CoT; and (2) leveraging these 203 new sampled positive solutions, pairs of different positive solutions and pairs of positive-negative solutions to perform post-training on the under-207 lying diverse rationales and to facilitate learning discrimination and preference from identified pairs. 208 Through the data synthesis and post-training, the MLLM is progressively improved from an initial single solving perspective to a diverse state. The 211

overall framework is depicted in Figure 2.

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

247

248

249

250

251

252

253

254

255

256

257

258

259

261

262

# 3.1 Data Synthesis

In vision-language reasoning tasks, given an image I and a corresponding question q, an MLLM is expected to perform joint reasoning over both modalities to generate a rationale r, followed by deriving a final answer a. However, constructing largescale datasets comprising high-quality (I, q, r, a)remains a significant challenge, primarily due to the scarcity of well-annotated rationale data. This data bottleneck hinders the post-training enhancement of MLLM reasoning capabilities. Although MLLMs possess a rudimentary ability for CoT reasoning and self-reflection, leveraging them to generate diverse and high-quality (I, q, r, a) samples from existing multimodal mathematical datasets is difficult. Recent advancements in language models, such as DeepSeek-R1, demonstrate strong capabilities in producing coherent, reflective reasoning across extended textual contexts. Formal languages, characterized by strict syntactic and semantic rules, provide a structured representation that eliminates ambiguity and enforces logical consistency. When visual content is described using formal language, it enables language models to see and reason over image elements more effectively. In our work, we utilize DeepSeek-R1 (Guo et al., 2025) to synthesize diverse detailed reasoning chains on samples from the MultiMath-300K dataset (Peng et al., 2024). This facilitates the construction of a richer and more diverse set of cross-modal mathematical reasoning samples, culminating in our proposed 40K MathV-DP dataset. The data generation pipeline is illustrated on the left side of Figure 2.

Data Source. We adopt MultiMath-300K (Peng et al., 2024) as the primary data source for our data synthesis. This dataset is a large-scale, multimodal, multilingual, multi-level, and multi-step mathematical reasoning benchmark, encompassing a wide range of K-12 level problems. It spans nearly the entire K-12 curriculum, covering a broad spectrum of mathematical domains, including arithmetic, algebra, geometry, functions, algorithms, and more. Compared to existing multimodal mathematics datasets (e.g., Geo170K (Gao et al., 2023) and MathV360K (Shi et al., 2024)), the problems in MultiMath-300K are newly curated and do not overlap with those in previously released datasets. Each instance is paired with a descriptive image caption to support vision-language alignment, as



Figure 2: The overall flowchart of the proposed multimodal question-solution data synthesis and post-training. Post-training consists of supervised fine-tuning and rule-based reinforcement learning (GRPO) to learn diverse and reflection reasoning manner.

well as a detailed step-by-step solution. The availability of formal visual descriptions and CoT annotations in MultiMath-300K with single solution per sample makes it particularly well-suited as seed data for synthesizing diverse solutions from multiple perspectives. Specifically, we randomly selected 10K samples from them as seed data  $\mathcal{D}$ .

263

265

267

269

270

274

276

277

278

279

283

287

290

**Diverse Solutions Construction.** Given an image, we prompt large language reasoning model (i.e., DeepSeek-R1) with its formal dense caption, question and limited original solution to construct more diverse CoT data with reflection. The prompt for generating new solutions *s* is shown in Figure 3. Two correct solutions and two incorrect solutions that differ from each other are generated at once for each source sample. They are organized into three formats to constitute MathV-DP dataset involving CoT with reflection thinking, discrimination of different correct solutions and preference of solutions.

The correct solution with reflection is first taken out separately with the original image and question. The rationale before the final answer in each solution is wrapped with *<think>* and *</think>* tags as  $r_{think}$  to form a new set  $\mathcal{D}_s^+$  totaling 20K:

$$\mathcal{D}_{s}^{+} = \{ (I_{i}, q_{i}, r_{think}, a_{i}) \}_{i=1}^{|\mathcal{D}|}$$
(1)

The generated different correct solutions are then concatenated with the instruction  $Ins_1$  (i.e., "Are the solution perspectives of the two solutions dis*similar?*") to form set  $\mathcal{D}_d$  totaling 10K:

$$\mathcal{D}_d = \left\{ \left( I_i, q_i, s_{i_1}^+, s_{i_2}^+, Ins_1, 1 \right) \right\}_{i=1}^{|\mathcal{D}|}$$
(2)

For the data format of correctness preference, one of each of the correct and incorrect solutions is randomly selected and both are concatenated together as a pair in a random back-and-forth order to construct set  $D_p$  totaling 10K. Instruction  $Ins_2$ is "Is the former/later solution the correct one?":

$$\mathcal{D}_p = \left\{ \left( I_i, q_i, s_i^+, s_i^-, Ins_2, 1 \right) \right\}_{i=1}^{|\mathcal{D}|}$$
(3)

## 3.2 Post-Training

To improve the multimodal mathematical reasoning capabilities of MLLMs, we propose a twostage post-training framework comprising supervised fine-tuning followed by rule-based reinforcement learning. In this pipeline, supervised finetuning serves to stabilize the model's reasoning ability and learn diverse solving process with reflection, while the subsequent reinforcement learning phase promotes better generalization, preference of solution correctness and diversity in multimodal mathematical reasoning task.

# 3.2.1 Supervised Fine-Tuning

Specifically, we utilize  $D_s^+$  with diverse solution perspectives during the supervised fine-tuning 314

292

293

294

296

297

298

299

300

301

302

303

304

305

306

307

308

310

311

312

291

387

388

389

390

391

392

345

346

**Prompt-Solutions Generation:** [Role] You are a math expert. [Image Description] formal image description [Original Question] question [Given Solution] original solution [Task] Please change the given solution to two solutions that are different and incorrect. Then change the given solution to two solutions that have different solution paths but the same final right answer. [Requirement] For each solution, please involve complete and detailed seeking thought process with planning, reflection, and verification. Please split each whole solution with '/\*\*\*\*\*/.' The solutions should be coherent and independent, and contain no information about the correct given solution."

Figure 3: The prompt template used in our DeepSeek-R1 API for generating additional solutions with reflection for each input image description, question and original CoT solution.

stage to guide the model  $\mathcal{M}$  toward generating coherent and diverse reasoning chains with a negative 316 log-likelihood objective: 317

315

318

319

323

324

$$\mathcal{L}_{\text{SFT}} = -\sum_{(I,q,r,a)\sim\mathcal{D}_s^+} \log \mathcal{M}(r,a \mid q,I) \quad (4)$$

Supervised fine-tuning not only aligns the model's outputs with desired formats but also encourages the emergence of more sophisticated multimodal mathematical reasoning reflection behaviors. This establishes a robust foundation for the subsequent RL phase, where rule-based feedback is employed to further refine the model's reasoning abilities.

## 3.2.2 Rule-Based Reinforcement Learning

Building upon the model fine-tuned via supervised 328 fine-tuning, we further optimize its structured reasoning capabilities, output validity and diversity of solutions through a rule-based reinforcement 331 learning framework. In particular, we design three 332 reward functions and employ group relative policy 333 optimization (GRPO) (Shao et al., 2024) for policy 334 updates.

Accuracy Reward. The accuracy-based reward 336 function assesses the correctness of the MLLM's final output by extracting the predicted answer using regular expressions and comparing it against the ground truth. We regard multimodal mathematical reasoning as deterministic tasks, the model 341 is required to present the final answer in a predefined format to facilitate consistent and rule-based evaluation.

**Think Format Reward.** To enforce the explicit presence of a reasoning process, the format-based reward function mandates that the MLLM's rationale be encapsulated within predefined delimiters, i.e., *<think>* and *</think>*. Regularization is used to verify the existence and correct ordering of these markers, thereby ensuring adherence to the required output structure.

Discrimination and Preference Reward. The discrimination/preference reward function can be viewed as a binary classification task. It is used to evaluate whether the MLLM correctly distinguishes the diversity of different solutions and whether it prefers the correct solution. This reward signal facilitates the model to learn the different perspectives of the solutions and the correctness preference.

Group Relative Policy Optimization. To ensure stable training with both consistent policy updates and informative reward signals, we adopt group relative policy optimization (GRPO) as our reinforcement learning algorithm. For each token in the generated sequence, GRPO computes the loglikelihoods under the current policy  $\pi(\theta)$  and a reference policy. The ratio between these probabilities is then calculated and clipped within the interval  $[1 - \epsilon, 1 + \epsilon]$  to mitigate the risk of overly aggressive updates. The reward, normalized to serve as an advantage estimate, is subsequently incorporated into a proximal policy optimization (PPO) objective function:

$$\mathcal{L}_{\text{clip}} = -\mathbb{E}[\min(\text{ratio}_t \cdot Ad_t, \text{clipratio}_t \cdot Ad_t)], (5)$$

where  $Ad_t$  represents the advantage estimate, quantifying the relative improvement of the chosen action over the expected value under the reference policy. To further constrain the updated policy from deviating excessively from the reference distribution, a Kullback-Leibler (KL) divergence term is incorporated into the objective, scaled by a coefficient  $\beta$ . The total loss function is defined as:

$$\mathcal{L}_{\mathrm{RL}}(\theta) = -\mathbb{E}[\min\left(\mathrm{ratio}_t \cdot Ad_t, \mathrm{clipratio}_t \cdot Ad_t\right) \\ -\beta \cdot \mathrm{KL}\left(\pi_\theta(y|x), \pi_{\mathrm{ref}}(y|x)\right)]$$
(6)

GRPO employs a clipping strategy that effectively mitigates drastic changes in the policy, while the incorporation of KL regularization enforces proximity between the updated and reference policies. This dual mechanism enables stable and efficient integration of rule-based rewards, preserving training robustness throughout the optimization process.

Model		MathVista											
Widder	ALL	FQA	GPS	MWP	TQA	VQA	ALG	ARI	GEO	LOG	NUM	SCI	STA
Heuristics Baselines													
Random Chance	17.9	18.2	21.6	3.8	19.6	26.3	21.7	14.7	20.1	13.5	8.3	17.2	16.3
Frequent Guess (Lu et al., 2023)	26.3	22.7	34.1	20.4	31.0	24.6	33.1	18.7	31.4	24.3	19.4	32.0	20.9
Human	60.3	59.7	48.4	73.0	63.2	55.9	50.9	59.2	51.4	40.7	53.8	64.9	63.9
Close-Source Multimodal Large Langugae Models (MLLMs)													
Gemini 1.0 Nano 2 (Team et al., 2023)	30.6	28.6	23.6	30.6	41.8	31.8	27.1	29.8	26.8	10.8	20.8	40.2	33.5
Qwen-VL-Plus (Bai et al., 2023)	43.3	54.6	38.5	31.2	55.1	34.1	39.1	32.0	39.3	18.9	26.4	59.0	56.1
Gemini 1.0 Pro (Team et al., 2023)	45.2	47.6	40.4	39.2	61.4	39.1	45.2	38.8	41.0	10.8	32.6	54.9	56.8
Claude 3 Haiku (Anthropic, 2024)	46.4	-	-	-	-	-	-	-	-	-	-	-	-
GPT-4V (OpenAI, b)	49.9	43.1	50.5	57.5	65.2	38.0	53.0	49.0	51.0	21.6	20.1	63.1	55.8
GPT-4o (OpenAI, a)	63.8	-	-	-	-	-	-	-	-	-	-	-	-
OpenAI o1 (OpenAI, c)	73.9	-	-	-	-	-	-	-	-	-	-	-	-
Open-Source Multimodal Large Langugae Models (MLLMs)													
mPLUG-Owl-7B (Ye et al., 2023)	22.2	22.7	23.6	10.2	27.2	27.9	23.6	19.2	23.9	13.5	12.7	26.3	21.4
miniGPT4-7B (Zhu et al., 2023)	23.1	18.6	26.0	13.4	30.4	30.2	28.1	21.0	24.7	16.2	16.7	25.4	17.9
LLaVAR-13B (Zhang et al., 2023a)	25.2	21.9	25.0	16.7	34.8	30.7	24.2	22.1	23.0	13.5	15.3	42.6	21.9
InstructBLIP-7B (Dai et al., 2024)	25.3	23.1	20.7	18.3	32.3	35.2	21.8	27.1	20.7	18.9	20.4	33.0	23.1
LLaVA-13B (Liu et al., 2023)	26.1	26.8	29.3	16.1	32.3	26.3	27.3	20.1	28.8	24.3	18.3	37.3	25.1
SPHINX-V1-13B (Lin et al., 2023b)	27.5	23.4	23.1	21.5	39.9	34.1	25.6	28.1	23.4	16.2	17.4	40.2	23.6
LLaVA-1.5-13B (Liu et al., 2024)	27.7	23.8	22.7	18.3	40.5	30.2	25.3	26.4	22.8	21.6	26.4	35.3	23.6
OmniLMM-12B (OpenBMB, 2024)	34.9	45.0	17.8	26.9	44.9	39.1	23.1	32.3	20.9	18.9	27.8	45.9	44.2
SPHINX-V2-13B (Lin et al., 2023b)	36.7	54.6	16.4	23.1	41.8	43.0	20.6	33.4	17.6	24.3	21.5	43.4	51.5
G-LLaVA-13B (Gao et al., 2023)	-	-	56.7	-	-	-	-	-	-	-	-	-	-
Math-LLaVA (Shi et al., 2024)	46.6	37.2	57.7	56.5	51.3	33.5	53	40.2	56.5	16.2	33.3	49.2	43.9
Qwen2-VL-7B (Wang et al., 2024c)	57.6	65.1	41.8	66.1	60.1	53.7	44.5	56.4	43.1	24.3	39.6	63.1	69.4
Qwen2.5-VL-7B (Bai et al., 2025)	68.2	72.5	66.8	76.9	66.7	54.3	70.1	68.7	66.9	26.9	43.0	65.7	76.1
Qwen2-VL-DP	60.9	70.7	56.4	69.8	64.6	48.7	50.9	60.4	47.2	25.4	40.3	65.6	71.1
Qwen2.5-VL-DP	70.4	72.8	72.6	77.2	68.5	54.5	71.1	69.6	69.3	27.0	43.1	66.9	77.2

Table 1: Comparison with baselines on the testmini set of MathVista benchmark. Baseline results are obtained from Lu et al. (2023). The best results in both the close-source and open-source MLLMs are in bold. MathVista is divided in two ways: task type or mathematical skill, and we report the accuracy under each subset.

# **4** Experiments

393

396

400

401

402

403

404

405

406

407

408

409

410

## 4.1 Model and Implementation

We utilize the Qwen2-VL (Wang et al., 2024c) and Qwen2.5-VL (Bai et al., 2025) series as our baseline architectures and focus our evaluation on the 7B parameter scale to assess the effectiveness of our proposed method. Both the projection layer and the language model parameters are trainable. Supervised fine-tuning stage is performed with a batch size of 16, a learning rate of 2e-5 over 1 epoch. During the reinforcement learning stage, we generate 4 rollouts per query with a sampling temperature of 1.0. The maximum sequence length is set to 1024 to ensure the model has sufficient capacity to produce complete reasoning solution. Both the policy and reference models are initialized from the same base model, with the reference model held frozen during RL training. The policy

model is fine-tuned using a learning rate of 1e-6 and a batch size of 4. The KL divergence regularization coefficient  $\beta$  in Eq. 6 is set to 0.04 by default. All experiments are conducted on NVIDIA H100 GPU with 80GB of memory. 411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

## 4.2 Evaluation and Metrics

We assess our model's performance in a zero-shot setting using the minitest subset of the MathVista benchmark (Lu et al., 2023). This subset comprises 1,000 items, including 540 multiple-choice problems and 460 free-response questions requiring answers in the form of integers, floating-point numbers, or lists. MathVista is designed to comprehensively evaluate the multimodal mathematical capabilities of MLLMs, covering diverse reasoning categories such as algebraic (ALG), arithmetic (ARI), geometric (GEO), logical (LOG), numeric commonsense (NUM), scientific (SCI), and statisti-

6

Model	Math-V																
	ALL	Alg	AnaG	Ari	CG	Comb	Cnt	DG	GT	Log	Angle	Area	Len	SG	Sta	Торо	TG
Heuristics Baselines																	
Human	68.8	55.1	78.6	99.6	98.4	43.5	98.5	91.3	62.2	61.3	33.5	47.2	73.5	87.3	93.1	99.8	69.0
Close-Source Multimodal Large Langugae Models (MLLMs)																	
Qwen-VL-Plus (Bai et al., 2023)	10.7	11.3	17.9	14.3	12.7	4.8	10.5	15.4	8.9	14.3	11.6	6.4	10.0	14.3	6.9	8.7	11.3
Qwen-VL-Max (Bai et al., 2023)	15.6	10.7	19.1	20.0	16.9	12.5	17.9	16.4	12.2	21.0	13.3	14.2	19.8	11.5	20.7	13.0	17.3
Gemini Pro (Team et al., 2023)	17.7	15.1	10.7	20.7	20.1	11.9	7.5	20.2	21.1	16.8	19.1	19.0	20.0	14.3	13.8	17.4	20.8
GPT-4V (OpenAI, b)	22.8	27.3	32.1	35.7	21.1	16.7	13.4	22.1	14.4	16.8	22.0	22.2	20.9	23.8	24.1	21.7	25.6
GPT-4o (OpenAI, a)	30.4	42.0	39.3	49.3	28.9	25.6	22.4	24.0	23.3	29.4	17.3	29.8	30.1	29.1	44.8	34.8	17.9
Ope	n-Sou	irce M	<i>Iultim</i>	odal .	Larg	e Lang	gugae	e Mod	dels (	(MLL	Ms)						
SPHINX-V2-13B (Lin et al., 2023b)	9.7	6.7	7.1	12.9	7.5	7.7	6.0	9.6	16.7	10.1	11.0	11.8	12.5	8.2	8.6	8.7	6.0
LLaVA-1.5-13B (Liu et al., 2024)	11.1	7.0	14.3	14.3	9.1	6.6	6.0	13.5	5.6	13.5	10.4	12.6	14.7	11.5	13.8	13.0	10.7
Math-LLaVA (Shi et al., 2024)	15.7	9.0	20.2	15.7	18.2	10.1	10.5	16.4	14.4	16.0	20.2	18.4	17.6	9.4	24.1	21.7	17.9
Qwen2-VL-7B (Wang et al., 2024c)	16.3	11.3	24.9	15.7	16.9	10.1	11.9	16.4	15.6	519.3	22.5	16.4	22.5	14.3	17.2	4.4	20.8
Qwen2.5-VL-7B (Bai et al., 2025)	25.0	22.0	29.8	32.1	19.5	18.5	16.4	22.1	11.1	25.2	29.3	27.6	28.5	22.9	34.5	17.4	22.0
Qwen2-VL-DP	17.7	15.2	20.8	20.8	20.2	12.0	7.9	20.3	21.2	16.9	19.2	19.1	23.2	14.4	13.9	17.5	20.9
Qwen2.5-VL-DP	26.9	23.3	30.8	32.2	20.6	27.3	17.4	23.9	22.9	28.6	28.9	30.9	28.8	28.7	37.9	18.4	23.2

Table 2: Performance Comparison on the Math-V benchmark with the accuracy metric across various mathmatical subjects. Baseline results are obtained from Wang et al. (2024a). The best results in both the close-source and open-source MLLMs are in bold.

cal reasoning (STA). Additionally, its questions are distributed across various subtypes, including Figure Question Answering (FQA), Geometry Problem Solving (GPS), Math Word Problem (MWP), Textbook Question Answering (TQA), and Visual Question Answering (VQA). For evaluation, we leverage GPT-4 (OpenAI, a) to extract final answers or selected choices from model responses and compute accuracy by verifying the correspondence between predicted and grounded answers. In addition, we perform evaluations on Math-V (Wang et al., 2024a). Math-V contains 3,040 visualcontext math problems curated from authentic math competitions.

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

Accuracy evaluation mainly depends on the final answer of the MLLM output, we also use the *effective semantic diversity* metric (Shypula et al., 2025) to assess the diversity of the MLLM's output solutions. For each input, model generates Kresponses  $G_i = \{g_i^1, g_i^2, \dots, g_i^K\}$ . We then adopt the following pairwise diversity score:

$$\operatorname{Div}_{\operatorname{pair}}\left(G_{i}\right) = \frac{1}{\binom{K}{2}} \sum_{1 \leq j < k \leq K} d_{\operatorname{sem}}\left(g_{i}^{j}, g_{i}^{k}\right),$$
(7)

where  $d_{\text{sem}}$  is semantic distance function. It is obtained by Sentence Transformer (Reimers and Gurevych, 2019), which is 1 if semantically dissimilar and 0 otherwise. This pairwise evaluation strategy incorporates normalization over the total number of candidate pairs, thereby ensuring robustness against fluctuations in the number of valid outputs generated for different prompts. The overall diversity of a model on the benchmark is then computed by averaging all pairwise diversity scores. 455

456

457

458

459

460

461

462

#### 5 Results and Analysis

#### 5.1 Main Comparison on Accuracy

We compare Qwen-VL-DP with other MLLMs on 463 the minitest split of the MathVista benchmark in Ta-464 ble 1. As shown in the table, open-source MLLMs 465 such as instructBLIP (Dai et al., 2024) and LLaVA-466 1.5 (Liu et al., 2023) have poor performance in mul-467 timodal mathematics, with overall accuracy lower 468 than 30%. Compared to the base model, Qwen2.5-469 VL-7B, with superior multimodal mathematical 470 ability, Qwen2.5-VL-DP achieves 70.4% overall 471 accuracy with a improvement of 2.2%. Qwen2-VL-472 DP also obtains improvement of 3.3% compared 473 with base model Qwen2-VL-7B. More surprisingly, 474 the proposed Qwen2.5-VL-DP model outperforms 475 close-source models GPT-4V and GPT-4o (Ope-476 nAI, b), even achieving comparable performance 477 to OpenAI o1 (OpenAI, c), the most powerful close-478 source MLLMs with the ability of detailed think-479 ing. The results on Math-V are shown in Table 2. 480 Qwen2.5-VL-DP demonstrates substantial perfor-481 mance gains over its base model, narrowing the gap with state-of-the-art models such as GPT-4V and GPT-40. The excellent performance of Qwen-VL-DP indicates that the high-quality data synthesis of solutions with diverse perspective and reflection is effective in improving MLLM's multimodal mathematical reasoning capabilities and performance.

#### 5.2 Comparision on Generation Diversity

The proposed Qwen-VL-DP model has demon-490 strated exceptional performance in multimodal 491 mathematical reasoning task. To assess its capabil-492 ity of generation diversity, we conduct evaluation 493 experiments using effective semantic diversity met-494 ric on MathVista's minitest subset. For each input 495 sample, the number of generated responses K is 496 taken as 3, 5, and 10 to calculate the corresponding 497 pairwise diversity score for final averaging. Ta-498 ble 3 presents comparison of the effective semantic 499 diversity among the Qwen-VL base model, the supervised fine-tuned model, the model tuned using only GRPO, and the post-training model after two stages using MathV-DP data. The results indicate that either supervised fine-tuning or reinforcement learning on MLLM using solution data with diverse perspectives can enhance the generative diversity of the base model. Through our synthesis of MathV-DP and proposed post-training, MLLM 508 can further enhance the accuracy performance of multimodal mathematical reasoning while improv-510 ing the diversity of output responses. The reason is 511 that Qwen-VL-DP has learnt diverse solution per-512 spectives after supervised fine-tuning and further 513 514 learnt the discriminative and preference of different solutions after rule-based reinforcement learning. 515

Model	Diver@3	Diver@5	Diver@10
Qwen2-VL-7B	27.64	30.18	31.33
Qwen2-VL-SFT	33.72	35.63	35.75
Qwen2-VL-GRPO	35.05	36.08	37.11
Qwen2-VL-DP	37.48	38.97	39.16
Qwen2.5-VL-7B	33.29	34.76	36.89
Qwen2.5-VL-SFT	39.46	39.78	39.81
Qwen2.5-VL-GRPO	39.02	39.49	39.73
Qwen2.5-VL-DP	40.42	41.44	41.58

Table 3: Effective semantic diversity scores for Qwen-VL models evaluated in our experiments.

We conduct ablation study across three training

paradigms: (1) supervised fine-tuning (SFT) on

# 516 5.3 Enhancements from SFT and RL

517

482

483

484

485

486

487

488

489

518



Figure 4: Accuracy of Qwen-VL model adopting different post-training strategies on MathVista.

519

520

521

522

523

524

525

526

527

528

529

530

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

our curated dataset, (2) SFT followed by GRPO, and (3) RL applied in isolation. As shown in Figure 4, MLLM by SFT demonstrates improvements on the MathVista. Applying RL to the SFT model yields further gains, suggesting that RL facilitates deeper and more varied deductive reasoning. These progressive enhancements underscore the complementary strengths of SFT and RL: while SFT provides a stable foundation by aligning the model with diverse high-quality reasoning perspectives, RL further strengthens these abilities by promoting advanced cognitive behaviors. In contrast, applying RL without prior SFT leads to suboptimal performance, likely due to the absence of a structured reasoning baseline. Overall, integrating SFT with RL emerges as an effective paradigm for enhancing the MLLM's mathematical reasoning ability.

# 6 Conclusions

In this work, we proposed MathV-DP, a novel dataset that enriched multimodal mathematical reasoning with diverse solving perspectives and reflective supervision. Building upon Qwen-VL, we introduced Qwen-VL-DP, trained via both supervised fine-tuning and group relative policy optimization (GRPO), a rule-based reinforcement learning method tailored to reward correctness, diversity, and discrimination of multiple solutions. Our experiments on MathVista's minitest and Math-V benchmarks demonstrated that incorporating diverse reasoning perspectives significantly enhanced both the accuracy and generative diversity of MLLMs. These findings highlight the importance of moving beyond one-to-one image-text supervision, advocating for a shift towards learning from multiple valid solving perspectives.

# 7 Limitations

554

557

558

560

562

563

564

566

570

571

573

574

575

578

582

583

584

585

586

588

589

590

591

592

594

595

596

599

By learning from synthetic CoT data with diverse solving perspectives and reflection, and preference data involving discrimination of solution diversity and correctness, MLLM could be enhanced in multimodal mathematical reasoning as well as generative diversity across multiple responses. Such diversity could not be controlled explicitly in a single response; a single generation tends to randomly be one of the multiple correct solution perspectives learned. In future work, our model will be guided or trained to controllably generate the expected solution perspective.

# 8 Ethics Statement

We do not envision that our work will result in any harm as defined in ethics policy. Qwen2-VL and Qwen2.5-VL base model use Apache License. For datasets, MultiMath-300K uses Apache License 2.0. The evaluation datasets use permissive Creative Commons Licenses. The intended use of these source datasets and evaluation datasets is to train and test the model's multimodal reasoning capability, which is consistent with our utilization of all these data. Our proposed MathV-DP dataset can improve the multimodal mathematical reasoning ability of the open-source Qwen-VL through post-training.

#### References

- AI Anthropic. 2024. The claude 3 model family: Opus, sonnet, haiku. Claude-3 Model Card.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A frontier large vision-language model with versatile abilities. CoRR, abs/2308.12966.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. 2025. Qwen2.5-vl technical report. <u>arXiv preprint</u> arXiv:2502.13923.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. <u>arXiv</u> preprint arXiv:2204.05862.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie 604 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind 605 Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot 607 learners. Advances in neural information processing 608 systems, 33:1877–1901. 609 610 Feng Chen and Yujian Feng. 2023. Chain-of-thought prompt distillation for multimodal named entity 611 and multimodal relation extraction. arXiv preprint 612 arXiv:2306.14122. 613 Jiaxing Chen, Yuxuan Liu, Dehu Li, Xiang An, Ziy-614 ong Feng, Yongle Zhao, and Yin Xie. 2024. Plug-615 and-play grounding of reasoning in multimodal large 616 language models. arXiv preprint arXiv:2403.19322. 617 Liang Chen, Lei Li, Haozhe Zhao, Yifan Song, and 618 Vinci. 2025. R1-v: Reinforcing super generaliza-619 tion ability in vision-language models with less than 620 \$3. https://github.com/Deep-Agent/R1-V. Ac-621 cessed: 2025-02-02. 622 Wenliang Dai, Junnan Li, Dongxu Li, Anthony 623 Meng Huat Tiong, Junqi Zhao, Weisheng Wang, 624 Boyang Li, Pascale N Fung, and Steven Hoi. 625 2024. Instructblip: Towards general-purpose vision-626 language models with instruction tuning. Advances 627 in Neural Information Processing Systems, 36. 628 K Anders Ericsson and Herbert A Simon. 1980. Verbal 629 reports as data. Psychological review, 87(3):215. 630 Jiahui Gao, Renjie Pi, Jipeng Zhang, Jiacheng Ye, Wan-631 jun Zhong, Yufei Wang, Lanqing Hong, Jianhua Han, 632 Hang Xu, Zhenguo Li, et al. 2023. G-llava: Solving 633 geometric problem with multi-modal large language 634 model. arXiv preprint arXiv:2312.11370. 635 Google. Gemini. https://gemini.google.com. 636 Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, 637 Yujiu Yang, Minlie Huang, Nan Duan, and Weizhu 638 Chen. 2023. Tora: A tool-integrated reasoning 639 agent for mathematical problem solving. CoRR, 640 abs/2309.17452. 641 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, 642 Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, 643 Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: In-644 centivizing reasoning capability in llms via reinforce-645 ment learning. arXiv preprint arXiv:2501.12948. 646 Yanyang Guo, Fangkai Jiao, Zhiqi Shen, Liqiang Nie, 647 and Mohan S. Kankanhalli. 2023. UNK-VQA: A 648 dataset and A probe into multi-modal large models' 649 abstention ability. CoRR, abs/2310.10942. 650 Cheng-Yu Hsieh, Chun-Liang Li, Chih-Kuan Yeh, 651 Hootan Nakhost, Yasuhisa Fujii, Alexander Ratner, 652 Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. 653 2023. Distilling step-by-step! outperforming larger 654 language models with less training data and smaller 655 model sizes. arXiv preprint arXiv:2305.02301. 656

- Yushi Hu, Otilia Stretcu, Chun-Ta Lu, Krishnamurthy Viswanathan, Kenji Hata, Enming Luo, Ranjay Krishna, and Ariel Fuxman. 2023. Visual program distillation: Distilling tools and programmatic reasoning into vision-language models. CoRR, abs/2312.03052.
  - Albert Qiaochu Jiang, Sean Welleck, Jin Peng Zhou, Timothée Lacroix, Jiacheng Liu, Wenda Li, Mateja Jamnik, Guillaume Lample, and Yuhuai Wu. 2023. Draft, sketch, and prove: Guiding formal theorem provers with informal proofs. In The Eleventh International Conference on Learning Representations.

664

670

673

674

675

677

678

679

704

707

710

- Lei Li, Yuqi Wang, Runxin Xu, Peiyi Wang, Xiachong Feng, Lingpeng Kong, and Qi Liu. 2024. Multimodal ArXiv: A dataset for improving scientific comprehension of large vision-language models. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics, pages 14369–14387.
- Hongzhan Lin, Ziyang Luo, Jing Ma, and Long Chen. 2023a. Beneath the surface: Unveiling harmful memes with multimodal reasoning distilled from large language models. In Findings of the Association for Computational Linguistics: EMNLP 2023, pages 9114-9128.
- Ziyi Lin, Chris Liu, Renrui Zhang, Peng Gao, Longtian Qiu, Han Xiao, Han Qiu, Chen Lin, Wenqi Shao, Keqin Chen, et al. 2023b. Sphinx: The joint mixing of weights, tasks, and visual embeddings for multi-modal large language models. arXiv preprint arXiv:2311.07575.
- ML Littman and AW Moore. 1996. Reinforcement learning: A survey, journal of artificial intelligence research 4.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Improved baselines with visual in-Lee. 2024. struction tuning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 26296–26306.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. In Advances in Neural Information Processing Systems.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2023. Mathvista: Evaluating math reasoning in visual contexts with gpt-4v, bard, and other large multimodal models. CoRR, abs/2310.02255.
- Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, Qingwei Lin, Shifeng Chen, and Dongmei Zhang. 2023. Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct. CoRR, abs/2308.09583.

Trung Quoc Luong, Xinbo Zhang, Zhanming Jie, Peng	711
Sun, Xiaoran Jin, and Hang Li. 2024. Reft: Rea-	712
soning with reinforced fine-tuning. <u>arXiv preprint</u>	713
<u>arXiv:2401.08967</u> , 3.	714
OpenAI.a. Chatgpt. https://chat.openai.com.	715
<pre>OpenAI. b. Gpt-4v(ision). https://openai.com/</pre>	716
research/gpt-4v-system-card.	717
<pre>OpenAI. c. Introducing openai o1. https://openai.</pre>	718
com/o1/.	719
OpenBMB. 2024. Large multi-modal models for strong	720
performance and efficient deployment. https://	721
github.com/OpenBMB/OmniLMM.	722
Shuai Peng, Di Fu, Liangcai Gao, Xiuqin Zhong, Hong-	723
guang Fu, and Zhi Tang. 2024. Multimath: Bridging	724
visual and mathematical reasoning for large language	725
models. <u>arXiv preprint arXiv:2409.00147</u> .	726
Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya	727
Sutskever, et al. 2018. Improving language under-	728
standing by generative pre-training.	729
Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	730
Dario Amodei, Ilya Sutskever, et al. 2019. Language	731
models are unsupervised multitask learners. <u>OpenAI</u>	732
<u>blog</u> , 1(8):9.	733
<ul> <li>Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. <u>Advances in Neural Information Processing Systems</u>, 36:53728–53741.</li> </ul>	734 735 736 737 738 739
Nils Reimers and Iryna Gurevych. 2019. Sentence-bert:	740
Sentence embeddings using siamese bert-networks.	741
In Proceedings of the 2019 Conference on Empirical	742
Methods in Natural Language Processing. Associa-	743
tion for Computational Linguistics.	744
John Schulman, Filip Wolski, Prafulla Dhariwal,	745
Alec Radford, and Oleg Klimov. 2017. Proxi-	746
mal policy optimization algorithms. <u>arXiv preprint</u>	747
<u>arXiv:1707.06347</u> .	748
<ul> <li>Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu,</li></ul>	749
Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan	750
Zhang, YK Li, Y Wu, et al. 2024. Deepseekmath:	751
Pushing the limits of mathematical reasoning in open	752
language models. <u>arXiv preprint arXiv:2402.03300</u> .	753
Wenhao Shi, Zhiqiang Hu, Yi Bin, Junhua Liu, Yang	754
Yang, See-Kiong Ng, Lidong Bing, and Roy Ka-Wei	755
Lee. 2024. Math-llava: Bootstrapping mathemati-	756
cal reasoning for multimodal large language models.	757
arXiv preprint arXiv:2406.17294.	758
Alexander Shypula, Shuo Li, Botong Zhang, Vishakh	759
Padmakumar, Kayo Yin, and Osbert Bastani. 2025.	760
Evaluating the diversity and quality of 1lm generated	761
content. arXiv preprint arXiv:2504.12522.	762

865

866

867

868

869

870

763

764

770

771

773

777

779

786

790

791

792

793

- 813 814

- 810

- 815
- 816 817

- Richard S Sutton, Andrew G Barto, et al. 1998. Reinforcement learning: An introduction, volume 1. MIT press Cambridge.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. arXiv preprint arXiv:2312.11805.
- Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, et al. 2025. Kimi k1. 5: Scaling reinforcement learning with llms. arXiv preprint arXiv:2501.12599.
- Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Mingjie Zhan, and Hongsheng Li. 2024a. Measuring multimodal mathematical reasoning with math-vision dataset. arXiv preprint arXiv:2402.14804.
- Lei Wang, Yi Hu, Jiabang He, Xing Xu, Ning Liu, Hui Liu, and Heng Tao Shen. 2024b. T-sciq: Teaching multimodal chain-of-thought reasoning via large language model signals for science question answering. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 38, pages 19162–19170.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Daviheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024c. Owen2-vl: Enhancing vision-language model's perception of the world at any resolution. arXiv preprint arXiv:2409.12191.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models. In The Eleventh International Conference on Learning Representations.
- Christopher JCH Watkins and Peter Dayan. 1992. Olearning. Machine learning, 8:279-292.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In Advances in Neural Information Processing Systems.
- Yi Yang, Xiaoxuan He, Hongkun Pan, Xiyan Jiang, Yan Deng, Xingtao Yang, Haoyu Lu, Dacheng Yin, Fengyun Rao, Minfeng Zhu, et al. 2025. R1onevision: Advancing generalized multimodal reasoning through cross-modal formalization. arXiv preprint arXiv:2503.10615.
- Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. 2023.

mplug-owl: Modularization empowers large language models with multimodality. arXiv preprint arXiv:2304.14178.

- Haoxuan You, Rui Sun, Zhecan Wang, Long Chen, Gengyu Wang, Hammad A. Ayyubi, Kai-Wei Chang, and Shih-Fu Chang. 2023. Idealgpt: Iteratively decomposing vision and language reasoning via large language models. In Findings of the Association for Computational Linguistics: EMNLP, pages 11289-11303.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi 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 Wenhu Chen. 2023a. MMMU: A massive multi-discipline multimodal understanding and reasoning benchmark for expert AGI. CoRR, abs/2311.16502.
- Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. 2023b. Mammoth: Building math generalist models through hybrid instruction tuning. CoRR, abs/2309.05653.
- Dan Zhang, Sining Zhoubian, Ziniu Hu, Yisong Yue, Yuxiao Dong, and Jie Tang. 2024. Rest-mcts\*: Llm self-training via process reward guided tree search. Advances in Neural Information Processing Systems, 37:64735-64772.
- Yanzhe Zhang, Ruiyi Zhang, Jiuxiang Gu, Yufan Zhou, Nedim Lipka, Diyi Yang, and Tong Sun. 2023a. Llavar: Enhanced visual instruction tuning for text-rich image understanding. arXiv preprint arXiv:2306.17107.
- Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. 2023b. Multimodal chain-of-thought reasoning in language models. arXiv preprint arXiv:2302.00923.
- Ge Zheng, Bin Yang, Jiajin Tang, Hong-Yu Zhou, and Sibei Yang. 2023. Ddcot: Duty-distinct chainof-thought prompting for multimodal reasoning in language models. Advances in Neural Information Processing Systems, 36:5168–5191.
- Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc V. Le, and Ed H. Chi. 2023. Least-to-most prompting enables complex reasoning in large language models. In The Eleventh International Conference on Learning Representations.
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. Minigpt-4: Enhancing vision-language understanding with advanced large language models. CoRR, abs/2304.10592.