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

#### **026 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 [24]. These tasks demand both perception 030 and 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 [51].

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 [64, 77] or042specialized vision models [25]. 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 [4, 31], 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 [79] for depth 051 estimation and GroundingDINO [44] for object recogni-052 tion-while unified models still lag behind [46]. Simi-053 larly, the human brain dedicates distinct regions to cat-054 egorical recognition (ventral stream) and spatial reason-055 ing (dorsal stream)[19], with the reasoning and language-056 processing frontal and temporal lobes occupying a smaller 057 volume[29]. By contrast, vision-language models remain 058 heavily skewed toward language, treating visual encoders 059 as an afterthought [15]. 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<br/>more). Finally, programmatically-generated LATTE-trace119can hurt model performance as a result of the worse reason-<br/>ing quality, suggesting that again that high-quality reason-<br/>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 [27], GETOBJECTS [88], LOCALIZEOBJECTS [44], ES-141 TIMATEOBJECTDEPTH [79], ESTIMATEREGIONDEPTH 142 [79], DETECTFACES [36], CROP, ZOOMIN, GETIMAGE-143 TOTEXTSSIMILARITY [61], GETIMAGETOIMAGESSIMI-144 LARITY [61], GETTEXTTOIMAGESSIMILARITY [61]. In-145 spired by prior works on multi-modal tool use [21, 24, 43, 146 51, 70], 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

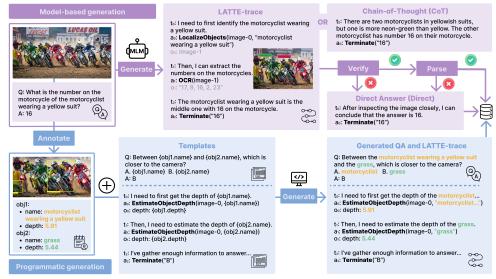


Figure 2. Data generation. We illustrate our model-based data generation (top) and programmatic generation (bottom) pipelines.

#### **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 [28, 32]. 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 [79]. 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 [84, 86] 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 [25], VisCoT [64], and LLaVa-CoT[77].

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	Percepti	on		Perception	Perception + Reasoning				
1.100100	BLINK	CV-Bench	RealWorldQA	Avg	MathVista	MMStar	MMVet	Avg	Avg
Vanilla IT	44.1	49.2	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	<u>44.2</u>	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	
	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 <b>60.2</b>	<b>50.3</b> 49.9	57.8 <b>58.8</b>	55.0 <b>56.3</b>	<b>42.4</b> 41.9	50.1 <b>51.0</b>	39.3 <b>50.9</b>	43.9 <b>48.0</b>	49.5 <b>52.1</b> (+2.7)

221 222

223

BLINK [17, 33, 63, 73], and the perception + reasoning ones are MathVista, MMStar, and MMVet [8, 48, 83]. Additional details can be found in the Appendix.

#### **3.1.** Main results

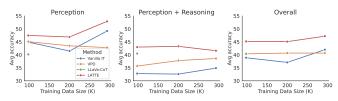


Figure 3. Performance of LATTE vs. Baselines across Training Data Scales. We find that our method leads to consistent gains across varying training data sizes – 98K, 200K and 293K.

225 Our method leads to substantial gains compared to 226 vanilla instruction-tuning on both perception and reasoning benchmarks, whereas other distillation baselines 227 228 result in smaller gains or even degradation on some perception 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-o might still be inferior to vision specialists on some perception tasks, as LLaVa-CoT distills purely from GPT4-o.

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 [28, 35].

#### 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

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

309

332

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] Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf 280 Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, 281 Samir Gadre, Shiori Sagawa, Jenia Jitsev, Simon Kornblith, Pang Wei Koh, Gabriel Ilharco, Mitchell Wortsman, and 282 283 Ludwig Schmidt. Openflamingo: An open-source frame-284 work for training large autoregressive vision-language mod-285 els, 2023. 5
- 286 [3] Jing Bi, Nguyen Manh Nguyen, Ali Vosoughi, and Chenliang Xu. Misar: A multimodal instructional system with 288 augmented reality, 2023. 5
  - Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, [4] Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qiushi Du, Zhe Fu, et al. Deepseek llm: Scaling opensource language models with longtermism. arXiv preprint arXiv:2401.02954, 2024. 1
  - [5] Baian Chen, Chang Shu, Ehsan Shareghi, Nigel Collier, Karthik Narasimhan, and Shunyu Yao. Fireact: Toward language agent fine-tuning, 2023. 4
  - [6] Guo Chen, Yin-Dong Zheng, Jiahao Wang, Jilan Xu, Yifei Huang, Junting Pan, Yi Wang, Yali Wang, Yu Qiao, Tong Lu, and Limin Wang. Videollm: Modeling video sequence with large language models, 2023. 5
  - [7] Houlun Chen, Xin Wang, Hong Chen, Zihan Song, Jia Jia, and Wenwu Zhu. Grounding-prompter: Prompting llm with multimodal information for temporal sentence grounding in long videos, 2023. 5
- 305 [8] Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang 306 Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, 307 Dahua Lin, et al. Are we on the right way for evaluating large 308 vision-language models? arXiv preprint arXiv:2403.20330, 2024. 4
- 310 [9] Xi Chen, Xiao Wang, Lucas Beyer, Alexander Kolesnikov, 311 Jialin Wu, Paul Voigtlaender, Basil Mustafa, Sebastian 312 Goodman, Ibrahim Alabdulmohsin, Piotr Padlewski, Daniel 313 Salz, Xi Xiong, Daniel Vlasic, Filip Pavetic, Keran Rong, 314 Tianli Yu, Daniel Kevsers, Xiaohua Zhai, and Radu Sori-315 cut. Pali-3 vision language models: Smaller, faster, stronger, 316 2023. 5
- 317 [10] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, 318 Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, 319 Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng 320 Dai. Internvl: Scaling up vision foundation models and 321 aligning for generic visual-linguistic tasks. arXiv preprint 322 arXiv:2312.14238, 2023. 5
- 323 [11] Zhipeng Chen, Kun Zhou, Beichen Zhang, Zheng Gong, 324 Wayne Xin Zhao, and Ji-Rong Wen. Chatcot: Tool-325 augmented chain-of-thought reasoning on chat-based large 326 language models. arXiv preprint arXiv:2305.14323, 2023. 4
- 327 [12] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhang-328 wei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng 329 Luo, Zheng Ma, et al. How far are we to gpt-4v? closing

the gap to commercial multimodal models with open-source 330 suites. arXiv preprint arXiv:2404.16821, 2024. 5 331

- [13] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. arXiv preprint arXiv:2110.14168, 2021. 4
- [14] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung, and Steven Hoi. Instructblip: Towards generalpurpose vision-language models with instruction tuning. Advances in Neural Information Processing Systems, 36, 2024. 5
- [15] Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, et al. Molmo and pixmo: Open weights and open data for state-of-the-art multimodal models. arXiv preprint arXiv:2409.17146, 2024. 1
- [16] Haodong Duan, Junming Yang, Yuxuan Qiao, Xinyu Fang, Lin Chen, Yuan Liu, Xiaoyi Dong, Yuhang Zang, Pan Zhang, Jiaqi Wang, Dahua Lin, and Kai Chen. Vlmevalkit: An opensource toolkit for evaluating large multi-modality models, 2024. 1, 6
- [17] Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A Smith, Wei-Chiu Ma, and Ranjay Krishna. Blink: Multimodal large language models can see but not perceive. arXiv preprint arXiv:2404.12390, 2024. 4
- [18] Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. Pal: Program-aided language models. In International Conference on Machine Learning, pages 10764-10799. PMLR, 2023. 4
- [19] Melvyn A Goodale and A David Milner. Separate visual pathways for perception and action. Trends in neurosciences, 15(1):20-25, 1992. 1
- [20] Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujiu Yang, Nan Duan, and Weizhu Chen. Critic: Large language models can self-correct with tool-interactive critiquing. arXiv preprint arXiv:2305.11738, 2023. 4
- [21] Tanmay Gupta and Aniruddha Kembhavi. Visual programming: Compositional visual reasoning without training. ArXiv, abs/2211.11559, 2022. 2, 4
- [22] Joy He-Yueya, Gabriel Poesia, Rose E Wang, and Noah D Goodman. Solving math word problems by combining language models with symbolic solvers. arXiv preprint arXiv:2304.09102, 2023. 4
- [23] 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, 2023. 5
- [24] 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, 5

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

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

- [25] 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,
  3
- [26] Bin Huang, Xin Wang, Hong Chen, Zihan Song, and Wenwu
  Zhu. Vtimellm: Empower llm to grasp video moments. *arXiv preprint arXiv:2311.18445*, 2(3):9, 2023. 5
- **397** [27] JadedAI. Easyocr, 2025. 2
- [28] Dongfu Jiang, Xuan He, Huaye Zeng, Con Wei, Max Ku,
  Qian Liu, and Wenhu Chen. Mantis: Interleaved multi-image
  instruction tuning. *arXiv preprint arXiv:2405.01483*, 2024.
  3, 4, 1, 5
- 402 [29] Georg B Keller, Tobias Bonhoeffer, and Mark Hübener. Sensorimotor mismatch signals in primary visual cortex of the behaving mouse. *Neuron*, 74(5):809–815, 2012. 1
- [30] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson,
  Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and
  Fei-Fei Li. Visual genome: Connecting language and vision
  using crowdsourced dense image annotations, 2016. 1
- [31] Nathan Lambert, Jacob Morrison, Valentina Pyatkin,
  Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester
  James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al.
  T\" ulu 3: Pushing frontiers in open language model posttraining. *arXiv preprint arXiv:2411.15124*, 2024. 1
- [32] Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor
  Sanh. What matters when building vision-language models?,
  2024. 3, 1
- [33] 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
- [34] Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Fanyi
  Pu, Jingkang Yang, Chunyuan Li, and Ziwei Liu. Mimic-it:
  Multi-modal in-context instruction tuning, 2023. 5
- [35] 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, 5
- [36] 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
- [37] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi.
  Blip: Bootstrapping language-image pre-training for unified
  vision-language understanding and generation. In *ICML*,
  2022. 5
- [38] 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 *ICML*,
  2023.
- [39] Kevin Lin, Faisal Ahmed, Linjie Li, Chung-Ching Lin,Ehsan Azarnasab, Zhengyuan Yang, Jianfeng Wang, Lin

Liang, Zicheng Liu, Yumao Lu, Ce Liu, and Lijuan Wang. Mm-vid: Advancing video understanding with gpt-4v(ision), 2023. 5

- [40] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2023.
- [41] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023.
- [42] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, 2024. 1, 5
- [43] 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, 5
- [44] 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
- [45] Shikun Liu, Linxi Fan, Edward Johns, Zhiding Yu, Chaowei Xiao, and Anima Anandkumar. Prismer: A vision-language model with multi-task experts, 2024. 5
- [46] 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
- [47] Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. Chameleon: Plug-and-play compositional reasoning with large language models, 2023. 5
- [48] 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
- [49] Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. Chameleon: Plug-and-play compositional reasoning with large language models. Advances in Neural Information Processing Systems, 36, 2024. 4
- [50] Chenyang Lyu, Minghao Wu, Longyue Wang, Xinting Huang, Bingshuai Liu, Zefeng Du, Shuming Shi, and Zhaopeng Tu. Macaw-llm: Multi-modal language modeling with image, audio, video, and text integration, 2023. 5
- [51] 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, 4
- [52] Michael Moor, Qian Huang, Shirley Wu, Michihiro Yasunaga, Cyril Zakka, Yash Dalmia, Eduardo Pontes Reis, Pranav Rajpurkar, and Jure Leskovec. Med-flamingo: a multimodal medical few-shot learner, 2023. 5
  498
- [53] Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu
   500

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

Jain, Vineet Kosaraju, William Saunders, et al. Webgpt:
Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021. 4

- 504 [54] Artemis Panagopoulou, Le Xue, Ning Yu, Junnan Li,
  505 Dongxu Li, Shafiq R. Joty, Ran Xu, Silvio Savarese, Caim506 ing Xiong, and Juan Carlos Niebles. X-instructblip: A
  507 framework for aligning x-modal instruction-aware represen508 tations to llms and emergent cross-modal reasoning. *ArXiv*,
  509 abs/2311.18799, 2023. 5
- [55] Bhargavi Paranjape, Scott Lundberg, Sameer Singh, Hannaneh Hajishirzi, Luke Zettlemoyer, and Marco Tulio
  Ribeiro. Art: Automatic multi-step reasoning and tool-use
  for large language models. *arXiv preprint arXiv:2303.09014*, 2023. 4
- 515 [56] Shishir G. Patil, Tianjun Zhang, Xin Wang, and Joseph E.
  516 Gonzalez. Gorilla: Large language model connected with
  517 massive apis. *arXiv preprint arXiv:2305.15334*, 2023. 4
- [57] Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan
  Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding
  multimodal large language models to the world, 2023. 5
- [58] Ji Qi, Ming Ding, Weihan Wang, Yushi Bai, Qingsong Lv,
  Wenyi Hong, Bin Xu, Lei Hou, Juanzi Li, Yuxiao Dong, and
  Jie Tang. Cogcom: Train large vision-language models diving into details through chain of manipulations, 2024. 5
- [59] Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan,
  Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian,
  Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie
  Zhou, Mark Gerstein, Dahai Li, Zhiyuan Liu, and Maosong
  Sun. Toolllm: Facilitating large language models to master
  16000+ real-world apis, 2023. 4
- [60] Changle Qu, Sunhao Dai, Xiaochi Wei, Hengyi Cai,
  Shuaiqiang Wang, Dawei Yin, Jun Xu, and Ji-Rong Wen.
  Tool learning with large language models: A survey. *arXiv preprint arXiv:2405.17935*, 2024. 4
- [61] 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
- 540 [62] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta
  541 Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Can542 cedda, and Thomas Scialom. Toolformer: Language models
  543 can teach themselves to use tools, 2023. 4
- 544 [63] Dustin Schwenk, Apoorv Khandelwal, Christopher Clark,
  545 Kenneth Marino, and Roozbeh Mottaghi. A-okvqa: A
  546 benchmark for visual question answering using world knowl547 edge, 2022. 4
- 548 [64] Hao Shao, Shengju Qian, Han Xiao, Guanglu Song, Zhuo549 fan Zong, Letian Wang, Yu Liu, and Hongsheng Li. Visual
  550 cot: Unleashing chain-of-thought reasoning in multi-modal
  551 language models, 2024. 1, 3, 5
- [65] 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. 4
- [66] Mustafa Shukor, Corentin Dancette, Alexandre Rame, and
  Matthieu Cord. Unival: Unified model for image, video, audio and language tasks, 2023. 5

- [67] Guangzhi Sun, Wenyi Yu, Changli Tang, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun Ma, and Chao Zhang. Finegrained audio-visual joint representations for multimodal large language models, 2023.
   562
- [68] Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang, Qiying Yu, Zhengxiong Luo, Yueze Wang, Yongming Rao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. Generative multimodal models are in-context learners, 2024.
- [69] Quan Sun, Qiying Yu, Yufeng Cui, Fan Zhang, Xiaosong Zhang, Yueze Wang, Hongcheng Gao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. Emu: Generative pretraining in multimodality, 2024. 5
- [70] Dídac Surís, Sachit Menon, and Carl Vondrick. Vipergpt: Visual inference via python execution for reasoning. *arXiv* preprint arXiv:2303.08128, 2023. 2, 4
- [71] Qiaoyu Tang, Ziliang Deng, Hongyu Lin, Xianpei Han, Qiao Liang, Boxi Cao, and Le Sun. Toolalpaca: Generalized tool learning for language models with 3000 simulated cases. arXiv preprint arXiv:2306.05301, 2023. 4
- [72] Yunlong Tang, Jinrui Zhang, Xiangchen Wang, Teng Wang, and Feng Zheng. Llmva-gebc: Large language model with video adapter for generic event boundary captioning, 2023.
- [73] 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
- [74] Junke Wang, Dongdong Chen, Chong Luo, Xiyang Dai, Lu Yuan, Zuxuan Wu, and Yu-Gang Jiang. Chatvideo: A tracklet-centric multimodal and versatile video understanding system, 2023. 5
- [75] Yi Wang, Kunchang Li, Xinhao Li, Jiashuo Yu, Yinan He, Guo Chen, Baoqi Pei, Rongkun Zheng, Jilan Xu, Zun Wang, Yansong Shi, Tianxiang Jiang, Songze Li, Hongjie Zhang, Yifei Huang, Yu Qiao, Yali Wang, and Limin Wang. Internvideo2: Scaling video foundation models for multimodal video understanding. *ArXiv*, abs/2403.15377, 2024.
- [76] Zhanyu Wang, Longyue Wang, Zhen Zhao, Minghao Wu, Chenyang Lyu, Huayang Li, Deng Cai, Luping Zhou, Shuming Shi, and Zhaopeng Tu. Gpt4video: A unified multimodal large language model for Instruction-followed understanding and safety-aware generation, 2023. 5
- [77] 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
- [78] Le Xue, Manli Shu, Anas Awadalla, Jun Wang, An Yan, Senthil Purushwalkam, Honglu Zhou, Viraj Prabhu, Yutong Dai, Michael S Ryoo, Shrikant B. Kendre, Jieyu Zhang, Can Qin, Shu Zhen Zhang, Chia-Chih Chen, Ning Yu, Juntao Tan, Tulika Awalgaonkar, Shelby Heinecke, Huan Wang, Yejin Choi, Ludwig Schmidt, Zeyuan Chen, Silvio Savarese, Juan Carlos Niebles, Caiming Xiong, and Ran Xu. xgen-mm (blip-3): A family of open large multimodal models. *ArXiv*, abs/2408.08872, 2024. 5
- [79] 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
- [80] Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin,
   Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu, Ce Liu,
   616

- 617Michael Zeng, and Lijuan Wang. Mm-react: Prompting chat-618gpt for multimodal reasoning and action, 2023. 4
- [81] Shunyu Yao, Howard Chen, John Yang, and Karthik
  Narasimhan. Webshop: Towards scalable real-world web interaction with grounded language agents. Advances in Neu-*ral Information Processing Systems*, 35:20744–20757, 2022.
  4
- [82] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran,
  Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models, 2023. 4
- [83] 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 learn- ing.* PMLR, 2024. 4
- [84] Jieyu Zhang, Weikai Huang, Zixian Ma, Oscar Michel, Dong
  He, Tanmay Gupta, Wei-Chiu Ma, Ali Farhadi, Aniruddha
  Kembhavi, and Ranjay Krishna. Task me anything. *arXiv preprint arXiv:2406.11775*, 2024. 3
- [85] Jianguo Zhang, Tian Lan, Ming Zhu, Zuxin Liu, Thai Hoang,
  Shirley Kokane, Weiran Yao, Juntao Tan, Akshara Prabhakar, Haolin Chen, Zhiwei Liu, Yihao Feng, Tulika Awalgaonkar, Rithesh Murthy, Eric Hu, Zeyuan Chen, Ran Xu,
  Juan Carlos Niebles, Shelby Heinecke, Huan Wang, Silvio
  Savarese, and Caiming Xiong. xlam: A family of large action models to empower ai agent systems, 2024. 4
- [86] 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
- [87] Kechi Zhang, Jia Li, Ge Li, Xianjie Shi, and Zhi Jin.
  Codeagent: Enhancing code generation with tool-integrated
  agent systems for real-world repo-level coding challenges. *arXiv preprint arXiv:2401.07339*, 2024. 4
- [88] 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

## Section 2018 LATTE: Learning to Reason with Vision Specialists

Supplementary Material

#### **658 A. Additional details on methods**

#### 659 A.1. Programmatic generation

We manually design five thought templates for each action 660 and randomly sample one during generation. As for the ac-661 tions, we manually select the specialized vision tools for 662 663 each type of questions (e.g., ESTIMATEOBJECTDEPTH for questions on objects' relative depths, and LOCALIZE for ob-664 ject counting questions, etc.) and compose templates with 665 them. We fill in the actions' inputs in the templates with 666 annotations obtained from the previous step. Similarly, we 667 define observation templates based on the action outputs' 668 formats, and fill them in with dense annotations. 669

#### 670 A.2. Data filtering and mixing

In total, we generate 815K data with GPT-40 for both 671 single-image and multi-image questions across 31 data 672 sources from Cauldron and Mantis-Instruct [28, 32]. We 673 also programmatically generate 1M+ QA pairs and corre-674 675 sponding LATTE-traces with VG images and annotations [30], which we randomly sample from to augment model-676 generated data. We then develop different data recipes with 677 678 3 filtering/mixing techniques, where we vary the distribution of: (1) data formats; (2) data sources; and (3) model-679 generated vs. program-generated reasoning traces. 680

681 Data format. As mentioned in the data generation pipeline, model-generated data can be categorized into two formats: 682 683 LATTE-trace or CoT examples (Figure 2). Additionally, they can be further grouped into LATTE-trace/CoT-pos and 684 LATTE-trace/CoT-neg examples where the final answers 685 686 are correct and wrong respectively (Figure 4). Note that 687 we convert both LATTE-trace-neg and CoT-neg examples 688 into the direct answer (Direct) format with ground-truth an-689 swers (Figure 2) so the final data format is one of LATTEtrace, CoT, and Direct. We experiment with data consisting 690 of only LATTE-trace, LATTE-trace + CoT, LATTE-trace + 691 Direct, and all three formats (Figure 4). 692

Data source. We also perform filtering based on data 693 694 sources as Cauldron and Mantis-Instruct cover a wide range of tasks, some of which benefit more from vision special-695 ists than others. To this end, we define LATTE-useless 696 datasets as the ones where GPT-40 either decides to out-697 698 put CoT much more often than LATTE-trace (i.e. % of 699 CoT-pos – LATTE-trace-pos > 10), or reaches wrong an-700 swers much more frequently than correct ones when using LATTE-trace (i.e. % LATTE-trace-neg - LATTE-trace-pos 701 702 > 10) (Figure 4). The remaining datasets are considered 703 LATTE-useful datasets, and we experiment with including 704 all datasets vs. only the LATTE-useful datasets (Figure 4).

Model-generated vs. program-generated data. Due to 705 the lack of controllability, the distribution of actions in 706 model-generated data is highly imbalanced, with a couple 707 of actions such as GETOBJECTS and OCR dominating the 708 dataset. Therefore, we also try increasing action diversity 709 by adding programmatic data involving underrepresented 710 actions such as LOCALIZEOBJECTS, ESTIMATEOBJECT-711 DEPTH and ESTIMATEREGIONDEPTH. We experiment 712 with data mixtures with different ratios of model-generated 713 vs. program-generated data, ranging from 1:1 to 1:0.1. 714

While the best data recipe varies with the benchmark,<br/>one recipe stands out, resulting in consistent gains in the<br/>average performance across 6 benchmarks for all 3 models.715This dataset consists of 293K LATTE-traces generated by<br/>GPT-40 across all 31 data sources.718

#### **B.** Additional Experimental Details

Baselines. For VPD, since their data and model are not re-721 leased, we reproduce their data by converting our LATTE-722 traces into CoTs in VPD's format, where we remove the ac-723 tion calls and add the observations to the reasoning chains. 724 For LLaVa-CoT, we directly use their data to finetune the 725 same model. For VisCoT, since it only has reasoning steps 726 for one data source, training with its data leads to much 727 worse performance. We include its results in the Appendix. 728 For fair comparison, we train our models and baselines with 729 the same model backbone (e.g. Mantis-SigLIP), the same 730 hyperparameters, and the same number of examples. 731

Training details. We fine-tune models starting from check-732 points at different stages - pretrained and instruction tuned 733 for Mantis-8B-SigLIP-LLaMA-3, and stage 1.5 for LLaVA-734 OneVision-7B - to investigate where LATTE-traces bring 735 the largest gains. We adopt the hyperparameters from 736 [28, 42] and fine-tune both the language model and the pro-737 jector with learning rate = 1e - 5 for 1 epoch with either 738 NVIDIA A100s 40GB or H100s 80GB. We additionally 739 perform hyperparameter tuning with LLaVA-OneVision-7B 740 and vary the tunable components, the language models' 741 learning rate, and number of epochs. We include this re-742 sult in the Appendix. 743

Evaluation details. We adapt VLMEvalKit [16] for our744evaluation, where an LLM judge (i.e. GPT-4-turbo) is745used to score predictions between 0 and 1 compared to the746groundtruth short answers for open-ended questions. Addi-747tional details about experimental setup can be found in the748Appendix.749

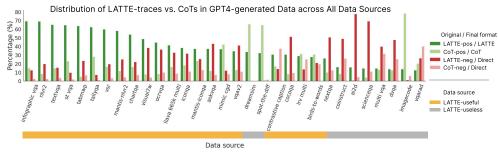


Figure 4. **Distribution of data formats and sources.** We visualize the frequency of data formats (i.e. LATTE-pos/neg, and CoT-pos/neg, pos = correct final answers, neg = incorrect) in the original GPT-4-generated data and in our training data (i.e. LATTE-trace, CoT, or Direct) across all data sources.

#### 750 C. Additional Results

#### 751 C.1. Main Results

Table 3. **Perception.** We highlight the large performance gains brought by our method on subsets of the perception benchmarks.

Method	BLINK				CV B	ench
method	Relative Depth*	Jigsaw	Multi-view Reasoning	Visual Similarity	2D	3D
Vanilla IT	51.6	51.3	51.1	65.2	38.9	59.4
VPD	46.0	44.0	55.6	59.3	40.4	57.3
LLaVa-CoT	53.2	52.0	54.1	47.4	45.1	35.7
LATTE	58.1	68.7	55.6	85.2	47.8	60.2

752 Learning to reason with vision specialists allows our 753 model to excel at vision-centric tasks that require finegrained perception. Taking a closer look at the perception-754 centric benchmarks, we note that our model does exception-755 ally well on the subsets that require fine-grained perception 756 757 such as depth estimation (Table 3). For example, our model 758 scores higher on the Relative Depth subset of BLINK and 3D subset of CV-Bench, most likely because it has learned 759 to utilize the ESTIMATEDEPTH tool. Similarly, learning to 760 use the localization specialists via LOCALIZEOBJECTS en-761 ables it to do better on the 2D tasks in CV-Bench by almost 762 763 9%; and the image-image and image-text similarity tools 764 boost its performance on Jigsaw, Multi-view reasoning, and Visual Similarity by significant margins of up to 20%. 765

#### 766 C.2. Ablations

We perform ablations with both model-generated and programmatic LATTE-traces. For model-generated data, we
explore two data filtering techniques to adjust the distribution of (1) data formats and (2) data sources (Figure 4).

771 Data quality matters more than quantity: our best data
772 recipe is the smallest and yet leads to better performance
773 than other recipes of larger scales. We find that 293K
774 LATTE-traces result in the biggest gain of 4.9% on aver775 age over the baseline across all benchmarks (Figure 5).
776 Adding CoT examples results in a smaller gain of 2.6%,
777 even though the training data size almost doubles (Fig-

Table 4. Ablation Results across Data Sources. We learn that including all data sources hurts model's perception and overall performance while including only LATTE-useful datasets helps.

Data source	Size	Method	Perception	Percept. + Reason.	Overall
All datasets	815K	Vanilla IT LATTE	50.7 47.7 (- <mark>3.0</mark> )	34.7 35.1 (+0.4)	42.7 41.4 (- <mark>1.4</mark> )
LATTE-useful datasets	566K	Vanilla IT LATTE	46.3 46.8 (+0.5)	33.3 35.6 (+2.3)	39.8 41.2 (+1.4)

Table 5. **Ablations on programmatic LATTE-traces.** We find that training with only or additional programmatic LATTE-traces hurts model performance.

M: P	Data format	Size	Perception	Percept. + Reason.	Overall
_	Direct		49.2	34.9	42.0
0:1	P-traces	293K	38.4	15.9	27.2
1:0	M-traces		52.9	<u>41.6</u>	47.2
1:0.1	+P-traces 29K	322K	47.8	40.1	44.0
1:0.25	+P-traces 73K	366K	50.5	42.1	46.3
1:0.5	+P-traces 147K	440K	<u>51.2</u>	39.7	45.5
1:1	+P-traces 293K	586K	50.3	36.2	43.2

ure 5). On the other hand, combining LATTE-trace and Di-778 rect examples hurts the model's performance compared to 779 LATTE-traces only, especially on the perception tasks (Fig-780 ure 5). This is likely because the mixture makes the model 781 more likely to rely on its own perceptual ability, which is 782 worse than the vision specialists'. We empirically observe 783 that models trained with a mix of LATTE-trace and Di-784 rect examples tend to adopt the Direct format more often 785 at inference time. Again, these results suggest that small 786 multi-modal language model's perceptual ability is weak, 787 and learning to reason with specialists helps the most with 788 its performance on perception tasks. 789

Data sources matter too: including all datasets hurts 790 performance whereas including only LATTE-useful 791 datasets brings performance gains. Similarly, we see that 792 including only the LATTE-useful datasets - where GPT-793 40 frequently chooses to reason with vision specialists and 794 reach correct final answers - improves the model's aver-795 age performance compared to the baseline, while including 796 all data sources does not (Table 4). Again, we see that a 797 smaller set of 566K LATTE-traces can lead to better per-798

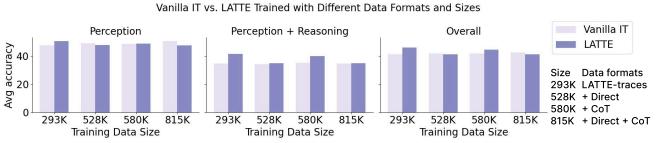


Figure 5. Ablation Results across Data Formats. The best data recipe is 293K LATTE-traces, which leads to the greatest gains over Vanilla IT and the highest overall performance. Adding either CoT or Direct doesn't bring additional gains despite the increased data size.

formance than a much larger dataset (815K), implying thatdata quality matters more than quantity.

Programmatically generated LATTE-traces lead to per-801 formance degradations, suggesting that we need both 802 high-quality perceptual information and reasoning to 803 achieve strong performance. We experiment with train-804 ing with only programmatic LATTE-traces and with a 805 mixture of model-generated and programmatic reasoning 806 traces, with different ratios ranging from 1:0.1 to 1:1. 807 Unfortunately, we find that training with only program-808 matic LATTE-traces results in large performance drops, 809 and adding programmatic LATTE-traces fails to bring ad-810 ditional gains despite the increase in overall data size (Ta-811 ble 5). This is likely due to the model's weaker reasoning 812 capability learned from templated reasoning traces, which 813 is further validated by qualitative examples (See Appendix). 814 These results suggest that reasoning over perceptual infor-815 816 mation is a challenging reasoning task that requires highquality reasoning data for the model to learn well, and that 817 model-generated reasoning traces are a lot more helpful for 818 this task than templated reasoning traces are. 819

820 Overall, our experiments suggest that the quality of per821 ceptual information and reasoning are both crucial to our
822 model's superior performance across diverse visual tasks.

#### 823 C.3. Additional qualitative examples

We present additional successful outputs of LATTE across
both single-image and multi-image examples in Figures 8
and 9 as well as failure cases in Figure 10.

#### 827 C.4. Qualitative error analysis

Why does adding programmatic LATTE-trace help on 828 MathVista but hurt MMVet performance? We observe 829 that adding programmatic LATTE-trace can result in up to 830 831 3% gain on MathVista and 9% drop on MMVet. Upon analysis, we discover that programmatic LATTE-trace improves 832 the general VQA split in MathVista the most by almost 9%. 833 This is because LOCALIZE is helpful for these questions, 834 and our programmatic data includes many LOCALIZE in-835 stances that allow LATTE to learn to use it effectively (Fig-836 837 ure 6). Conversely, programmatic data hurts LATTE's performance on MMVet most likely due to the model's worse838reasoning ability as a result of the simple and rigid thoughts839generated with templates in our programmatic data (Figure 6).840

#### C.5. Additional quantitative results

We report additional quantitative results of data ablations on<br/>Mantis-CLIP in Table 11, where we see the same trends we<br/>observe with Mantis-SigLIP: the smallest dataset of 293K<br/>LATTE-trace examples leads to the highest absolute perfor-<br/>mance and gain compared to other datasets with a mix of<br/>LATTE-trace, CoT, and/or Direct examples at larger scales.843<br/>844

#### Table 6. VisCoT results

Method	A-OKVQA	BLINK	MathVista	MMMU	MMStar	MMVet	MMVP	RealWorldQA	Avg
VPD-LLaVA	80.7	43.9	35.1	40.8	40.1	40.6	61.3	45.1	48.5
VisCoT-LLaVA	67.9	39.4	12.9	36.1	36.0	26.7	61.3	34.0	39.3

Visual-CoT Performance. We experimentally com-849 pare LATTE to Visual-CoT. We finetune Mantis-LLaVA-850 Pretrained (LLama3+SigLIP) with Visual CoT and compare 851 its performance with LATTE (Table 6). We use 413K ex-852 amples where the bounding boxes are valid and within the 853 image. We find that the models trained with Visual CoT 854 data achieve an average accuracy of 39.3% (much lower as 855 Visual COT's data are mostly Text/Doc images and con-856 tain only bboxes without natural language thoughts) on the 857 benchmarks. 858

#### Table 7. LATTE vs. Direct inference

Test prompt	A-OKVQA	BLINK	MathVista	MMMU	MMStar	MMVet	MMVP	RealWorldQA	Avg
LATTE	82.1	45.9	37.0	39.8	43.4	46.7	66.0	41.6	50.3
Direct	81.9	46.6	31.2	42.1	39.1	35.0	68.3	40.2	48.0

Performance gain with LATTE inference. We compare 859 the model's performance when trained with a random mix 860 of 293K LATTE-traces and Direct data (1:1) and tested 861 with LATTE format vs. Direct prompt. We find that the 862 model achieves an average of 50.3% when tested following 863 LATTE format vs. 48% with the Direct prompt (Table 7), 864 suggesting that reasoning with vision specialists at infer-865 ence time improves model's performance. 866

Hyperparameter tuningAdditional gains can be867achieved by tuning the vision encoder, training with<br/>a smaller learning rate or for more epochs. Last but<br/>not least, our hyperparameter tuning experiments with868

CVPR #36

842

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907 908

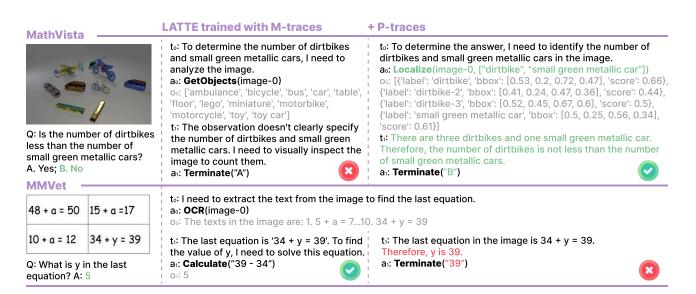


Figure 6. We present examples of LATTE success and failure after adding programmatic data to the fine-tuning dataset.

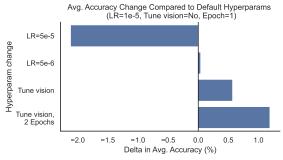


Figure 7. **Hyperparameter ablations.** Additional gains can be achieved with a smaller learning rate for the language model, tuning the vision encoder, and training for more epochs.

LLaVa-OV-Stage1.5 suggest that we can further improve
the model's absolute performance by tuning the vision encoder, training with a smaller learning rate and/or for longer
epochs (Figure 7).

#### 875 D. Dataset and model comparison

We summarize the differences between our work and the
other CoT datasets including ScienceQA, M<sup>3</sup>COT, Visual
CoT, VPD, and LLaVA-CoT (which is concurrent and unpublished work), in Table 8.

#### 880 E. Related work

We contexualize our work on language-only and multi-modal tool use and multi-modal language models.

Language-only tool use. Augmenting LLMs with external tools and APIs can significantly enhance their problemsolving capabilities, especially for tasks that require specialized knowledge or interaction with external resources
[60]. Examples include web searching [53, 81], mathematical calculations [13, 22], code interpretation [18, 87], and other domain-specific operations that are challenging for LLMs to accomplish solely with their intrinsic knowledge.

Researchers have mainly explored two approaches to enable LLMs to effectively leverage external tools. The first approach is through in-context learning, where instructions on tool usage and contextually relevant examples are provided directly in the prompt [11, 20, 49, 55, 65, 82]. The second approach involves training models to specialize in calling external functions [5, 56, 59, 62, 71, 85]. For instance, Toolformer [62] proposed training LLMs to call APIs precisely to solve complex tasks in a self-supervised manner. More recently, xLAM unified data formats across various environments, and achieved more powerful LLMs specialized in function calling that are also referred to as large action models [85]. These works demonstrated the effectiveness of fine-tuning LLMs to execel at calling tools, which inspired our work on tuning multi-modal language models to be multi-modal action models.

Multi-modal tool-use. Unlike language-only tool use, 909 most works in the multi-modal space have only explored 910 training-free methods, which adapt existing LLMs to use 911 multi-modal tools via zero-shot/few-shot prompting [21, 912 70, 80]. Visprog first demonstrates the potential of leverag-913 ing LLMs for visual question answering by prompting them 914 to generate programs of predefined modules for image un-915 derstanding and manipulation [21]. Simlarly, ViperGPT de-916 fines a comprehensive list of APIs and prompts code LLMs 917 to generate Python programs for solving visual queries [70]. 918 In addition to question answering, m&ms further inves-919 tigated LLMs' abilities in planning and executing open-920 ended multi-modal tasks with external tools [51]. Be-921 yond LLMs, more recent efforts such as VisualSketchPad 922

970

973

CVPR 2025 Submission #36.	CONFIDENTIAL	<b>REVIEW COPY.</b>	<b>DO NOT DISTRIBUTE.</b>
---------------------------	--------------	---------------------	---------------------------

Title			Datas	set			Mod	lel
The	Training set size	Domains	Data source number	Tools involved in reasoning chain	Multi- image questions?	Multimodal reasoning chain?	Inference- time action calling?	Multi- image support?
Science QA	12.6K	Science	1	×	×	×	×	×
МЗСоТ	7.8K	Science, Math, Commonsense	2	×	×	×	×	×
VPD	90K	General, OCR, Compositional, Knowledge, Count- ing	6	OCR, Object De- tection, Depth Esti- mation, VQA, Cap- tioning, Knowledge	×	×	×	×
VisualCoT	98K (+340K with bboxes but no thoughts)	General, Doc, Fine- grained, Relation, Chart	12	Localize, Crop	×	×	✓	1
LLaVA- o1	100K	Science, Math, General, Science, Doc	10	×	×	X	×	×
LATTE	293K (+over 1M program gener- ated reasoning traces)	Includes all of the above	31	15 tools (see Sec- tion 2)	<b>v</b>	<b>v</b>		<b>√</b>

Table 8.	Dataset	and	model	comparison.
----------	---------	-----	-------	-------------

also augmented multi-modal models with external tools and
showcased success [24]. However, in-context tool use only
works well with powerful proprietary multi-modal models
such as GPT-40; all open-source multi-modal models fail to
perform tool use well without fine-tuning.

928 There are only a few works on training multi-modal models to be better at tool use. LLaVa-Plus was the first 929 930 work that shows the possibility of training a multi-modal model to use vision specialists [43]. However, it exhibits 931 poor performance due to the weakness of the model back-932 bone and the small scale and suboptimal quality of the data 933 934 [43]. Visual Program Distillation distills tool-use and rea-935 soning abilities into a multi-modal model with chain-of-936 thought (CoT) data obtained from programs [23]. Nonetheless, this model is unable to use tools at inference and is 937 limited to simple question answering tasks [23]. Similarly, 938 the more recent work Visual CoT introduces a new syn-939 940 thetic CoT dataset for training multi-modal models for en-941 hanced reasoning [64]. The closest work to ours is Cog-CoM, which identifies 6 useful manipulations and trains 942 943 multi-modal models with synthetic chain-of-manipulation (CoM) data [58]. Nonetheless, the manipulations are lim-944 945 ited and useful for only detailed visual question answering 946 and visual grounding, and the authors have only experimented with adding 70K CoM data to 500K QA examples. 947 None of these works studied or improved the quantity and 948 quality of these chains to enhance multi-modal models' per-949 formance across diverse complex multi-modal tasks, which 950 951 our work focuses on.

Multi-modal language models. Most recently, there have 952 been many advances made on open-source multi-modal 953 models [2, 3, 6, 7, 9, 10, 12, 14, 26, 34, 37-39, 39-954 42, 45, 47, 50, 52, 54, 57, 66-69, 72, 74-76, 78]. These 955 efforts include training multi-modal models to take in mul-956 tiple images, engage in multi-turn conversations, and even 957 understand videos [28, 35, 42]. For example, LLaVA-Next 958 achieves strong multi-image understanding through large-959 scale interleaved visual instruction tuning with M4-Instruct 960 [42]. Similarly, Mantis introduces a new large-scale multi-961 image instruction tuning dataset Mantis-Instruct for multi-962 image training [28]. These efforts pave the foundation for 963 our work on training multi-modal models to learn from 964 multi-image interleaved chains-of-thought-and-action data. 965

## F. Model-based data generation 966

#### **F.1.** Generation prompt

We present the full data generation prompt used in our 968 model-based data generation pipeline in Listing 2. 969

#### F.2. Dataset statistics

We present a table with detailed statistics of the LATTE-	971
trace 293K dataset in Table 12.	972

#### G. Action implementation

Our Python implementation of all actions can be found in Uisting 1. 975

976 H. Programmatic data generation

#### 977 H.1. QA and action templates

We present the question-answer and corresponding action
templates used in our programatic data generation in Table 13. We design 16 different question templates for both
single-image and multi-image examples that cover 5 capabilities: attribute recognition, counting, 2D and 3D spatial
understanding, and multi-image understanding.

### 984 H.2. Thought templates

We also present the five thought templates in Listing 3 wedefine for each action, where one of them is randomly sam-pled and used during generation.

#### 988 H.3. Example action distribution

We plot example distributions of all actions before and afteradding programmatic LATTE-trace 73K data in Figure 11.

#### 991 I. Additional training details

We report additional training hyperparameters for Mantismodels and LLaVA-OV in Table 9 and 10 respectively.

# Table 9. Additional training hyperparameters for Mantis-SigLIP and Mantis-CLIP.

Name	Value
bf16	TRUE
tf32	True
per_device_train_batch_size	1
per_device_eval_batch_size	1
gradient_accumulation_steps	16
weight_decay	0
warmup_ratio	0.03
lr_scheduler_type	cosine
lora_enabled	FALSE
qlora_enabled	FALSE
max_seq_len	8192

Table 10. Additional training hyperparameters for LLaVA-OV.

Name	Value
bf16	TRUE
tf32	True
mm_vision_tower_lr	2.00E-06
mm_projector_type	mlp2x_gelu
mm_vision_select_layer	-2
image_aspect_ratio	anyres_max_9
image_grid_pinpoints	"(1x1),,(6x6)"
mm_patch_merge_type	spatial_unpad
per_device_train_batch_size	1
per_device_eval_batch_size	1
gradient_accumulation_steps	16
weight_decay	0
warmup_ratio	0.03
lr_scheduler_type	cosine
model_max_length	8192

### J. Additional evaluation details

994

We present additional inference and evaluation details in Table 14 and the LLM judge prompts used for MMVet and MathVista from VLMEvalKit [16] in Listings 4 and 5. 997

	e												
Data source	Final data format	Size	Model	A-OKVQA	BLINK	MathVista	MMMU	MMStar	MMVet	MMVP	RealWorldQA	Avg	Delta
	Direct LATTE-trace	293K	Mantis-CLIP LATTE	80.7 81.1	45.8 <b>49.6</b>	33.1 <b>36.6</b>	42.2 42.8	36.7 <b>40.8</b>	28.9 <b>45.2</b>	62.7 63.3	<b>52.3</b> 51.1	47.8 <b>51.3</b>	3.5
All datasets	Direct LATTE-trace+ CoT	580K	Mantis-CLIP LATTE	82.0 <b>82.6</b>	47.2 47.7	31.5 35.9	40.6 38.3	38.3 39.6	31.4 43.4	<b>63.3</b> 61.0	49.0 51.1	47.9 49.9	2.0
	Direct LATTE-trace+ CoT+ Direct	815K	Mantis-CLIP LATTE	81.0 81.3	47.3 43.1	32.9 32.7	47.3 <b>48.0</b>	38.1 37.7	29.1 31.0	63.0 58.7	51.0 47.6	48.7 47.5	-1.2
LATTE-useful datasets	Direct LATTE-trace+ CoT+ Direct	566K	Mantis-CLIP LATTE	79.1 80.7	44.1 41.6	30.9 31.5	39.4 41.0	35.5 38.2	30.1 36.1	62.7 59.0	49.8 48.2	46.5 47.0	0.6

Table 11. Additional Results on Model-generated data ablations with Mantis-CLIP. We observe similar results of data ablations on Mantis-CLIP as on Mantis-SigLIP.

Table 12. Detailed statistics of the LATTE-trace 293K dataset.

Original Dataset	Subset	Instance #	Max Image #	Avg Image #	Max Turn #	Avg Turn #
	ai2d	694	2	1	6	2
	aokvqa	6322	5	1.1	7	2.1
	chartqa	4726	2	1	10	2.1
	cocoqa	13289	3	1.1	4	2
	dvqa	2158	2	1	7	2.5
	iconqa	3791	3	1.1	5	2.2
	infographic_vqa	3822	3	1	9	2.3
	mimic_cgd	6899	6	2.1	7	2.8
	nlvr2	9716	4	2.1	6	2.5
Cauldron	ocrvqa	22991	2	1	7	2
	scienceqa	850	2	1	6	2.3
	st_vqa	11322	3	1	8	2
	tabmwp	14548	1	1	10	2.5
	tallyqa	16171	3	1.4	5	2.1
	textvqa	15475	5	1	6	2.1
	visual7w	4773	3	1.1	5	2.1
	vqarad	115	2	1	4	2.2
	vqav2	13394	5	1.2	6	2.1
	vsr	1864	2	1.2	4	2.1
	birds-to-words	742	4	2	5	2.7
	coinstruct	31773	8	2.3	8	2.2
	contrastive_caption	4296	8	3.6	6	2 2
	dreamsim	1738	3	3	3	2
	iconqa	6660	7	2.6	6	2.2
Mantis	imagecode	559	18	10.1	10	3.1
Waltus	lrv_multi	3401	9	3.3	6	2.2
	multi_vqa	2089	7	3.8	8	2.6
	nlvr2	5436	4	2	5	2.5
	spot-the-diff	2591	5	2.8	8	3
	nextqa	3057	15	8.2	9	2.3
	llava_665k_multi	77843	11	2.2	10	2.1
]	Fotal	293105	18	2.2	10	2.3

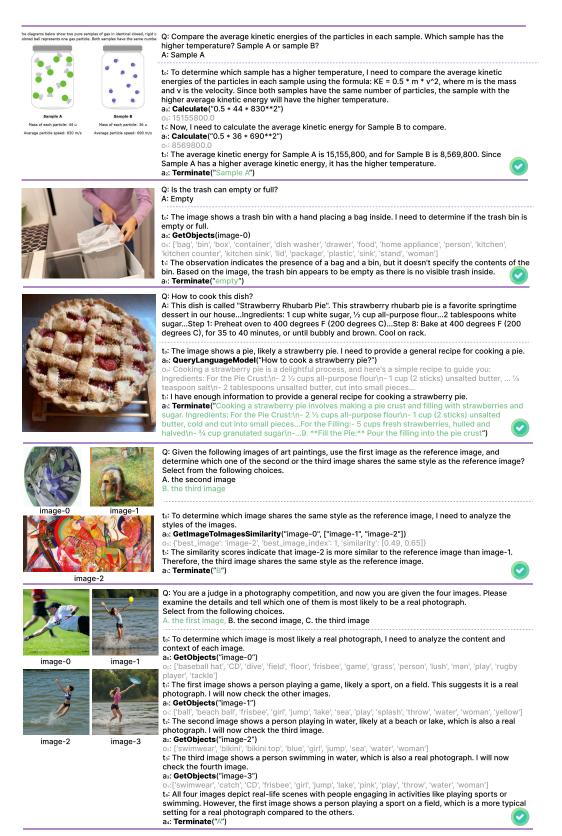


Figure 8. We present additional successful examples of LATTE.

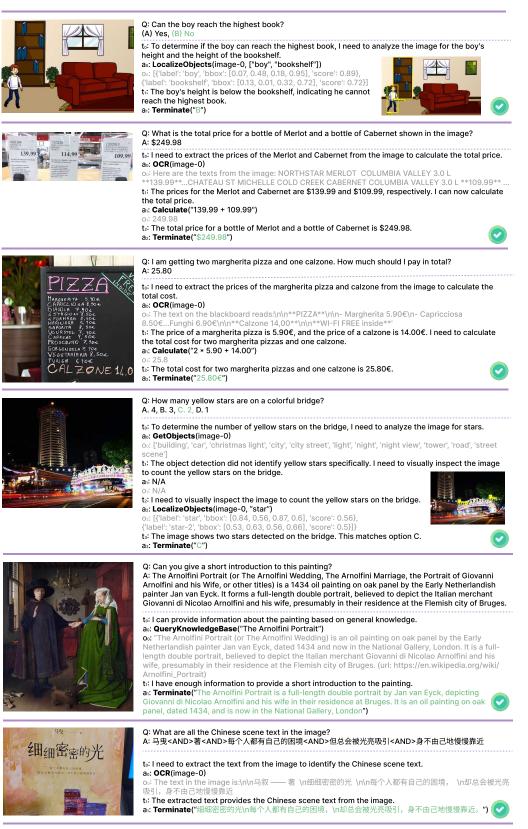


Figure 9. We present additional successful examples of LATTE.

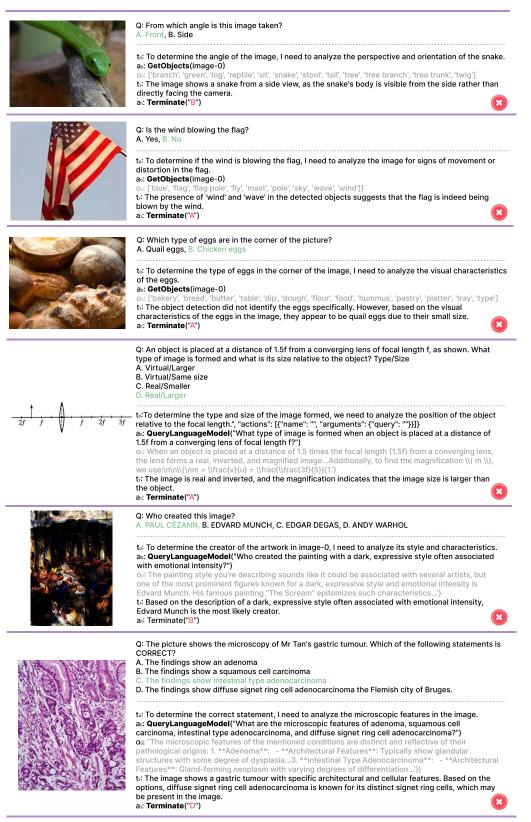


Figure 10. We present additional failure cases of LATTE.

#### Table 13. Templates for programmatic data generation.

# of input images	Capabilities	Question Template	Action Template		
	Counting	How many {object} are there? inting Among {objects}, which is the most frequent object? Among {objects}, which object appears the least?			
	Counting, Attribute recognition	How many {attribute} {object} are there?	LocalizeObjects		
1	2D spatial reasoning	Among {objects}, which is on the most left side? Among {objects}, which is on the most right side? Among {objects}, which is on the most top side? Among {objects}, which is on the most bottom side?	-		
	3D spatial reasoning	Which of {objects} is closer? Which of {objects} is farther?	LocalizeObjects, EstimateRegionDepth x2 OR, EstimateObjectDepth x2		
2-3	Multi-image understanding Multi-image understanding, Counting Multi-image understanding, Counting Multi-image understanding, Counting Multi-image understanding, Attribute recognition Multi-image understanding, Attribute recognition, Counting	Which image has {object}? How many {object} are in in these images? Which image has most {object}? Which image has least {object}? Which image has {attribute} {object}? How many {attribute} {object} in these images?	LocalizeObjects x N		

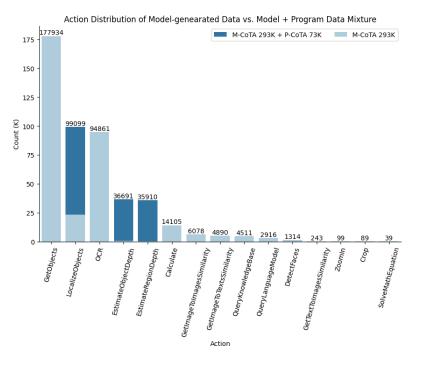


Figure 11. Action distribution of model-generated data vs. model and program data mixtures.

Table 14. Additional inference and evaluation details.

Stage	Name	Value
Inference	do_sample temperature max_new_tokens max_consecutive_auto_reply	FALSE 0 2000 10
Evaluation	llm judge for multiple choice & yes/no questions llm judge for short answer questions (i.e. MMVet, MathVista) llm judge max_new_tokens llm judge retry	gpt-3.5-turbo-0125 gpt-4-1106-preview 2048 5

```
998
      1 class BaseAction:
999
            .....
      2
1000
            This is the Action class for agent to use.
      3
1001
            Using this Action class to wrap APIs, tools, models as an Action of an agent
1002
            .....
      5
1003
      6
1004
            def ___init___(
      7
1005
                self,
      8
1006
                id: int,
      9
1007
               description: str = "",
     10
1008
     11
               args_spec: dict = {},
1009
               rets_spec: dict = {},
     12
1010 13
               examples: List = []
1011 14
           ) -> None:
1012 15
                ....
1013
                the agent action should be connected with data and env
     16
1014 17
                Args:
1015 18
                   id: the id of the action
1016 19
                   description: the description of the action
1017 20
1018 21
                   args_spec: the specification of the arguments
                    rets_spec: the specification of the returns
1019 22
                   examples: a list of examples of the action
                .....
1020 23
1021 24
                self.name = self.__class__._name__
1022
     25
                self.id = id
1023
                self.description = description
     26
1024 27
                self.args_spec = args_spec
1025 28
                self.rets_spec = rets_spec
1026 29
                self.examples = examples
1027
                self.device = "cuda:0" if torch.cuda.is_available() else "cpu"
      30
1028
     31
1029 32
            def __call__(self, **kwargs) -> str:
1030 33
1031 34
                implement the Action as
1032
                ......
     35
1033
     36
                raise NotImplementedError
1034 37
1035 38
1036 39 class OCR(BaseAction):
          def __init__(self, id) -> None:
1037
     40
1038
                description = "Extract texts from an image or return an empty string if no text is in the image
     41
1039
            . Note that the texts extracted may be incorrect or in the wrong order. It should be used as a
1040
            reference only."
1041 42
               args_spec = {"image": "the image to extract texts from."}
1042
                rets_spec = {"text": "the texts extracted from the image."}
     43
1043 44
                examples = [{"name": "OCR", "arguments": {"image": "image-0"}}]
1044 45
1045 46
                super().__init__(
1046
     47
                    id=id, description=description, args_spec=args_spec, rets_spec=rets_spec, examples=examples
1047
     48
                )
1048 49
1049 50
            def __call__(self, image, tool_version=LATEST_GPT_MODEL_ID):
1050 51
               if tool_version == "easyocr":
1051
     52
                    import easyocr
1052
     53
                    import io
1053
                    reader = easyocr.Reader(["en"]) # Load the OCR model into memory
     54
1054 55
                   image = image_processing(image)
1055 56
                    if isinstance(image, str):
1056
     57
                        # If image is a path, use it directly
1057
     58
                        image_path_or_bytes = (
1058 59
                            image if os.path.exists(image) else get_full_path_data(image)
1059 60
                        )
1060 61
                    else:
1061
                        # If image is an Image object, convert it to a bytes stream
     62
1062
     63
                        buffer = io.BytesIO()
1063 64
                        image.save(buffer, format="JPEG")
1064
                       buffer.seek(0)
      65
```

```
image_path_or_bytes = buffer
                                                                                                                      1065
66
                                                                                                                      1066
67
               result = reader.readtext(image_path_or_bytes)
                                                                                                                      1067
68
69
               result_text = [text for _, text, _ in result]
                                                                                                                      1068
               result_formatted = {"text": ", ".join(result_text)}
                                                                                                                      1069
70
           else:
                                                                                                                      1070
                                                                                                                      1071
72
               from openai import OpenAI
               import base64
                                                                                                                      1072
73
               client = OpenAI(api_key=os.getenv("OPENAI_API_KEY"))
                                                                                                                      1073
74
                                                                                                                      1074
76
               def encode_image(image_path):
                                                                                                                      1075
                    with open(image_path, "rb") as image_file:
                                                                                                                      1076
77
                                                                                                                      1077
                        return base64.b64encode(image_file.read()).decode('utf-8')
78
                                                                                                                      1078
79
                                                                                                                      1079
80
               image_path = image_processing(image, return_path=True)
               base64_image = encode_image(image_path)
                                                                                                                      1080
81
                                                                                                                      1081
82
               response = client.chat.completions.create(
                                                                                                                      1082
83
                   model=tool_version,
                                                                                                                      1083
84
                                                                                                                      1084
                    messages=[
85
86
                        {
                                                                                                                      1085
                             "role" : "user",
                                                                                                                      1086
87
                             "content": [
                                                                                                                      1087
88
                                 {"type": "text", "text": f"What are the texts in the image?"},
                                                                                                                      1088
89
                                                                                                                      1089
90
                                 {
                                     "type"
                                               : "image_url",
                                                                                                                      1090
91
                                     "image_url": {
                                                                                                                      1091
92
                                          "url": f"data:image/jpeg;base64,{base64_image}",
                                                                                                                      1092
93
                                                                                                                      1093
94
                                     },
                                                                                                                      1094
95
                                 },
                                                                                                                      1095
                            ],
96
                                                                                                                      1096
                        }
97
                                                                                                                      1097
98
                    1.
                                                                                                                      1098
                    max_tokens=300.
99
                                                                                                                      1099
               )
100
                                                                                                                      1100
101
               result_formatted = {"text": response.choices[0].message.content}
                                                                                                                      1101
102
103
           return result_formatted
                                                                                                                      1102
                                                                                                                      1103
104
105
                                                                                                                      1104
                                                                                                                      1105
   class GetObjects(BaseAction):
106
                                                                                                                      1106
107
      def __init__(self, id) -> None:
           description = "Using this function to get objects in an image."
                                                                                                                      1107
108
           args_spec = {"image": "the image to get objects from."}
                                                                                                                      1108
109
           rets_spec = {"objects": "the objects detected in the image."}
                                                                                                                      1109
110
           examples = [{"name": "GetObjects", "arguments": {"image": "image-0"}}]
                                                                                                                      1110
                                                                                                                      1111
           super().__init__(
                                                                                                                      1112
               id=id, description=description, args_spec=args_spec, rets_spec=rets_spec, examples=examples
                                                                                                                      1113
114
                                                                                                                      1114
           )
                                                                                                                      1115
116
       def __call__(self, image, tool_version="https://huggingface.co/xinyu1205/recognize-anything-plus-
                                                                                                                      1116
       model/resolve/main/ram_plus_swin_large_14m.pth?download=true"):
                                                                                                                      1117
                                                                                                                      1118
           from ram.models import ram_plus
118
                                                                                                                      1119
119
           from ram import get_transform, inference_ram_openset as inference
                                                                                                                      1120
120
           model_path_or_url = tool_version
                                                                                                                      1121
                                                                                                                      1122
           image_size = 384
                                                                                                                      1123
123
           transform = get_transform(image_size=image_size)
                                                                                                                      1124
124
           vit size = "swin l"
                                                                                                                      1125
           # load model
                                                                                                                      1126
126
                                                                                                                      1127
           model = ram_plus(pretrained=model_path_or_url,
                                                                                                                      1128
                            image_size=image_size,
128
129
                            vit=vit_size)
                                                                                                                      1129
                                                                                                                      1130
           model.eval()
130
           model = model.to(self.device)
                                                                                                                      1131
```

**1132** 132 image = image\_processing(image) **1133** 133 image = transform(image).unsqueeze(0).to(self.device) **1134** 134 tags = inference(image, model) **1135** 135 objs = tags.split(" | ") **1136** 136 return {"objects": objs} 1137 137 **1138** 138 1139 139 class VisualizeRegionsOnImage(BaseAction): **1140** 140 def \_\_init\_\_(self, id) -> None: **1141** 141 description = "Using this function to label regions on an image." 1142 142 args\_spec = {"image": "the image to label.", **1143** 143 "regions": "the regions to label on the image, where each region is represented by 1144 a dictionary with the region's bounding box and label text (can be empty string).", **1145** 144 "color": "an optional argument that specifies the color of the bounding box." **1146** 145 } 1147 rets\_spec = {"image": "the image with regions labeled."} 146 **1148** 147 examples = [ **1149** 148 {"name": "VisualizeRegionsOnImage", "arguments": {"image": "image-0", "regions": [{"label": 1150 "", "bbox": [0.3, 0.2, 0.5, 0.4]}]}}, **1151** 149 {"name": "VisualizeRegionsOnImage", "arguments": {"image": "image-0", "regions": [{"label": 1152 "cat", "bbox": [0.3, 0.2, 0.5, 0.4]}], "color": "red"}} **1153** 150 1 **1154** 151 **1155** 152 super().\_\_init\_\_( **1156** 153 id=id, description=description, args\_spec=args\_spec, rets\_spec=rets\_spec, examples=examples **1157** 154 **1158** 155 **1159** 156 def \_\_\_call\_\_ (self, image, regions: List[Region], color='yellow', width=4): **1160** 157 image = image\_processing(image) 1161 158 text\_color = 'black' **1162** 159 W,H = image.size **1163** 160 img1 = image.copy() **1164** 161 draw = ImageDraw.Draw(img1) **1165** 162 font = ImageFont.truetype('/usr/share/fonts/truetype/dejavu/DejaVuSansMono-Bold.ttf', 16) **1166** 163 for i, obj in enumerate(regions): **1167** 164 bbox = obj['bbox'] **1168** 165 bbox = bbox[0] \* W, bbox[1] \* H, bbox[2] \* W, bbox[3] \* H**1169** 166 draw.rectangle(bbox, outline=color, width=width) **1170** 167 x1, y1, x2, y2 = bbox 1171 168 label = obj['label'] if "label" in obj else "" **1172** 169 w,h = font.getsize(label) **1173** 170 if x1 + w > W or y2 + h > H: **1174** 171 draw.rectangle((x1, y2 - h, x1 + w, y2), fill=color) **1175** 172 draw.text((x1, y2-h),label,fill=text\_color,font=font) **1176** 173 else: **1177** 174 draw.rectangle((x1, y2, x1 + w, y2 + h), fill=color) **1178** 175 draw.text((x1, y2),label,fill=text\_color,font=font) **1179** 176 return {"image": img1} **1180** 177 1181 178 1182 179 class LocalizeObjects(BaseAction): **1183** 180 def \_\_init\_\_(self, id) -> None: **1184** 181 description = "Localize one or multiple objects/regions with bounding boxes. This tool may 1185 output objects that don't exist or miss objects that do. You should use the output only as weak 1186 evidence for reference. When answering questions about the image, you should double-check the 1187 detected objects. You should be especially cautious about the total number of regions detected, 1188 which can be more or less than the actual number." **1189** 182 args\_spec = { **1190** 183 "image": "the image to localize objects/regions in.", "objects": "a list of object names to localize. e.g. ['dog', 'cat', 'person']. the model **1191** 184 1192 might not be able to detect rare objects or objects with complex descriptionriptions." **1193** 185 } **1194** 186 rets\_spec = {"image": "the image with objects localized and visualized on it.", "regions": "the 1195 regions of interests localized in the image, where each region is represented by a dictionary with 1196 the region's label text, bounding box and confidence score. The confidence score is between 0 and 1197 1, where 1 means the model is very confident. Note that both the bounding boxes and confidence 1198 scores can be unreliable and should only be used as reference."}

```
examples = [{"name": "LocalizeObjects", "arguments": {"image": "image-0", "objects": ["dog", "
                                                                                                                    1199
187
       cat"]}}]
                                                                                                                    1200
                                                                                                                    1201
188
189
           super().__init__(
                                                                                                                    1202
                                                                                                                    1203
               id=id, description=description, args_spec=args_spec, rets_spec=rets_spec, examples=examples
190
                                                                                                                    1204
191
           )
                                                                                                                    1205
192
       def call (self, image, objects: List[str]):
                                                                                                                    1206
193
           from groundingdino.util.inference import load_model, load_image, predict, annotate
                                                                                                                    1207
194
                                                                                                                    1208
195
           import cv2
           text = ". ".join(objects)
                                                                                                                    1209
196
           model = load_model("/user/mma/mma/GroundingDINO/groundingdino/config/GroundingDINO_SwinT_OGC.py
                                                                                                                    1210
197
       н.,
                                                                                                                    1211
                               "/user/mma/mma/GroundingDINO/weights/groundingdino_swint_ogc.pth",
                                                                                                                    1212
198
                                                                                                                    1213
199
                               device=self.device)
           BOX_TRESHOLD = 0.35
                                                                                                                    1214
200
                                                                                                                    1215
           TEXT TRESHOLD = 0.25
201
           image_path = image_processing(image, return_path=True)
                                                                                                                    1216
202
           original_image = image_processing(image)
                                                                                                                    1217
203
           image_source, image = load_image(image_path)
                                                                                                                    1218
204
205
                                                                                                                    1219
                                                                                                                    1220
206
           boxes, logits, phrases = predict(
               model=model.
                                                                                                                    1221
207
                                                                                                                    1222
               image=image,
208
               caption=text,
                                                                                                                    1223
209
                                                                                                                    1224
               box_threshold=BOX_TRESHOLD,
210
               text_threshold=TEXT_TRESHOLD
                                                                                                                    1225
                                                                                                                    1226
           )
                                                                                                                    1227
                                                                                                                    1228
           objects = []
214
                                                                                                                    1229
           obj_cnt = \{\}
                                                                                                                    1230
           for i in range(len(boxes)):
216
               xyxy = box_convert(boxes=boxes[i], in_fmt="cxcywh", out_fmt="xyxy").numpy()
                                                                                                                    1231
                                                                                                                    1232
218
               bbox = [round(val, 2) for val in list(xyxy)]
                                                                                                                    1233
               score = round(logits[i].item(), 2)
219
               phrase = phrases[i]
                                                                                                                    1234
220
               obj_cnt[phrase] = obj_cnt.get(phrase, 0) + 1
                                                                                                                    1235
               phrase = f"{phrase}-{obj_cnt[phrase]}" if obj_cnt[phrase] > 1 else phrase
                                                                                                                    1236
               objects.append({"label": phrase, "bbox": bbox, "score": score})
                                                                                                                    1237
224
           visualize = VisualizeRegionsOnImage(0)
                                                                                                                    1238
                                                                                                                    1239
           results = visualize(image=original_image, regions=objects)
                                                                                                                    1240
226
           tagged_image = results["image"]
           results_formatted = {"regions": objects, "image": tagged_image}
                                                                                                                    1241
           return results_formatted
                                                                                                                    1242
228
                                                                                                                    1243
229
                                                                                                                    1244
230
                                                                                                                    1245
231 class Crop(BaseAction):
                                                                                                                    1246
       def __init__(self, id) -> None:
           description = "Crop an image with the bounding box. It labels the cropped region with a
                                                                                                                    1247
       bounding box and crops the region with some margins around the bounding box to help with contextual
                                                                                                                    1248
                                                                                                                    1249
        understanding of the region."
                                                                                                                    1250
           args_spec = {
234
               "image": "the image to crop.",
                                                                                                                    1251
               "bbox": "the bbox to crop. It should be a list of [left, top, right, bottom], where each
                                                                                                                    1252
236
                                                                                                                    1253
       value is a float between 0 and 1 to represent the percentage of the image width/height and how far
                                                                                                                    1254
       it is from the top left corner at [0, 0].",
                                                                                                                    1255
           }
                                                                                                                    1256
           rets_spec = {"image": "the cropped image."}
238
           examples = [{"name": "Crop", "arguments": {"image": "image-0", "bbox": [0.33, 0.21, 0.58,
                                                                                                                    1257
239
                                                                                                                    1258
       0.46]}}]
                                                                                                                    1259
240
           super().__init__(
                                                                                                                    1260
241
                                                                                                                    1261
               id=id, description=description, args_spec=args_spec, rets_spec=rets_spec, examples=examples
242
                                                                                                                    1262
243
           )
244
                                                                                                                    1263
       def __call__(self, image, bbox):
                                                                                                                    1264
245
           image = image_processing(image)
                                                                                                                    1265
246
```

```
1266 247
1267 248
                if isinstance(bbox, str):
1268 249
                     try:
1269 250
                        bbox = ast.literal_eval(bbox)
1270 251
                    except:
1271 252
                        bbox = []
1272 253
1273 254
                use_percent = (all(x \le 1.0 \text{ for } x \text{ in } bbox))
1274 255
                if not use_percent:
1275 256
                    raise ValueError("Bounding box coordinates must be between 0 and 1.")
1276 257
1277 258
                visualize = VisualizeRegionsOnImage(0)
1278 259
                results = visualize(image=image, regions=[{"label": "", "bbox": bbox}])
1279 260
                image = results["image"]
1280 261
1281 262
                W, H = image.size
1282 263
                bbox = [bbox[0] * W, bbox[1] * H, bbox[2] * W, bbox[3] * H]
                bbox = expand_bbox(bbox, image.size)
1283 264
1284 265
                out_img = image.crop(bbox)
1285 266
                return {"image": out_img}
1286 267
1287 268
1288 269 class ZoomIn(BaseAction):
1289 270
          def __init__(self, id) -> None:
1290 271
                description = "Zoom in on a region of the input image. This tool first crops the specified
1291
            region from the image with the bounding box and then resizes the cropped region to create the zoom
1292
            effect. It also adds some margins around the cropped region to help with contextual understanding
1293
            of the region."
1294 272
                args\_spec = {
1295 273
                     "image": "the image to zoom in on.",
1296 274
                    "bbox": "The bbox should be a list of [left, top, right, bottom], where each value is a
1297
            float between 0 and 1 to represent the percentage of the image width/height and how far it is from
1298
            the top left corner at [0, 0].",
1299 275
                    "zoom_factor": "the factor to zoom in by. It should be greater than 1.",
1300 276
                }
1301 277
                rets_spec = {"image": "the zoomed in image."}
1302 278
                examples = [
1303 279
                    {"name": "ZoomIn", "arguments": {"image": "image-0", "bbox": [0.4, 0.3, 0.5, 0.4], "
1304
            zoom_factor": 2}},
1305 280
                1
1306 281
1307 282
                super(). init (
1308 283
                    id=id, description=description, args_spec=args_spec, rets_spec=rets_spec, examples=examples
1309 284
                )
1310 285
1311 286
            def __call__(self, image, bbox, zoom_factor):
1312 287
                if zoom_factor <= 1:</pre>
1313 288
                    raise ValueError ("Zoom factor must be greater than 1 to zoom in")
1314 289
1315 290
                image = image_processing(image)
1316 291
                use_percent = (all(x <= 1.0 for x in bbox))</pre>
1317 292
                if not use_percent:
1318 293
                    raise ValueError("Bounding box coordinates must be between 0 and 1.")
1319 294
1320 295
                crop = Crop(0)
1321 296
                cropped_image = crop(image, bbox)["image"]
1322 297
1323 298
                W, H = cropped_image.size
1324 299
1325 300
                # Calculate the size of the zoomed image
1326 301
                new_width = int(W * zoom_factor)
1327 302
                new_height = int(H * zoom_factor)
1328 303
1329 304
                # Resize the cropped image to create the zoom effect
1330 305
                zoomed_image = cropped_image.resize((new_width, new_height), Image.LANCZOS)
1331 306
                return {'image': zoomed_image}
1332 307
```

class GetImageToImagesSimilarity(BaseAction): def \_\_init\_\_(self, id) -> None: description = "Get the similarity between one image and a list of other images. Note that this similarity score may not be accurate and should be used as a reference only." args\_spec = { "image": "the reference image.", "other\_images": "the other images to compare to the reference image.", } rets\_spec = {"similarity": "the CLIP similarity scores between the reference image and the other images.", "best\_image\_index": "the index of the most similar image."} examples = [ {"name": "GetImageToImagesSimilarity", "arguments": {"image": "image-0", "other\_images": [" image-1", "image-2"]}} ] super().\_\_init\_\_( id=id, description=description, args\_spec=args\_spec, rets\_spec=rets\_spec, examples=examples def \_\_call\_\_(self, image, other\_images, tool\_version='ViT-H-14-378-quickgelu', other\_images\_raw= None): import torch import open\_clip original\_images = other\_images\_raw model, \_, preprocess = open\_clip.create\_model\_and\_transforms(tool\_version, pretrained='dfn5b') model.eval() image = image\_processing(image) images = [image\_processing(image) for image in other\_images] image = preprocess(image).unsqueeze(0) images = torch.stack([preprocess(image) for image in images]) with torch.no\_grad(), torch.cuda.amp.autocast(): image1\_features = model.encode\_image(image) image2\_features = model.encode\_image(images) image1\_features /= image1\_features.norm(dim=-1, keepdim=True) image2\_features /= image2\_features.norm(dim=-1, keepdim=True) probs = image1\_features @ image2\_features.T sim\_scores = [round(sim\_score, 2) for sim\_score in probs[0].tolist()] best\_image\_match = torch.argmax(probs).item() return {'similarity': sim\_scores, "best\_image\_index": best\_image\_match, "best\_image": original\_images[best\_image\_match]} class GetImageToTextsSimilarity(BaseAction): def \_\_init\_\_(self, id) -> None: description = "Get the similarity between one image and a list of texts. Note that this similarity score may not be accurate and should be used as a reference only." args\_spec = { "image": "the reference image.", "texts": "a list of texts to compare to the reference image.", } rets\_spec = {"similarity": "the CLIP similarity between the image and the texts.", " best\_text\_index": "the index of the most similar text.", "best\_text": "the most similar text."} examples = [ {"name": "GetImageToTextsSimilarity", "arguments": {"image": "image-0", "texts": ["a cat", "a dog"]}} super().\_\_init\_\_( id=id, description=description, args\_spec=args\_spec, rets\_spec=rets\_spec, examples=examples def \_\_call\_\_(self, image, texts, tool\_version='ViT-H-14-378-quickgelu'): 

```
1400 367
                import torch
1401 368
                import open_clip
1402 369
1403 370
                model, _, preprocess = open_clip.create_model_and_transforms(tool_version, pretrained='dfn5b')
1404 371
                model.eval() # model in train mode by default, impacts some models with BatchNorm or
1405
            stochastic depth active
1406 372
                tokenizer = open_clip.get_tokenizer(tool_version)
1407 373
1408 374
                image = preprocess(image_processing(image)).unsqueeze(0)
1409 375
                text = tokenizer(texts)
1410 376
1411 377
                with torch.no_grad(), torch.cuda.amp.autocast():
1412 378
                    image_features = model.encode_image(image)
1413 379
                    text_features = model.encode_text(text)
1414 380
                    image_features /= image_features.norm(dim=-1, keepdim=True)
1415 381
                    text_features /= text_features.norm(dim=-1, keepdim=True)
1416 382
1417 383
                    probs = image_features @ text_features.T
1418 384
                sim_scores = [round(sim_score, 2) for sim_score in probs[0].tolist()]
1419 385
                best_text_match = torch.argmax(probs).item()
1420 386
                return {'similarity': sim_scores, "best_text_index": best_text_match, "best_text": texts[
1421
            best text match]}
1422 387
1423 388
1424 389 class GetTextToImagesSimilarity(BaseAction):
1425 390
            def __init__(self, id) -> None:
1426 391
                description = "Get the similarity between one text and a list of images. Note that this
1427
            similarity score may not be accurate and should be used as a reference only."
1428 392
                args_spec = {
1429 393
                     "text": "the reference text.",
1430 394
                    "images": "a list of images to compare to the reference text.",
1431 395
1432 396
                rets_spec = {"similarity": "the CLIP similarity between the image and the texts.", "
1433
            best_image_index": "the index of the most similar image."}
1434 397
                examples = [
1435 398
                    {"name": "GetTextToImagesSimilarity", "arguments": {"text": "a black and white cat", "
1436
            images": ["image-0", "image-1"]}}
1437 399
                ]
1438 400
1439 401
                super().__init__(
1440 402
                    id=id, description=description, args_spec=args_spec, rets_spec=rets_spec, examples=examples
1441 403
                )
1442 404
            def __call__(self, text, images, tool_version='ViT-H-14-378-quickgelu'):
1443 405
1444 406
                import torch
1445 407
                import open_clip
1446 408
                original_images = images
1447 409
                model, _, preprocess = open_clip.create_model_and_transforms(tool_version, pretrained='dfn5b')
1448 410
                model.eval() # model in train mode by default, impacts some models with BatchNorm or
1449
            stochastic depth active
1450 411
                tokenizer = open_clip.get_tokenizer(tool_version)
1451 412
1452 413
                text = tokenizer([text])
1453 414
                images = [image_processing(image) for image in images]
1454 415
                images = torch.stack([preprocess(image) for image in images])
1455 416
1456 417
                with torch.no_grad(), torch.cuda.amp.autocast():
1457 418
                    image_features = model.encode_image(images)
1458 419
                    text_features = model.encode_text(text)
1459 420
                    image_features /= image_features.norm(dim=-1, keepdim=True)
1460 421
                    text_features /= text_features.norm(dim=-1, keepdim=True)
1461 422
1462 423
                    probs = text_features @ image_features.T
1463 424
                sim_scores = [round(sim_score, 2) for sim_score in probs[0].tolist()]
1464 425
                best_image_match = torch.argmax(probs).item()
1465 426
                return {'similarity': sim_scores, "best_image_index": best_image_match, "best_image":
1466
            original_images[best_image_match] }
```

class EstimateObjectDepth(BaseAction): def \_\_init\_\_(self, id) -> None: description = "Estimate the depth of an object in an image using DepthAnything model. It returns an estimated depth value of the object specified by the a brief text description. The smaller the value is, the closer the object is to the camera, and the larger the farther. This tool may help you to better reason about the spatial relationship, like which object is closer to the camera." args\_spec = { "image": "the image to get the depth from.", "object": "a short description of the object to get the depth from.", } rets\_spec = {"depth": "the estimated depth of the object."} examples = [ {"name": "EstimateObjectDepth", "arguments": {"image": "image-0", "object": "a black cat" }}, super().\_\_init\_\_( id=id, description=description, args\_spec=args\_spec, rets\_spec=rets\_spec, examples=examples ) def \_\_call\_\_(self, image, object, mode="mean"): action = LocalizeObjects(0) results = action(image=image, objects=[object]) if len(results["regions"]) == 0: return {"depth": "Object not found."} else: # use the best match object's bbox best\_match = np.argmax([region["score"] for region in results["regions"]]) bbox = results["regions"][best match]["bbox"] depth\_estimator = EstimateRegionDepth(0) return depth\_estimator(image=image, bbox=bbox, mode=mode) class EstimateRegionDepth(BaseAction): def \_\_init\_\_(self, id) -> None: description = "Estimate the depth of a region in an image using DepthAnything model. It returns an estimated depth value of the region specified by the input bounding box. The smaller the value is, the closer the region is to the camera, and the larger the farther. This tool may help you to better reason about the spatial relationship, like which object is closer to the camera. args spec = { "image": "the image to get the depth from.", "bbox": "the bbox of the region to get the depth from. It should be a list of [left, top, right, bottom], where each value is a float between 0 and 1 to represent the percentage of the image width/height and how far it is from the top left corner at [0, 0].", # "mode": "the mode to get the depth. It should be one of 'center' or 'average'. 'center' returns the depth of the center of the region. 'average' returns the average depth of the region.", } rets\_spec = {"depth": "the estimated depth of the region."} examples = [ {"name": "EstimateRegionDepth", "arguments": {"image": "image-0", "bbox": [0.3, 0.2, 0.5, 0.4]}}, ] super().\_\_init\_\_( id=id, description=description, args\_spec=args\_spec, rets\_spec=rets\_spec, examples=examples ) def \_\_call\_\_(self, image, bbox: List[str], mode="mean"): import numpy as np from scipy import stats image = image\_processing(image) depth\_model = pipeline(task="depth-estimation", model="depth-anything/Depth-Anything-V2-Small-hf", device=self.device) result = depth\_model(image) depth = result["predicted\_depth"][0].numpy() 

```
1534 481
                depth = depth.max() - depth # smaller values in depth map are farther from the camera so
1535
             reversing the values
1536 482
                H, W = depth.shape
1537 483
1538 484
                use_percent = all(x <= 1.0 for x in bbox)</pre>
1539 485
                if not use percent:
1540 486
                    raise ValueError("Bounding box coordinates must be between 0 and 1.")
1541 487
                bbox = [bbox[0] * W, bbox[1] * H, bbox[2] * W, bbox[3] * H]
1542 488
                if mode == "center":
1543 489
                    x, y = (bbox[0] + bbox[2]) / 2, (bbox[1] + bbox[3]) / 2
1544 490
                    x, y = int(x), int(y)
1545 491
                    depth_value = depth[y, x]
1546 492
                elif mode == "mean":
1547 493
                    x1, y1, x2, y2 = bbox
1548 494
                    x1, y1, x2, y2 = int(x1), int(y1), int(x2), int(y2)
1549 495
                    depth_value = np.mean(depth[y1:y2, x1:x2])
1550 496
                elif mode == "mode":
1551 497
                    x1, y1, x2, y2 = bbox
1552 498
                    x1, y1, x2, y2 = int(x1), int(y1), int(x2), int(y2)
1553 499
                    mode_result = stats.mode(depth[y1:y2, x1:x2])
1554 500
                    depth_value = mode_result.mode[0]
1555 501
                else:
1556 502
                    raise NotImplementedError(f"Depth mode {mode} is not supported.")
1557 503
                return {"depth": round(depth_value, 2)}
1558 504
1559 505
1560 506 class Calculate(BaseAction):
1561 507
            def __init__(self, id) -> None:
1562 508
                description = "Calculate a math expression."
1563 509
                args_spec = {"expression": "the math expression to calculate."}
                rets_spec = {"result": "the result of the math expression."}
1564 510
1565 511
                examples = [
1566 512
                     {"name": "Calculate", "arguments": {"expression": "2 + 2"}},
1567 513
                     {"name": "Calculate", "arguments": {"expression": "4*9*84"}},
1568 514
                     {"name": "Calculate", "arguments": {"expression": "5-4/2"}},
1569 515
                1
1570 516
1571 517
                super().__init__(
1572 518
                    id=id, description=description, args_spec=args_spec, rets_spec=rets_spec, examples=examples
1573 519
                )
1574 520
1575 521
           def __call__(self, expression):
1576 522
                result = eval(expression)
1577 523
                return {"result": result}
1578 524
1579 525
1580 526 class SolveMathEquation(BaseAction):
1581 527
            def __init__(self, id) -> None:
1582 528
                description = "Using this action to solve a math problem with WolframAlpha."
1583 529
                args_spec = {"query": "a question that involves a math equation to be solved."}
                rets_spec = {"result": "the result of the query."}
1584 530
1585 531
                examples = [
1586 532
                    {"name": "SolveMathEquation", "arguments": {"query": "2 + 2=?"}},
1587
                     {"name": "SolveMathEquation", "arguments": {"query": "x^2 + 2x + 1 = 0, what is x?"}},
     533
1588 534
                1
1589 535
1590 536
                self.client = wolframalpha.Client(os.getenv("WOLFRAM_ALPHA_API_KEY"))
1591 537
                super().__init__(
1592 538
                     id=id, description=description, args_spec=args_spec, rets_spec=rets_spec, examples=examples
1593 539
                )
1594 540
1595 541
            def __call__(self, query):
1596 542
                from urllib.error import HTTPError
1597
     543
1598 544
                is_success = False
1599 545
1600 546
                res = self.client.query(query)
```

```
1601
547
           if not res["@success"]:
                                                                                                                        1602
548
                                                                                                                        1603
                return (
549
550
                    "Your Wolfram query is invalid. Please try a new query for wolfram.",
                                                                                                                        1604
                                                                                                                        1605
551
                    is_success,
               )
                                                                                                                        1606
552
                                                                                                                        1607
553
           assumption = next(res.pods).text
           answer = ""
                                                                                                                        1608
554
           for result in res["pod"]:
                                                                                                                        1609
555
               if result["@title"] == "Solution":
                                                                                                                        1610
556
557
                    answer = result["subpod"]["plaintext"]
                                                                                                                        1611
                if result["@title"] == "Results" or result["@title"] == "Solutions":
                                                                                                                        1612
558
                    for i, sub in enumerate(result["subpod"]):
                                                                                                                        1613
559
                        answer += f"ans {i}: " + sub["plaintext"] + "\n"
                                                                                                                        1614
560
                                                                                                                        1615
561
                    break
           if answer == "":
                                                                                                                        1616
562
                                                                                                                        1617
               answer = next(res.results).text
563
                                                                                                                        1618
564
           if answer is None or answer == "":
                                                                                                                        1619
565
               return {"result": "No good Wolfram Alpha Result was found"}
                                                                                                                        1620
566
567
           else:
                                                                                                                        1621
                                                                                                                        1622
568
                return {"result": answer}
                                                                                                                        1623
569
                                                                                                                        1624
570
571 class DetectFaces (BaseAction):
                                                                                                                        1625
       def __init__(self, id) -> None:
                                                                                                                        1626
572
           description = "Using this function to detect faces in an image."
                                                                                                                        1627
573
           args_spec = {"image": "the image to detect faces from."}
                                                                                                                        1628
574
           rets_spec = {"image": "the image with objects localized and visualized on it.", "regions": "the
                                                                                                                        1629
575
                                                                                                                        1630
        regions of the faces detected, where each regin is represented by a dictionary with the region's
       label text and bounding box."}
                                                                                                                        1631
           examples = [
                                                                                                                        1632
576
                {"name": "DetectFaces", "arguments": {"image": "image-0"}}
                                                                                                                        1633
577
                                                                                                                        1634
578
           1
                                                                                                                        1635
           import face_detection
579
           ckpt_path = f"/root/.cache/torch/hub/checkpoints/WIDERFace_DSFD_RES152.pth"
                                                                                                                        1636
580
                                                                                                                        1637
           if not os.path.exists(ckpt_path):
581
582
                from huggingface_hub import hf_hub_download
                                                                                                                        1638
                hf_hub_download(repo_id="user/mma", filename="WIDERFace_DSFD_RES152.pth", local_dir="/root
                                                                                                                        1639
583
       /.cache/torch/hub/checkpoints/")
                                                                                                                        1640
                                                                                                                        1641
584
                                                                                                                        1642
           self.model = face_detection.build_detector(
585
                "DSFDDetector", confidence_threshold=.5, nms_iou_threshold=.3)
                                                                                                                        1643
586
           super().__init__(
                                                                                                                        1644
587
                                                                                                                        1645
588
                id=id, description=description, args_spec=args_spec, rets_spec=rets_spec, examples=examples
                                                                                                                        1646
589
           )
                                                                                                                        1647
590
                                                                                                                        1648
       def enlarge_face(self,box,W,H,f=1.5):
591
           x1, y1, x2, y2 = box
                                                                                                                        1649
592
           w = int((f-1) * (x2-x1)/2)
                                                                                                                        1650
593
           h = int((f-1) * (y2-y1)/2)
                                                                                                                        1651
594
           x1 = max(0, x1-w)
                                                                                                                        1652
595
                                                                                                                        1653
596
           y1 = max(0, y1-h)
                                                                                                                        1654
597
           x^2 = \min(W, x^2 + w)
                                                                                                                        1655
598
           y_{2} = \min(H, y_{2}+h)
                                                                                                                        1656
           return [x1,y1,x2,y2]
599
                                                                                                                        1657
600
                                                                                                                        1658
601
       def __call__(self, image):
                                                                                                                        1659
602
            import numpy as np
                                                                                                                        1660
           image = image_processing(image)
603
                                                                                                                        1661
604
           with torch.no_grad():
                                                                                                                        1662
605
                                                                                                                        1663
               faces = self.model.detect(np.array(image))
606
                                                                                                                        1664
607
608
           W,H = image.size
                                                                                                                        1665
                                                                                                                        1666
           objs = []
609
           for i,box in enumerate(faces):
                                                                                                                        1667
610
```

```
1668 611
                     x1, y1, x2, y2, c = [int(v) for v in box.tolist()]
1669 612
                     normalized_bbox = [x1/W, y1/H, x2/W, y2/H]
1670 613
                     objs.append(dict(
1671 614
                         bbox=[round(num, 2) for num in normalized_bbox],
1672 615
                         label=f'face {i+1}' if i > 0 else 'face',
1673 616
                    ))
1674 617
                visualize = VisualizeRegionsOnImage(0)
1675 618
                results = visualize(image=image, regions=objs)
1676 619
                tagged_image = results["image"]
1677 620
                results_formatted = {"regions": objs, "image": tagged_image}
1678 621
                return results_formatted
1679 622
1680 623
1681 624 class QueryLanguageModel (BaseAction):
1682 625
            def __init__(self, id) -> None:
1683 626
                description = "Using this function to ask a language model a question."
1684 627
                args_spec = {"query": "the question to ask the language model."}
1685 628
                rets_spec = {"result": "the response from the language model."}
1686 629
                examples = [
1687 630
                     {"name": "QueryLanguageModel", "arguments": {"query": "What is the capital of France?"}},
1688 631
1689 632
                super().__init__(
1690 633
                     id=id, description=description, args_spec=args_spec, rets_spec=rets_spec, examples=examples
1691 634
                )
1692 635
1693 636
            def __call__(self, query):
                 from openai import OpenAI
1694 637
1695 638
                client = OpenAI(api_key=os.getenv("OPENAI_API_KEY"))
1696 639
1697 640
                response = client.chat.completions.create(
1698 641
                    model=LATEST_GPT_MODEL_ID,
1699 642
                    messages=[
1700 643
                        {
1701 644
                             "role" : "user",
1702 645
                             "content": [
1703 646
                                 {"type": "text", "text": f"{query}"},
1704 647
                             1,
1705 648
                         }
1706 649
                    ],
1707 650
                    max_tokens=300,
1708 651
                )
1709 652
1710 653
                return {'result': response.choices[0].message.content}
1711 654
1712 655
1713 656 class QueryKnowledgeBase(BaseAction):
1714 657
           def __init__(self, id) -> None:
1715 658
                description = "Using this function to query a knowledge base."
1716 659
                args_spec = {"query": "the query to search in a knowledge base such as wikipedia."}
1717 660
                rets_spec = {"result": "the answer from the knowledge base."}
1718 661
                examples = [
1719 662
                     {"name": "QueryKnowledgeBase", "arguments": {"query": "Paris"}},
1720 663
                1
1721 664
1722 665
                super().__init__(
1723 666
                    id=id, description=description, args_spec=args_spec, rets_spec=rets_spec, examples=examples
1724 667
1725 668
1726 669
            def __call__(self, query, lang="en", sentences=2, knowledge_base="wikipedia"):
1727 670
                 if knowledge_base == "wikipedia":
1728 671
                     # Set the language for Wikipedia (default is 'en' for English)
1729 672
                     wikipedia.set_lang(lang)
1730 673
1731 674
                     # Search Wikipedia for pages related to the query
1732 675
                     search_results = wikipedia.search(query)
1733 676
                    if not search_results:
1734 677
                       return {"No results found."}
```

678		1735
679	# Get the summary of the first search result	1736
680	<pre>page = wikipedia.page(search_results[0])</pre>	1737
681	<pre>summary = wikipedia.summary(page.title, sentences=sentences)</pre>	1738
682	result = {	1739
683	"title": page.title,	1740
684	"url": page.url,	1741
685	"summary": summary	1742
686	}	1743
687	return result	1744
688	else:	1745
689	<pre>raise NotImplementedError(f"Knowledge base {knowledge_base} is not supported.")</pre>	1746
690		1747
691		1748
692	class Terminate(BaseAction):	1749
693	<pre>definit(self, id) -&gt; None:</pre>	1750
694	description = "Using this function to finish the task."	1751
695	<pre>args_spec = {"answer": "the final answer."}</pre>	1752
696	<pre>rets_spec = {"answer": "the final answer."}</pre>	1753
697	<pre>examples = [{"name": "Terminate", "arguments": {"answer": "yes"}}]</pre>	1754
698		1755
699	<pre>super()init(</pre>	1756
700	<pre>id=id, description=description, args_spec=args_spec, rets_spec=rets_spec, examples=examples</pre>	1757
701		1758
702		1759
703	<pre>defcall(self, answer):</pre>	1760
704	return {"answer": answer}	1761

Listing 1. Python implementation of all actions

1762 BEGIN OF GOAL 1763 2 You are a helpful assistant, and your goal is to solve the # USER REQUEST #. You can either rely on 1764 your own capabilities or perform actions with external tools to help you. A list of all available 1765 actions are provided to you in the below. 1766 3 [END OF GOAL] 1767 4 1768 5 [BEGIN OF ACTIONS] 1769 6 Name: OCR 1770 7 Description: Extract texts from an image or return an empty string if no text is in the image. Note 1771 that the texts extracted may be incorrect or in the wrong order. It should be used as a reference 1772 only. 1773 8 Arguments: {'image': 'the image to extract texts from.'} 1774 9 Returns: {'text': 'the texts extracted from the image.'} 1775 10 Examples: 1776 II {"name": "OCR", "arguments": {"image": "image-0"}} 1777 12 1778 13 Name: LocalizeObjects 1779 14 Description: Localize one or multiple objects/regions with bounding boxes. This tool may output objects 1780 that don't exist or miss objects that do. You should use the output only as weak evidence for 1781 reference. When answering questions about the image, you should double-check the detected objects. 1782 You should be especially cautious about the total number of regions detected, which can be more or 1783 less than the actual number. 1784 15 Arguments: {'image': 'the image to localize objects/regions in.', 'objects': "a list of object names to 1785 localize. e.g. ['dog', 'cat', 'person']. the model might not be able to detect rare objects or 1786 objects with complex descriptionriptions." } 1787 16 Returns: {'image': 'the image with objects localized and visualized on it.', 'regions': "the regions of 1788 interests localized in the image, where each region is represented by a dictionary with the region 1789 's label text, bounding box and confidence score. The confidence score is between 0 and 1, where 1 1790 means the model is very confident. Note that both the bounding boxes and confidence scores can be 1791 unreliable and should only be used as reference."} 1792 17 Examples: 1793 18 {"name": "LocalizeObjects", "arguments": {"image": "image-0", "objects": ["dog", "cat"]}} **1794** 19 1795 20 Name: GetObjects 1796 21 Description: Using this function to get objects in an image. 1797 22 Arguments: {'image': 'the image to get objects from.'} 1798 23 Returns: {'objects': 'the objects detected in the image.'} 1799 24 Examples: 1800 25 {"name": "GetObjects", "arguments": {"image": "image-0"}} 1801 26 1802 27 Name: EstimateRegionDepth 1803 28 Description: Estimate the depth of a region in an image using DepthAnything model. It returns an 1804 estimated depth value of the region specified by the input bounding box. The smaller the value is, 1805 the closer the region is to the camera, and the larger the farther. This tool may help you to 1806 better reason about the spatial relationship, like which object is closer to the camera. 1807 29 Arguments: {'image': 'the image to get the depth from.', 'bbox': 'the bbox of the region to get the 1808 depth from. It should be a list of [left, top, right, bottom], where each value is a float between 1809 0 and 1 to represent the percentage of the image width/height and how far it is from the top left 1810 corner at [0, 0].'} 1811 30 Returns: {'depth': 'the estimated depth of the region.' } 1812 31 Examples: 1813 32 {"name": "EstimateRegionDepth", "arguments": {"image": "image-0", "bbox": [0.3, 0.2, 0.5, 0.4]}} **1814** 33 1815 34 Name: EstimateObjectDepth 1816 35 Description: Estimate the depth of an object in an image using DepthAnything model. It returns an 1817 estimated depth value of the object specified by the a brief text description. The smaller the 1818 value is, the closer the object is to the camera, and the larger the farther. This tool may help 1819 you to better reason about the spatial relationship, like which object is closer to the camera. 1820 36 Arguments: {'image': 'the image to get the depth from.', 'object': 'a short description of the object 1821 to get the depth from.' } 1822 37 Returns: {'depth': 'the estimated depth of the object.' } 1823 38 Examples: 1824 39 {"name": "EstimateObjectDepth", "arguments": {"image": "image-0", "object": "a black cat"}} 1825 40 1826 41 Name: Crop 1827 42 Description: Crop an image with the bounding box. It labels the cropped region with a bounding box and 1828 crops the region with some margins around the bounding box to help with contextual understanding of

the region. 1829 43 Arguments: {'image': 'the image to crop.', 'bbox': 'the bbox to crop. It should be a list of [left, top 1830 , right, bottom], where each value is a float between 0 and 1 to represent the percentage of the 1831 image width/height and how far it is from the top left corner at [0, 0].' } 1832 1833 44 Returns: {'image': 'the cropped image.'} 45 Examples: 1834 46 {"name": "Crop", "arguments": {"image": "image-0", "bbox": [0.33, 0.21, 0.58, 0.46]}} 1835 1836 47 48 Name: ZoomIn 1837 1838 49 Description: Zoom in on a region of the input image. This tool first crops the specified region from the image with the bounding box and then resizes the cropped region to create the zoom effect. It 1839 1840 also adds some margins around the cropped region to help with contextual understanding of the 1841 region. so Arguments: {'image': 'the image to zoom in on.', 'bbox': 'The bbox should be a list of [left, top, 1842 right, bottom], where each value is a float between 0 and 1 to represent the percentage of the 1843 image width/height and how far it is from the top left corner at [0, 0].', 'zoom\_factor': 'the 1844 1845 factor to zoom in by. It should be greater than 1.' } 51 Returns: {'image': 'the zoomed in image.'} 1846 52 Examples: 1847 s3 {"name": "ZoomIn", "arguments": {"image": "image-0", "bbox": [0.4, 0.3, 0.5, 0.4], "zoom\_factor": 2}} 1848 54 1849 1850 55 Name: QueryLanguageModel 56 Description: Using this function to ask a language model a question. 1851 57 Arguments: {'query': 'the question to ask the language model.'} 1852 58 Returns: {'result': 'the response from the language model.' } 1853 59 Examples: 1854 60 {"name": "QueryLanguageModel", "arguments": {"query": "What is the capital of France?"}} 1855 1856 61 1857 62 Name: GetImageToImagesSimilarity 1858 63 Description: Get the similarity between one image and a list of other images. Note that this similarity 1859 score may not be accurate and should be used as a reference only. 64 Arguments: {'image': 'the reference image.', 'other\_images': 'the other images to compare to the 1860 reference image.'} 1861 1862 65 Returns: {'similarity': 'the CLIP similarity scores between the reference image and the other images.', 'best\_image\_index': 'the index of the most similar image.'} 1863 1864 66 Examples: {"name": "GetImageToImagesSimilarity", "arguments": {"image": "image-0", "other images": ["image-1", " 1865 67 image-2"]}} 1866 1867 68 69 Name: GetImageToTextsSimilarity 1868 1869 70 Description: Get the similarity between one image and a list of texts. Note that this similarity score 1870 may not be accurate and should be used as a reference only. n Arguments: {'image': 'the reference image.', 'texts': 'a list of texts to compare to the reference 1871 image.'} 1872 72 Returns: {'similarity': 'the CLIP similarity between the image and the texts.', 'best\_text\_index': 'the 1873 index of the most similar text.', 'best\_text': 'the most similar text.'} 1874 1875 73 Examples: 74 {"name": "GetImageToTextsSimilarity", "arguments": {"image": "image-0", "texts": ["a cat", "a dog"]}} 1876 1877 75 76 Name: GetTextToImagesSimilarity 1878 1879 Description: Get the similarity between one text and a list of images. Note that this similarity score 77 1880 may not be accurate and should be used as a reference only. 1881 78 Arguments: {'text': 'the reference text.', 'images': 'a list of images to compare to the reference text 1882 .'} 1883 79 Returns: {'similarity': 'the CLIP similarity between the image and the texts.', 'best\_image\_index': ' 1884 the index of the most similar image.' } 1885 80 Examples: 1886 81 {"name": "GetTextToImagesSimilarity", "arguments": {"text": "a black and white cat", "images": ["image -0", "image-1"]}} 1887 1888 82 83 Name: DetectFaces 1889 84 Description: Using this function to detect faces in an image. 1890 1891 85 Arguments: {'image': 'the image to detect faces from.'} 86 Returns: {'image': 'the image with objects localized and visualized on it.', 'regions': "the regions of 1892 the faces detected, where each regin is represented by a dictionary with the region's label text 1893 and bounding box." } 1894 87 Examples: 1895

```
1896 88 {"name": "DetectFaces", "arguments": {"image": "image-0"}}
1897
1898 90 Name: QueryKnowledgeBase
1899 91 Description: Using this function to query a knowledge base.
1900 92 Arguments: {'query': 'the query to search in a knowledge base such as wikipedia.'}
1901
     93 Returns: {'result': 'the answer from the knowledge base.'}
1902
     94 Examples:
1903 95 {"name": "QueryKnowledgeBase", "arguments": {"query": "Paris"}}
1904 96
1905 97 Name: Calculate
1906
     98 Description: Calculate a math expression.
1907 99 Arguments: {'expression': 'the math expression to calculate.'}
1908 100 Returns: { 'result': 'the result of the math expression.' }
1909 IOI Examples:
1910 102 {"name": "Calculate", "arguments": {"expression": "2 + 2"}}
1911 103 {"name": "Calculate", "arguments": {"expression": "4*9*84"}}
1912 104 {"name": "Calculate", "arguments": {"expression": "5-4/2"}}
1913 105
1914 Name: SolveMathEquation
1915 107 Description: Using this action to solve a math problem with WolframAlpha.
1916 ION Arguments: {'query': 'a question that involves a math equation to be solved.'}
1917 109 Returns: {'result': 'the result of the query.'}
1918 110 Examples:
1919 III {"name": "SolveMathEquation", "arguments": {"query": "2 + 2=?"}}
1920 112 {"name": "SolveMathEquation", "arguments": {"guery": x^2 + 2x + 1 = 0, what is x?"}}
1921 113
1922 114 Name: Terminate
1923 US Description: Using this function to finish the task.
1924 II6 Arguments: {'answer': 'the final answer.'}
1925 117 Returns: {'answer': 'the final answer.' }
1926 II8 Examples:
1927 119 {"name": "Terminate", "arguments": {"answer": "yes"}}
1928 120
1929 121 [END OF ACTIONS]
1930 122
1931 123 [BEGIN OF TASK INSTRUCTIONS]
1932 124 1. You must only select actions from # ACTIONS #.
1933 125 2. You can only call one action at a time.
1934 126 3. If no action is needed, please make actions an empty list (i.e. ''actions'': []).
1935 127 4. You must always call Terminate with your final answer at the end.
1936 128 [END OF TASK INSTRUCTIONS]
1937 129
1938 130 [BEGIN OF FORMAT INSTRUCTIONS]
1939 I31 Your output should be in a strict JSON format as follows:
1940 132 {"thought": "the thought process, or an empty string", "actions": [{"name": "action1", "arguments": {"
1941
            argument1": "value1", "argument2": "value2"}}]
1942 133 [END OF FORMAT INSTRUCTIONS]
1943 134
1944 135 [BEGIN OF EXAMPLES]:
1945 136 # USER REQUEST #:
1946 137 In image-0, Which of the two objects on the plate is the biggest?
1947 138 A. The pile of scrambled eggs is the biggest.
1948 139 B. The strawberries are the biggest object.
1949 140 Please answer directly with only the letter of the correct option and nothing else.
1950 141 # RESPONSE #:
1951 42 {"thought": "To determine which of the two objects on the plate is larger, I need to analyze the size
1952
            of the scrambled eggs, and the strawberries", "actions": [{"name": "LocalizeObjects", "arguments":
1953
             {"image": "image-0", "objects": ["scrambled eggs", "strawberries"]}}]
1954 143 OBSERVATION:
1955 144 {"image": "image-1", "regions": [{"label": "eggs", "bbox": [0.5, 0.6, 0.6, 0.8], "score": 0.85}, {"
1956
            label": "strawberries", "bbox": [0.4, 0.5, 0.45, 0.7], "score": 0.54}]}
1957 445 {"thought": "To calculate the area of a bounding box, we can use the formula: area = (x_max - x_min)
1958
            (y_max - y_min). We first get the area of the scrambled eggs.", "actions": [{"name": "Calculate", "
1959
            arguments": {"expression": "(0.6-0.5) * (0.8-0.6)"}}]}
1960 146 OBSERVATION:
1961 147 {"result": "0.02"}
1962 148 {"thought": "Then, we also calculate the area of the strawberries.", "actions": [{"name": "Calculate",
```

	"arguments": {"expression": "(0.45-0.4) * (0.7-0.5)"}}]}	1963
149	OBSERVATION:	1964
150	{"result": "0.01"}	1965
151	$\{$ "thought": "Since 0.02 > 0.01, it is apparent that the eqgs cover a larger area within their bounding	1966
	box.", "actions": [{"name": "Terminate", "arguments": {"answer": "A"}}]}	1967
152		1968
153	# USER REQUEST #:	1969
154	Given the input image image-0, How many pedestrians are there in the image? Please answer directly	1970
	with a single word or number.	1971
155	# RESPONSE #:	1972
156	{"thought": "To determine the number of pedestrians, I need to first localize them on the image.", "	1973
	actions": [{"name": "LocalizeObjects", "arguments": {"image": "image-0", "objects": ["person"]}}]}	1974
157	OBSERVATION:	1975
158	{"image": "image-1", "regions": [{"label": "person", "bbox": [0.77, 0.47, 0.79, 0.54], "score": 0.83},	1976
	{"label": "person-2", "bbox": [0.69, 0.49, 0.7, 0.52], "score": 0.43}]}	1977
159	{"thought": "The LocalizeObjects action returns two regions for \"person\", but one of the regions has	1978
	a lower confidence score. Upon a closer look at the output image image-1, we can see that there is	1979
	actually only one pedestrian in the image.", "actions": [{"name": "Terminate", "arguments": {"	1980
	answer": "1"}}]}	1981
160		1982
161	# USER REQUEST #:	1983
162	Based on $image-0$ , is the object on top bigger than the object below?	1984
163	A. The object on the bottom is bigger.	1985
164	B. The object on top is bigger.	<b>1986</b>
165	C. Both objects are the same size.	1987
166	Please answer directly with only the letter of the correct option and nothing else.	<b>1988</b>
167	# RESPONSE #:	1989
168	{"thought": "By looking at the image, we can see that both objects are game consoles of the same brand	1990
	and size.", "actions": [{"name": "Terminate", "arguments": {"answer": "C"}}]}	1991
169		1992
170	# USER REQUEST #:	1993
171	What is x in the image?	1994
172	# RESPONSE #:	1995
173	{"thought": "To get the result of the equation, I need to first extract the equation from the image.",	1996
	"actions": [{"name": "OCR", "arguments": {"image": "image-0"}}]}	1997
	OBSERVATION:	1998
	{"text": "x-2^3=0"}	1999
176	{"thought": "The math equation is $'x-2^3=0'$ , and I need to find x. I can solve it with a math-related	2000
	tool.", "actions": [{"name": "SolveMathEquation", "arguments": {"query": "x-2^3=0, what is x?"}}]}	2001
	OBSERVATION:	2002
	{"result": "x = 8"}	2003
179	<pre>{"thought": "As suggested in the last observation, the answer is 8.", "actions": [{"name": "Terminate",</pre>	2004
	"arguments": {"answer": "8"}}]}	2005
180		2006
181	[END OF EXAMPLES]	2007
	Listing 2 Model have distance with a surface and and	

### Listing 2. Model-based data generation system prompt

```
2008
      def GetObjects template():
2009
            thought_templates = ["I need to check what objects are present in the {image_kw}.",
      2
2010
                                 "I need to analyze the {image kw} for context."
      3
2011
      4
                                 "I need to identify the objects in the {image_kw}.",
2012
                                 "To answer the question, let's first analyze the {image_kw}.",
      5
2013
                                 "To answer the question, analyzing the objects in the {image_kw} is necessary."
      6
2014
            ]
2015
      7
            return thought_templates
2016
2017
      9 def LocalizeObjects_template():
2018
     10
            thought_templates = ["I need to analyze the positions of {objects} in the {image_kw}.",
2019
                                 "I need to analyze the locations of {objects} in the {image_kw}.",
     11
2020
                                 "I need to localize the {objects} based on the {image_kw}.",
     12
2021
                                 "I'll identify the positions of {objects} in the {image_kw}.",
     13
2022
                                 "I need to determine the positions of {objects} by analyzing the {image_kw}."]
     14
2023
     15
            return thought_templates
2024
     16
2025
     17 def EstimateObjectDepth_template():
2026
            thought_templates = ["I should estimate the depth of {object} to determine whether it is closer or
     18
2027
            farther.",
2028
     19
                                 "I will estimate the depth of {object}.",
                                  "I need to estimate the depth for {object} to make a comparison.",
2029
     20
2030
     21
                                  "To determine how far {object} is, I need to evaluate the distance to it.",
2031
                                  "I now need to estimate the depth for {object}."]
     22
2032
     23
            return thought_templates
2033
     24
2034
     25
2035
     26 def EstimateRegionDepth_template():
2036
            thought_templates = ["I should estimate the objects' depths to determine which one is closer.",
     27
2037
                                  "I need to estimate the region's depth in the image.",
     28
2038
                                  "I need to determine the depths of the detected objects based on their
     29
2039
            positions.".
2040
                                 "I need to estimate the depth of the objects to make a comparison.",
     30
2041
                                 "To determine the relative proximity of the objects in the image, I need to
     31
2042
            estimate the depth of each object."]
2043
     32
            return thought_templates
2044
     33
2045
     34 def Terminate_template():
2046
            thought_templates = ["Based on the information above, I can conclude that the answer is {answer}",
     35
2047
     36
                                  "Based on a close analysis of the {image_kw} and additional information above,
2048
             I believe the answer is {answer}.",
2049
                                  "I have analyzed the {image_kw} and the information above, and I believe the
     37
2050
            answer is {answer}.",
2051
                                 "The {image_kw} and the information above suggest that the answer is {answer}.
     38
2052
            ",
2053
                                 "According to the content of the {image_kw} and the observations, I can
     39
2054
            conclude that the answer is {answer}."]
2055
           return thought_templates
      40
```

Listing 3. Thought templates for each action

1 Compare the ground truth and prediction from AI models, to give a correctness score for the prediction. 2056 <AND> in the ground truth means it is totally right only when all elements in the ground truth are 2057 2058 present in the prediction, and <OR> means it is totally right when any one element in the ground truth is present in the prediction. The correctness score is 0.0 (totally wrong), 0.1, 0.2, 0.3, 2059 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, or 1.0 (totally right). Just complete the last space of the 2060 2061 correctness score. 2062 2 Question | Ground truth | Prediction | Correctness 3 ---- | ---- | ----2063 4 What is x in the equation?  $| -1 \langle AND \rangle -5 | x = 3 | 0.0$ 2064 5 What is x in the equation?  $| -1 \langle AND \rangle -5 | x = -1 | 0.5$ 2065 6 What is x in the equation?  $| -1 \langle AND \rangle -5 | x = -5 | 0.5$ 2066 7 What is x in the equation?  $| -1 \langle AND \rangle -5 | x = -5 \text{ or } 5 | 0.5$ 2067 2068  $_8$  What is x in the equation?  $\mid$  -1 <AND> -5  $\mid$  x = -1 or x = -5  $\mid$  1.0 9 Can you explain this meme? | This meme is poking fun at the fact that the names of the countries 2069 2070 Iceland and Greenland are misleading. Despite its name, Iceland is known for its beautiful green landscapes, while Greenland is mostly covered in ice and snow. The meme is saying that the person 2071 has trust issues because the names of these countries do not accurately represent their landscapes. 2072 2073 | The meme talks about Iceland and Greenland. It's pointing out that despite their names, Iceland 2074 is not very icy and Greenland isn't very green. | 0.4 2075 10 Can you explain this meme? | This meme is poking fun at the fact that the names of the countries Iceland and Greenland are misleading. Despite its name, Iceland is known for its beautiful green 2076 landscapes, while Greenland is mostly covered in ice and snow. The meme is saying that the person 2077 2078 has trust issues because the names of these countries do not accurately represent their landscapes. | The meme is using humor to point out the misleading nature of Iceland's and Greenland's names. 2079 2080 n Iceland, despite its name, has lush green landscapes while Greenland is mostly covered in ice and snow. The text 'This is why I have trust issues' is a playful way to suggest that these contradictions 2081 2082 can lead to distrust or confusion. The humor in this meme is derived from the unexpected contrast 2083 between the names of the countries and their actual physical characteristics. | 1.0 Listing 4. LLM judge prompt for MMVet 1 Please read the following example. Then extract the answer from the model response and type it at the 2084 2085 end of the prompt. 2086 3 Hint: Please answer the question requiring an integer answer and provide the final value, e.g., 1, 2, 2087 2088 3. at the end. 4 Question: Which number is missing? 2089 2090 5 Model response: The number missing in the sequence is 14. 6 Extracted answer: 14 2091 2092 2093 8 Hint: Please answer the question requiring a floating-point number with one decimal place and provide 2094 the final value, e.g., 1.2, 1.3, 1.4, at the end. 2095 9 Question: What is the fraction of females facing the camera? 10 Model response: The fraction of females facing the camera is 0.6, 2096 II which means that six out of ten females in the group are facing the camera. 2097 2098 12 Extracted answer: 0.6 2099 14 Hint: Please answer the question requiring a floating-point number with two decimal places and provide 2100 the final value, e.g., 1.23, 1.34, 1.45, at the end. 2101 2102 15 Question: How much money does Luca need to buy a sour apple candy and a butter-scotch candy? (Unit: \$) 16 Model response: Luca needs \$1.45 to buy a sour apple candy and a butterscotch candy. 2103 2104 17 Extracted answer: 1.45 2105 18 2106 19 Hint: Please answer the question requiring a Python list as an answer and provide the final list, e.g., 2107 [1, 2, 3], [1.2, 1.3, 1.4], at the end. 20 Question: Between which two years does the line graph saw its maximum peak? 2108 2109 21 Model response: The line graph saw its maximum peak between 2007 and 2008. 2110 22 Extracted answer: [2007, 2008] 23 2111 2112 24 Hint: Please answer the question and provide the correct option letter, e.g., A, B, C, D, at the end. 2113 25 Question: What fraction of the shape is blue? 2114 26 Choices: (A) 3/11 (B) 8/11 (C) 6/11 (D) 3/5 27 Model response: The correct answer is (B) 8/11. 2115 28 Extracted answer: B 2116

Listing 5. LLM judge prompt for MathVista