062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

Section 2018 LATTE: Learning to Reason with Vision Specialists

Anonymous CVPR submission

Paper ID 36

Abstract

001 While open-source vision-language models perform well on 002 simple question-answering, they still struggle with complex questions that require heterogeneous vision capabili-003 ties. Unfortunately, we have yet to develop methods that 004 005 infuse fine-grained recognition, visual grounding, depth estimation, and 3D reasoning into a single vision-language 006 007 model. Instead of forcing smaller models to learn both perception and reasoning, we propose LATTE, a family of 008 vision-language models that have LeArned to Think wiTh 009 vision spEcialists. By offloading perception to state-of-the-010 art vision models, our approach enables vision-language 011 012 models to focus solely on reasoning over high-quality per-013 ceptual information. To train LATTE, we create and filter a large dataset of 273K high-quality synthetic reasoning 014 traces over perceptual outputs of vision specialists. LATTE 015 trains on this data and brings significant gains across 6 016 benchmarks covering both perception and reasoning abil-017 018 ities, compared to baselines instruction-tuned with direct answers. On the other hand, models trained by distill-019 020 ing both perception and reasoning from larger models lead to smaller gains or even degradation on some perception 021 022 tasks. Further, our method results in a 2% to 5% improve-023 ment on average across all benchmarks over the vanilla 024 instruction-tuned baseline regardless of model backbones, with gains up to 16% in MMVet. 025

1. Introduction

The landscape of real-world vision-language tasks is vast, 027 028 spanning from basic visual question answering [1] and finegrained object recognition to complex multi-step geomet-029 ric reasoning [8]. These tasks demand both perception and 030 reasoning. For instance, a user might photograph a gas 031 032 price panel and ask how much fuel they can afford within 033 a given budget (Figure 1). Solving this requires a vision-034 language model with strong perception-localizing prices via OCR-and multi-step reasoning to compute the answer. 035 While large proprietary models like GPT-40 excel due to 036 extensive data and model size scaling, smaller open-source 037 038 models still struggle [22].

To narrow the gap between large proprietary models and039smaller open-source counterparts within a reasonable bud-040get, researchers have explored distilling both perception and041reasoning from larger vision-language models [25, 28] or042specialized vision models [9]. Despite these efforts, open-043source models continue to lag behind.044

We argue that the primary reason for this lag is the per-045 ception limitations of open-source vision-language models. 046 While open-source language models have largely caught up 047 with their proprietary counterparts [2, 13], vision remains 048 a complex fusion of heterogeneous capabilities. The com-049 puter vision community has historically tackled these ca-050 pabilities separately-e.g., DepthAnything [29] for depth 051 estimation and GroundingDINO [19] for object recogni-052 tion-while unified models still lag behind [20]. Simi-053 larly, the human brain dedicates distinct regions to cat-054 egorical recognition (ventral stream) and spatial reason-055 ing (dorsal stream)[6], with the reasoning and language-056 processing frontal and temporal lobes occupying a smaller 057 volume[12]. By contrast, vision-language models remain 058 heavily skewed toward language, treating visual encoders 059 as an afterthought [4]. 060

We depart from the learning to perceive and reason paradigm to propose a new approach: learning to reason with vision specialists. Rather than expecting a small model to master both perception and reasoning, we leverage decades of advancements in computer vision by relying on specialized vision models to provide perceptual information. This allows the vision-language model to focus exclusively on acquiring perceptual information from vision specialists and reasoning over them-enabling it to 'see further by standing on the shoulders of giants.' Such a paradigm reduces the burden on models to extract low-level perceptual signals, allowing them to concentrate on higher-level reasoning while benefiting from the robust capabilities of dedicated vision specialists, which is particularly important for small open-source models because of their limited capacity to effectively learn both perception and reasoning.

To implement this paradigm, we curate high-quality training data in the form of multi-step reasoning traces that integrate perceptual information from vision specialists. We

125

126

127

128

129

130

131

134

135

136

137



Figure 1. Example outputs of LATTE vs. SoTA multi-modal large language models. Our LATTE model is able to answer challenging visual questions by reasoning over perceptual information output by vision specialists.

080 formulate the multi-step reasoning traces as LATTE-trace, 081 where each step consists of: (1) a *thought* for verbalized reasoning; (2) an action to retrieve perceptual information 082 083 from a specific vision specialist; and (3) an observation of the returned data. Since obtaining these traces at scale with 084 human annotators is costly, we develop two data engines 085 for synthetic data generation. First, we leverage GPT-4o's 086 strong multimodal reasoning and state-of-the-art vision spe-087 cialists' precise perception to generate large-scale synthetic 088 089 reasoning traces across diverse image sources, applying aggressive filtering and mixing techniques. Second, we gen-090 erate reasoning traces using Python programs and struc-091 tured reasoning templates, comparing them against GPT-092 093 generated traces to evaluate reasoning quality. In total, we produce over 1M reasoning traces across 31 datasets with 094 095 GPT-40 and handcrafted programs.

With this data, we finetune small multi-modal language
models to reason with vision specialists and evaluate our
models on 6 benchmarks covering both perception and reasoning skills. We compare our model to two types of baselines: (1) multi-modal language models trained with vanilla
instruction tuning with only direct answers; and (2) models
trained by distilling both perception and reasoning.

Finally, we highlight four major takeaways from our ex-103 periments: First, learning to reason with vision specialists 104 enables our model to outperform vanilla instruction-tuned 105 baseline by significant margins on both perception and rea-106 soning benchmarks, with an overall average gain of 6.4%. 107 By contrast, the other distillation methods lead to smaller 108 gains or even degradation in the perception performance. 109 This trend holds as we scale the training data. Second, 110 111 our method consistently outperforms the vanilla instructiontuned baseline by 2 - 5% on average across all bench-112 marks regardless of model backbones, with staggering per-113 formance gains of 10 - 20% on MMVet. Third, through 114 data ablations, we confirm that the quality of LATTE-trace 115 matters more than quantity: our best data recipe consists of 116 117 only 293K LATTE-trace which GPT-40 generated and answered correctly, and it leads to larger performance gains118than all other data recipes of larger scales (up to 2x larger or
more). Finally, programmatically-generated LATTE-trace119can hurt model performance as a result of the worse reason-
ing quality, suggesting that again that high-quality reason-
ing is crucial to the model's performance.122

2. Learning to Think with Vision Specialists

Our goal is to train vision-language models to reason about complex multi-modal tasks with the help of vision specialists. To train such models, we need reasoning traces that involve (1) invoking vision specialists and (2) reasoning over their outputs. We refer to such data as LATTE-trace. One LATTE-trace \mathcal{T} is a sequence of steps S_i , where each step consists of thought t_i , action a_i and observation o_i :

$$\mathcal{T} = (S_0, S_1, ..., S_n) = (S_i)_{i=0}^n \tag{1}$$

$$S_i = (t_i, a_i, o_i), t_i \in L, a_i \in A$$
 (2) 133

where L represents language space, and A is the action space consisting of vision specialists. Note that the model only generates t_i and a_i , which the training loss is applied on, whereas o_i is obtained from the vision specialists.

Action space. The action space A of our model consists 138 of vision tools that are either specialized vision models or 139 image processing tools. Concretely, these include OCR 140 [10], GETOBJECTS [33], LOCALIZEOBJECTS [19], ES-141 TIMATEOBJECTDEPTH [29], ESTIMATEREGIONDEPTH 142 [29], DETECTFACES [17], CROP, ZOOMIN, GETIMAGE-143 TOTEXTSSIMILARITY [23], GETIMAGETOIMAGESSIMI-144 LARITY [23], GETTEXTTOIMAGESSIMILARITY [23]. In-145 spired by prior works on multi-modal tool use [7, 8, 18, 146 22, 26], we include a few additional tools to help with 147 reasoning: QUERYLANGUAGEMODEL, QUERYKNOWL-148 EDGEBASE, CALCULATE, and SOLVEMATHEQUATION. 149 We also include TERMINATE as a tool for the model to out-150 put a final answer in the same action format. Our final ac-151 tion space consists of 15 tools, and their full implementation 152 details can be found in the Appendix. 153

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217



Figure 2. Data generation. We illustrate our model-based data generation (top) and programmatic generation (bottom) pipelines.

154 2.1. LATTE-trace generation

We generate synthetic LATTE-trace data with two automatic approaches: Model-based generation and Programmatic data generation.

Model-based generation. The model-based data generation pipeline consists of three steps (Figure 2 top):

1. GENERATE. First, we leverage images and QA ex-160 amples in existing visual instruction tuning datasets and 161 generate LATTE-traces to solve the questions with GPT-162 163 40 (2024-08-06). We include diverse questions on both single-image and multi-image examples from two large-164 165 scale instruction tuning datasets, Cauldron and Mantis-166 Instruct [11, 14]. We feed the images and questions to GPT-40 and prompt it to answer the questions by following a 167 LATTE-trace or just CoT when it is not necessary (e.g., the 168 question is straightforward) or helpful (e.g., the question re-169 quires domain-specific knowledge out of the scope of avail-170 able tools) to call specialized vision tools (Figure 2). 171

VERIFY. Second, we verify GPT-4o's generated answers against the ground-truth. We force GPT-4o to always
end with TERMINATE(answer) and compare its prediction
to the ground-truth. If the final answer following a reasoning trace is correct, we move this LATTE-trace to the next
stage. Otherwise, we convert this example into the direct
answer (Direct) format with the ground-truth (Figure 2).

3. PARSE. Finally, we check the JSON syntax of each step
of the LATTE-trace. Similar to the previous stage, we again
keep the LATTE-traces free of syntax errors and turn the
others into the Direct format with ground-truth answers.

Programmatic data generation. While model-based data
generation distills reasoning from proprietary models, we
are curious if reasoning with vision specialists can be
learned in another manner without reliance on proprietary
models. To study this perspective, we implement a pro-

traces (Figure 2 bottom). This pipeline involves two steps: 189 1. ANNOTATE. First, we gather existing dense annotations 190 of images. We adopt Visual Genome (VG) as it contains 191 rich human annotations of objects, attributes, and relation-192 ships of the images. In addition, we obtain depth maps of 193 the VG images with Depth-Anything-v2 [29]. 194 2. **GENERATE**. Next, we programmatically generate both 195 the QA pairs and the corresponding LATTE-traces with 196 manually written templates and the dense annotations of the 197 images. We reuse the pipeline from [31, 32] for generating 198 diverse QA pairs that cover various vision capabilities such 199 as counting and spatial understanding. To generate LATTE-200 traces, we define templates for thoughts, actions, and obser-201 vations across all steps. See Appendix for more details. 202

grammatic data generation engine for synthesizing LATTE-

3. Experiments

We perform extensive experiments with small multi-modal models and 9 data recipes on 6 benchmarks.

Models. We adopt models that support multi-image inputs as our data includes reasoning traces with multiple images. For most of our experiments, we use Mantis-8B-SigLIP-LLaMA-3 as the base model. We additionally experiment with Mantis-8B-CLIP-LLaMA-3, and LLaVA-OneVision-7B (Qwen2-7B and SigLIP) in our ablations.

Baselines. We compare our model to two types of baselines: (1) vanilla instruction-tuning (Vanilla IT) – instruction-tuning data with only direct answers – and (2) distillation methods that train small models by distilling both perception and reasoning from larger models, including VPD [9], VisCoT [25], and LLaVa-CoT[28].

Evaluation setup.We select 6 multi-modal benchmarks218covering both perception and reasoning.The perception-219focused benchmarks include RealWorldQA, CV-Bench and220

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

Method	Perception				Perception	Overall			
	BLINK	CV-Bench	RealWorldQA	Avg	MathVista	MMStar	MMVet	Avg	Avg
Vanilla IT	44.1	<u>49.2</u>	41.4	44.9	31.0	39.7	27.8	32.8	38.9
VPD	41.6	48.8	44.8	<u>45.1</u> (+0.2)	33.0	41.1	32.8	35.7 (+2.8)	<u>40.4</u> (+1.5)
LLaVa-CoT	42.2	40.4	38.0	40.2 (-4.7)	<u>36.7</u>	44.6	40.2	<u>40.5</u> (+7.7)	<u>40.4</u> (+1.5)
LATTE	46.4	54.0	42.0	47.5 (+2.6)	36.9	44.2	47.9	43.0 (+10.2)	45.2 (+6.4)

Table 1. LATTE vs. Baselines on Perception and Reasoning Benchmarks. Our method LATTE brings substantial gains over the vanilla instruction-tuned (Vanilla IT) baseline on both perception and perception + reasoning benchmarks.

Table 2. LATTE vs. Vanilla IT with Different Models. We learn that LATTE leads to performance gains over Vanilla IT regardless of the base models. The gains are 2-5% on average across all 6 benchmarks and up to 16% on MMVet.

Language / Vision	Starting	Method	Perception				Perception + Reasoning				Overall
88.	checkpoint		CV-Bench	BLINK	RealWorldQA	Avg	MathVista	MMStar	MMVet	Avg	Avg
LLaMA3-8B / CLIP	Mantis	Vanilla IT LATTE	52.6 56.9	45.8 49.6	52.3 51.1	50.2 52.6	33.1 36.6	36.7 40.8	28.9 45.2	32.9 40.8	41.6 46.7 (+5.1)
LLaMA3-8B / SigLIP	Pretrained	Vanilla IT LATTE	52.3 57.2	43.7 47.8	51.8 53.7	49.3 52.9	31.1 34.9	40.5 44.6	33.0 45.2	34.9 41.6	42.1 47.2 (+5.2)
	Mantis Instruct-tuned	Vanilla IT LATTE	50.6 51.7	46.7 47.6	54.8 56.5	50.7 51.9	36.2 36.3	40.7 42.5	29.7 45.7	35.5 41.5	43.1 46.7 (+3.6)
Qwen2-7B / SigLIP	LLaVa-OV Stage 1.5	Vanilla IT LATTE	56.8 60.2	50.3 49.9	57.8 58.8	55.0 56.3	42.4 41.9	50.1 51.0	39.3 50.9	43.9 48.0	49.5 52.1 (+2.7)

221 222

223

BLINK [5, 15, 24, 27], and the perception + reasoning ones are MathVista, MMStar, and MMVet [3, 21, 30]. Additional details can be found in the Appendix.

3.1. Main results



Figure 3. Performance of LATTE vs. Baselines across Training Data Scales. We find that our method leads to consistent gains across varying training data sizes – 98K, 200K and 293K.

225 Our method leads to substantial gains compared to 226 vanilla instruction-tuning on both perception and reasoning benchmarks, whereas other distillation baselines 227 result in smaller gains or even degradation on some per-228 ception tasks. We find that learning to reason with vision 229 specialists enables our model to achieve consistent gains 230 231 on perception-focused VQA benchmarks as well as benchmarks that require both perception and reasoning, with av-232 233 erage gains of 2.6% and 10.2% respectively (Table 1). By contrast, both distillation baselines VPD and LLaVa-CoT 234 235 bring much smaller gains, with an average of 1.5% across 236 all benchmarks, compared to ours (6.4%). Further, we observe that the same trend holds as we scale the training 237 data size from 98K to 200K and 293K, where our method 238 consistently brings larger gains on both perception and per-239 ception + reasoning benchmarks (Figure 3). Interestingly, 240 241 LLaVa-CoT even hurts the model's performance on perception benchmarks, even though it increases the performance on the perception + reasoning benchmarks (Table 1). This result suggests that GPT4-0 might still be inferior to vision specialists on some perception tasks, as LLaVa-CoT distills purely from GPT4-0.

Our method beats the vanilla instruction-tuning baseline on average across all benchmarks regardless of the base model and checkpoint, with significant gains of 10-16% on MMVet. We fine-tune 3 different multi-modal models with all 293K LATTE-traces starting from different checkpoints. We observe that our method leads to consistent gains of 2-5% in the model's average accuracy across 6 benchmarks compared to the baselines instruction-tuned with the same examples in the Direct format (Table 2). We note that our method results in staggering gains of 10-16% on MMVet, which covers a wide range of perceptual and reasoning capabilities. Moreover, we find that our data results in larger gains on earlier pretrained checkpoints than on later-stage instruction-tuned checkpoints, likely due to the relatively small size of our data compared to Mantis' and LLaVa-OV's instruction-tuning data (1.2M and 4.5M) and some overlap in the images and questions [11, 16].

4. Conclusion

We propose to learn multi-modal language models to reason with vision specialists instead of becoming both vision specialists and reasoning experts.

Limitations and Future Work. First, our method requires268customized implementations of the specialized vision tools.269Second, reasoning with the vision specialists also requires270additional compute at inference time. Future work can optimize and enhance the implementations of vision specialists.272

293

295

298

303

304

305

306

307

308

309

310

311

312

313

314

315

316

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

386

References 273

- 274 [1] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret 275 Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 276 Vga: Visual question answering. In Proceedings of the IEEE 277 international conference on computer vision, pages 2425-278 2433, 2015. 1
- 279 [2] Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, 280 Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, 281 Qiushi Du, Zhe Fu, et al. Deepseek llm: Scaling open-282 source language models with longtermism. arXiv preprint 283 arXiv:2401.02954, 2024. 1
- 284 [3] Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang 285 Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, 286 Dahua Lin, et al. Are we on the right way for evaluating large 287 vision-language models? arXiv preprint arXiv:2403.20330, 288 2024. 4
- 289 [4] Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tri-290 pathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, et al. Molmo 292 and pixmo: Open weights and open data for state-of-the-art multimodal models. arXiv preprint arXiv:2409.17146, 2024. 294
- [5] Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, 296 Xudong Lin, Dan Roth, Noah A Smith, Wei-Chiu Ma, and 297 Ranjay Krishna. Blink: Multimodal large language models can see but not perceive. arXiv preprint arXiv:2404.12390, 299 2024.4
- 300 [6] Melvyn A Goodale and A David Milner. Separate visual 301 pathways for perception and action. Trends in neurosciences, 302 15(1):20-25, 1992. 1
 - [7] Tanmay Gupta and Aniruddha Kembhavi. Visual programming: Compositional visual reasoning without training. ArXiv, abs/2211.11559, 2022. 2
 - [8] Yushi Hu, Weijia Shi, Xingyu Fu, Dan Roth, Mari Ostendorf, Luke Zettlemoyer, Noah A Smith, and Ranjay Krishna. Visual sketchpad: Sketching as a visual chain of thought for multimodal language models, 2024. 1, 2
 - [9] Yushi Hu, Otilia Stretcu, Chun-Ta Lu, Krishnamurthy Viswanathan, Kenji Hata, Enming Luo, Ranjay Krishna, and Ariel Fuxman. Visual program distillation: Distilling tools and programmatic reasoning into vision-language models. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 9590-9601, 2024. 1,
- 317 [10] JadedAI. Easyocr, 2025. 2
- 318 [11] Dongfu Jiang, Xuan He, Huaye Zeng, Con Wei, Max Ku, 319 Qian Liu, and Wenhu Chen. Mantis: Interleaved multi-image 320 instruction tuning. arXiv preprint arXiv:2405.01483, 2024. 321 3,4
- 322 [12] Georg B Keller, Tobias Bonhoeffer, and Mark Hübener. Sen-323 sorimotor mismatch signals in primary visual cortex of the 324 behaving mouse. Neuron, 74(5):809-815, 2012. 1
- 325 [13] Nathan Lambert, Jacob Morrison, Valentina Pyatkin, 326 Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester 327 James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. 328 T'' ulu 3: Pushing frontiers in open language model post-329 training. arXiv preprint arXiv:2411.15124, 2024. 1

- [14] Hugo Laurencon, Léo Tronchon, Matthieu Cord, and Victor 330 Sanh. What matters when building vision-language models?, 331 2024. 3 332
- [15] Bohao Li, Yuying Ge, Yixiao Ge, Guangzhi Wang, Rui Wang, Ruimao Zhang, and Ying Shan. Seed-bench-2: Benchmarking multimodal large language models. arXiv preprint arXiv:2311.17092, 2023. 4
- [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. arXiv preprint arXiv:2408.03326, 2024. 4
- [17] Jian Li, Yabiao Wang, Changan Wang, Ying Tai, Jianjun Qian, Jian Yang, Chengjie Wang, Jilin Li, and Feiyue Huang. Dsfd: Dual shot face detector. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019. 2
- [18] Shilong Liu, Hao Cheng, Haotian Liu, Hao Zhang, Feng Li, Tianhe Ren, Xueyan Zou, Jianwei Yang, Hang Su, Jun Zhu, Lei Zhang, Jianfeng Gao, and Chunyuan Li. Llava-plus: Learning to use tools for creating multimodal agents, 2023. 2
- [19] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. arXiv preprint arXiv:2303.05499, 2023. 1, 2
- [20] Jiasen Lu, Christopher Clark, Sangho Lee, Zichen Zhang, Savya Khosla, Ryan Marten, Derek Hoiem, and Aniruddha Kembhavi. Unified-io 2: Scaling autoregressive multimodal models with vision language audio and action. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 26439-26455, 2024. 1
- [21] 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 International Conference on Learning Representations (ICLR), 2024. 4
- [22] Zixian Ma, Weikai Huang, Jieyu Zhang, Tanmay Gupta, and Ranjay Krishna. m&m's: A benchmark to evaluate tool-use for multi-step multi-modal tasks. EECV 2024, 2024. 1, 2
- [23] Alec Radford, Jong Wook Kim, Chris Hallacy, A. Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In ICML, 2021. 2
- [24] Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. A-okvqa: A benchmark for visual question answering using world knowledge, 2022. 4
- [25] Hao Shao, Shengju Qian, Han Xiao, Guanglu Song, Zhuofan Zong, Letian Wang, Yu Liu, and Hongsheng Li. Visual cot: Unleashing chain-of-thought reasoning in multi-modal language models, 2024. 1, 3
- [26] Dídac Surís, Sachit Menon, and Carl Vondrick. Vipergpt: Visual inference via python execution for reasoning. arXiv preprint arXiv:2303.08128, 2023. 2

- [27] Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann
 LeCun, and Saining Xie. Eyes wide shut? exploring the
 visual shortcomings of multimodal llms, 2024. 4
- [28] Guowei Xu, Peng Jin, Hao Li, Yibing Song, Lichao Sun, and
 Li Yuan. Llava-cot: Let vision language models reason stepby-step, 2025. 1, 3
- [29] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2. arXiv:2406.09414, 2024. 1, 2, 3
- [30] Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang.
 Mm-vet: Evaluating large multimodal models for integrated capabilities. In *International conference on machine learning.* PMLR, 2024. 4
- 401 [31] Jieyu Zhang, Weikai Huang, Zixian Ma, Oscar Michel, Dong
 402 He, Tanmay Gupta, Wei-Chiu Ma, Ali Farhadi, Aniruddha
 403 Kembhavi, and Ranjay Krishna. Task me anything. *arXiv*404 *preprint arXiv:2406.11775*, 2024. 3
- [32] Jieyu Zhang, Le Xue, Linxin Song, Jun Wang, Weikai Huang, Manli Shu, An Yan, Zixian Ma, Juan Carlos Niebles, Silvio Savarese, Caiming Xiong, Zeyuan Chen, Ranjay Krishna, and Ran Xu. Provision: Programmatically scaling vision-centric instruction data for multimodal language models. *Preprint*, 2024. 3
- [33] Youcai Zhang, Xinyu Huang, Jinyu Ma, Zhaoyang Li,
 Zhaochuan Luo, Yanchun Xie, Yuzhuo Qin, Tong Luo,
 Yaqian Li, Shilong Liu, et al. Recognize anything: A strong
 image tagging model. *arXiv preprint arXiv:2306.03514*,
 2023. 2