Self-Enhancing Programming-Driven Reasoning for Visual Question Answering

Anonymous ACL submission

Abstract

In Visual Question Answering (VQA) tasks, program-driven reasoning methods have made significant progress by transforming solutions into executable code. However, existing approaches often fall short due to their reliance on a single code generation iteration, lacking the flexibility to handle unforeseen errors. To address this limitation, we propose the Self-Enhancing Programming-driven Reasoning framework for VQA (Seper). Seper combines 011 a code generator and evaluator to decompose questions into multistep instructions, dynamically generating Python code based on input. The code evaluator performs both forward and backward evaluations, triggering an iterative code regeneration process that enables continu-018 ous optimization. Additionally, we introduce 019 prompt tuning to improve the quality of the generated code. Our experiments on the GQA and OK-VQA datasets demonstrate Seper's superior performance, highlighting its potential to advance VOA programming methods. Code: https://anonymous.4open.science/r/Seper-5540/

1 Introduction

027

034

042

The pursuit of artificial general intelligence (AGI) has led to the development of large-scale visual models, which often rely on task-specific datasets, such as object grounding, visual question answering, or image segmentation. These models struggle with complex, real-world tasks that require multistep reasoning. For example, answering "What was the color of the first horse that crossed the finish line?" involves identifying the finish line, recognizing the first horse, and determining its color. Such intricate tasks challenge the effectiveness and scalability of end-to-end models.

Early approaches, such as Modular Visual Question Answering(Andreas et al., 2016), decompose tasks into atomic units but require extensive supervision, limiting their adaptability across domains. The emergence of large language models has driven interest in automatic module integration, giving rise to program-driven reasoning methods (Schick et al., 2023; Suróis et al., 2023; Subramanian et al., 2023; Gupta and Kembhavi, 2022), which tackle complex problems through step-by-step code solutions. By addressing subproblems via defined APIs, these methods eliminate the training costs of traditional modular approaches and provide enhanced interpretability. 043

045

047

051

057

061

062

063

064

065

067

068

069

070

071

072

074

075

076

078

079

081

Existing programming-driven reasoning methods often lack iterative and reflective processes, typically performing only a single iteration of code generation and execution. These methods fail to address unforeseen issues, such as syntax errors or erroneous API results, which can affect the final outputs. While human programmers can debug code manually, this process is time-consuming and lacks real-time efficiency. Additionally, large language models (LLMs) usually require specific prompts to generate code, yet they struggle with complex visual language tasks due to difficulties in comprehending detailed task requirements, often resulting in suboptimal code generation.

Self-iterative models have emerged as a powerful approach in AI, particularly for code generation tasks that demand high precision. These models rely on unit tests to detect errors and use compiler outputs to inform corrections. By iteratively refining the generated code based on test results, they can effectively handle complex programming challenges, ensuring both correctness and optimization.

When applying the self-iterative framework to visual programming, several key challenges must be addressed during its implementation:

- Without unit tests, how can we effectively detect errors in the generated code?
- How can we revise the generated code based on the detected errors?

To address these limitations, we propose the Self-Enhancing Programming-Driven Reasoning (Seper) framework for Visual Question Answering (VQA), as illustrated in Figure 2. The Seper framework goes beyond traditional visual programming approaches by integrating multiple code evaluation components that automatically assess the generated code and produce detailed reports for further improvement.

083

087

091

097

100

101

102

103

104

105

106

107

108

109

110

111

112

129

To address these limitations, we propose the **S**elf-enhancing programming-Driven Reasoning (Seper) framework for Visual Question Answering (VQA), as illustrated in Figure 2. The Seper framework goes beyond traditional visual programming approaches by integrating multiple code evaluation components that automatically assess the generated code and produce detailed reports for further improvement.

Our main contributions include:

- We introduce the first iterative framework for Visual Question Answering that integrates both code generation and evaluation, marking a significant advancement in the field.
- We propose a supervised prompt tuning method for automatically adjusting prompt templates used by the code generator and reviewers during code evaluation.
- Experimental results on real benchmarks demonstrate the superior performance of our approach.

2 Related Work

The early method for the Visual Question Answer-113 ing is the Modular Visual Question Answering 114 (Modular VQA), which refers to a framework used 115 in Visual Question Answering (VQA) tasks that de-116 composes the question-answering process into dis-117 tinct, modular components (Andreas et al., 2016). 118 For example, Promptcap(Hu et al., 2023) is the first 119 to leverage the text reasoning capabilities of large language models in conjunction with generated 121 image captions to address visual questions. Sim-122 ilarly, Prophet(Shao et al., 2023) employs a large 123 language model to aggregate judgments from mul-124 tiple VQA models, improving the overall accuracy 125 of answers by combining insights from different 126 systems. 127

> In contrast to the methods described above, recent research has introduced **Visual Programming**

frameworks that uses the incontext learning abil-130 ity of large language models to generate python-131 like modular programs, which are then executed 132 to get both the solution and a comprehensive and 133 interpretable rationale. Each line of the generated 134 program may invoke one of several off-the-shelf 135 computer vision models, image processing subrou-136 tines, or python functions to produce intermediate 137 outputs that may be consumed by subsequent parts 138 of the program. This allows for the seamless inte-139 gration of new components, enhancing the system's 140 versatility and adaptability to a broad range of VQA 141 tasks. Several studies have explored visual pro-142 gramming for VQA (Suróis et al., 2023; Subrama-143 nian et al., 2023; Gupta and Kembhavi, 2022; Khan 144 et al., 2024; Dou et al., 2024), but these approaches 145 typically generate code in a single pass, without 146 iterative refinement. In contrast, our method con-147 tinuously evaluates and revises the generated code, 148 enabling ongoing optimization through repeated 149 iterations, which ensures improved performance 150 and adaptability. 151

152

153

154

155

156

157

158

159

160

161

162

163

165

166

167

168

169

170

171

172

173

174

175

176

177

178

However, the code generated by visual programming may not always run as expected. Even if execution is possible, issues such as incorrect interface calls or parameter mismatches can arise, leading to inaccurate results. Code refinement refers to the process of improving, optimizing, or cleaning up code while preserving its overall functionality. Statistical and learning-based techniques for code refinement have a long history in both programming languages and machine learning, typically focused on repairing human-written code (Long and Rinard, 2016; Bader et al., 2019; Le Goues et al., 2021). Self-Edit (Zhang et al., 2023) explores self-repair with natural language feedback on apps. Similarly, T-eval (Chen et al.) assesses Codex's ability to self-repair across a variety of tasks.

However, the aforementioned work is specifically designed for text-based tasks and relies on correction through unit testing (Olausson et al., 2023), which requires knowing the correct output in advance—an approach that is difficult to apply to Visual Question Answering (VQA) due to the challenge of labeling correct outputs. To address this issue, our framework incorporates a self-enhancing mechanism that autonomously identifies problems in the generated code and revises it accordingly.

3 Error Analysis Of Existing Methods

179

181

186

190

191

192

194

195

196

199

200

201

205

207

210

212

213

214

215

216

217

218

219

To examine the influence of code generation quality and pre-trained model capabilities on programming-driven reasoning, we analyzed 200 randomly selected error cases from Viper (Suróis et al., 2023) using the GQA (Hudson and Manning, 2019) test-dev dataset. Figure 1 presents the evaluation results, categorizing the failure reasons into three main groups:

- 1. **Code Error:** These errors arise during code execution and include syntax issues (e.g., missing semicolons or mismatched parentheses) and runtime problems (e.g., invalid input, memory overflow, or division by zero), which can lead to program termination.
- 2. **API Limitation:** This refers to situations where API calls, particularly those involving pre-trained models, fail to produce expected results. This category includes: Image-Text-Contrastive (ITC) models, which use models like CLIP (Radford et al., 2021) and XVLm (Zeng et al., 2021) which assess image attributes and are affected by object names and properties; Visual Question Answering (VQA) models, where slight variations in question phrasing can impact outcomes despite identical semantics; and Location (LOC) models, such as GroundingDINO (Liu et al., 2023), which detect and crop objects in images, with object name sensitivity influencing detection accuracy.
 - 3. Logic Error: Can be classified into answer logic errors and potential logic errors.
 - **Answer Inaccuracies:** Return values misaligned with question format (e.g., yes/no for multiple-choice).
 - **Potential Logic Errors:** Reflect mismatches between the question-solving process and code logic, often arising with more complex images than anticipated, such as missing target objects, as illustrated in Figure 5.

The insights from Figure 1 reveal that Viper mainly faces API limitations and logic errors, with code errors being rare due to GPT's robustness. VQA models notably contribute to API limitations, highlighting the importance of question representation in visual question answering.



Figure 1: Error analysis results.

227

228

229

230

231

232

234

235

236

237

238

240

241

242

243

244

245

246

247

248

249

250

251

253

254 255

256

257

258

259 260

261

262 263

264

266

268

270

4 Methodology

The Seper framework combines code generation, execution, and evaluation. Given an image and question, the code generator breaks down the task, generates code, and passes it to the Python compiler. The code is evaluated in two stages: forward evaluation checks input, and backward evaluation traces errors. Language Model-based reviewers detect issues, and if errors are found, a review report is generated to refine the code iteratively until it passes or the iteration limit is reached.

4.1 Code Generation

As shown in Figure 2(a), the code generation process consists of two stages using a Code Large Language Model (CodeLLM): rationale generation and Python code generation. To enhance reasoning interpretability and soundness, CodeLLM first generates rationales that break down the task into subtasks, guiding the subsequent code generation. These subtasks are then embedded as comments in the code to improve readability.

A predefined prompt, as detailed in Listing 1, guides the CodeLLM in generating accurate Python code. This prompt combines the question \hat{q} , API specifications, examples, and a review report, providing a comprehensive input to steer the code generation process.

Listing 1. I fompt for code generativ	: Prompt for code ger	generatio	ation	tion
---------------------------------------	-----------------------	-----------	-------	------

Here are the APIs available for
calling
<pre>\${API specifications}</pre>
<pre># Learn from the following examples</pre>
<pre>\${Related examples}</pre>
The question
\${Question}
The review result from previous
iteration
<pre>\${Rationale}</pre>
Generate python code to answer the
question with variable "image" as
innut
Tubac

Prompt Composition: To answer a visual question using an image, the goal is to leverage a Large



Figure 2: Our proposed Seper framework.



Figure 3: Supervised Prompt Tuning framework.

Language Model to generate Python code by enhancing the question representation with three key elements:

271

272

273

275

284

- 1. API Specifications: We use the ImagePatch API specification and text query \hat{q} as input prompts for the program generation model (Suróis et al., 2023). While following the original implementation, our approach differs in categorizing the APIs into VQA, Location, and Other types. Detailed API specifications are in Appendix B.
- 2. **Related examples**: The process of example selection is crucial. We follow the setup of mainstream methods by using manually constructed examples as context in experiments.
- 3. **Rationale**: Including the three subtasks decomposed by the LLM, as well as the possible review report from the review process detailing code execution results and error

corrections, is vital for the code generator. CodeLLM uses this feedback to iteratively refine and regenerate code.

290

291

292

293

294

296

297

298

299

301

302

303

304

305

306

307

309

310

311

312

313

314

315

316

317

4.2 Code Execution

Upon code generation, we employ the Python compiler to execute the code and gather all necessary variables and intermediate data for thorough review. Ultimately, a comprehensive execution result report is generated, as illustrated in Figure 2(b).

Execute Code: Upon code generation, snippets will be injected to preserve variables, which will be stored in a Python dictionary for later analysis. Generate Execution Result Report: Multiple execution result reports will be generated in structured format, each covering five key components comprehensively: (1) Source code: The Python code generated by the code generator; (2) Error Message: Returns detailed error information (location, type, stack trace) for the first encountered error, or "No error" if none detected; (3) Return Value: The final text-based return value or an empty string in case of encountering errors; (4) Runtime variables: Variables are categorized as ImagePatch objects with BLIP-2 descriptions, basic strings, and collections with text-converted elements and lengths; (5) Logical judgement result: Extracted conditional statements trace execution flow through loops and branches, using transformed variables to

- 318
- 319 320

323

325

327

330

331

332

334

335

336

337

347

348

355

359

363

indicate decisions.

4.3 Code Evaluation

The code evaluation process, as shown in Figure 2(c), involves extracting the relevant execution results for each sub-task and conducting evaluations in two sub-stages: forward and backward.

- Forward Evaluation: Verifying whether each line of code ascertain its successful execution, followed by validating whether the output of each line conforms to the expected results.
- **Backward Evaluation:** Assessment of code output against intended objectives, followed by retrospective analysis of potential errors and intermediate results.

Unlike traditional code generation tasks (Long and Rinard, 2016; Bader et al., 2019; Le Goues et al., 2021), which can be validated using test cases, our approach faces the challenge of missing test cases, leading to uncertain and unverifiable outcomes. To address this, we utilize heuristic-based reviews focused on the three components outlined in the previous section: **code error checking** to identify syntax and runtime issues, **API checking** to verify API calls and their results, and **code logic checking** to detect logical errors causing incorrect outputs. Different Language Model-based reviewers evaluate each component, following a two-step process using the prompts detailed in Appendix C:

- **Result assessment:** This process identifies errors by analyzing execution result reports, examining explicit and implicit errors from various sources. Each report is reviewed by different reviewers who provide detailed responses and a true or false answer. A majority vote then determines the final acceptance of the response.
- **Review generation:** Reviewers compile assessment results to identify error types and underlying causes, offering specific suggestions for code improvement. The goal is to provide a clear and thorough analysis to guide code modifications and enhance quality.

Finally, we merge the reports of each sub-task into a comprehensive report and provide feedback to the code generator.

4.3.1 Forward Evaluation.

The forward evaluation stage verifies that each line of code executes correctly and produces the expected output. It focuses on local issues, such as runtime errors, syntax errors, and the proper functioning of API calls.

Code Error Reviewer

Code error reviewer assesses the basic errors in the code and generate reviews as follows:

Result assessment: Leveraging the error information and location captured during the code execution, we directly utilize the error message in the execution result report. If no error message is present, the review generation step is skipped.

Review generation: In the prompt for the Code Error Reviewer, we first provide definitions for syntax errors and runtime errors. Subsequently, we instruct the reviewer to address syntax errors directly and to add exception handling for runtime errors by generating new code for the erroneous part.

API Reviewer

The API reviewer's main objective is to systematically identify and evaluate errors in API calls, consolidating these results into a comprehensive review. Given that VQA-type and location-type APIs exhibit the most errors, the reviewer employs specialized strategies tailored to these API types. Result Assessment: Each API is treated as a black box without disclosing intermediate results. The output variable's name indicates the expected output type. For example, "clothing_kind" in line 9 suggests the answer should be a clothing type. Results are assessed by checking alignment with the expected type. This involves two stages: extracting key details from the execution report (e.g., objects, parameters, and output values) and guiding the reviewer to evaluate acceptability by prompt template. For vga APIs, reviewers assess the plausibility and clarity of the answer. For location-type APIs, reviewers confirm the presence of the object as per the image caption and its relevance to the original question.

Review Generation: Inspired by the Exploration-Exploitation Dilemma in reinforcement learning, we design two instruction sets for guiding the API reviewer in generating comments. The first encourages CodeLLM to adjust parameters of existing API calls, such as using more specific parameters if a "simple_query(.)" output is vague. The second promotes exploration by suggesting alternative 402

403

404

405

406

407

408

409

410

411

412

413

414

364



Figure 5: Example of lacking structure type logical error.

API calls. For example, as shown in Figure 4, when
vqa-type APIs alone are insufficient, Seper prompts
the code generator to combine them with loc-type
APIs.

4.3.2 Backward Evaluation.

In neural networks, backpropagation calculates final error and updates gradients in preceding layers. Similarly, we propose a backward evaluation mechanism to trace final errors up the source code, using a logic reviewer to identify and analyze error origins.

Logic Reviewer

419

420

421

422

423

494

425

426

497

428

As discussed in section 3, logic errors can be categorized into answer errors and potential errors.

Result Assessment: Unlike APIs, which are 429 treated as black boxes, a logic reviewer can ac-430 cess all code lines and runtime variables, enabling 431 a thorough evaluation of the code's logic. The as-432 433 sessment is conducted from two perspectives: 1) Answer Inaccuracies: These errors are similar to 434 those in Visual Question Answering (VQA) sys-435 tems, such as misinterpretations of queries like "Is 436 he to the right or left of the mirror?". Reviewers as-437

sess whether the response aligns with the intended question. 2) **Potential Errors:** These errors stem from issues in key algorithmic constructs, including sequential, branching, and looping structures. Such errors may arise from missing structures (e.g., absent loops) or inaccuracies within existing ones (e.g., incorrect branch conditions). Reviewers identify these issues by analyzing the variables and source code to ensure the logic is correctly implemented. 438

439

440

441

442

443

444

445

446

447

Review Generation: Reviewers analyze the code 448 to identify logic errors and provide actionable feed-449 back for improvement. The evaluation is conducted 450 from two perspectives: 1) Answer Inaccuracies: 451 For answer errors, reviewers trace the return value 452 back to the first incorrect statement to identify the 453 root cause, which may involve an API or logical 454 error. This process helps pinpoint the exact issue 455 and informs further analysis. 2) Potential Errors: 456 For potential errors, such as incorrect spatial prop-457 erties in left-right judgment (as shown in Figure 458 5), reviewers assess the logic structure as a whole. 459 They evaluate whether the algorithm's constructs 460 are correctly implemented and aligned with the 461 intended behavior, suggesting modifications if necessary. Finally, reviewers output "pass" if no errors
are found; otherwise, the aggregated report is sent
to CodeLLM for further iteration.

4.4 Supervised Prompt Tuning

466

467

468

469

470

471

472

473

474

475

476

477 478

479

480

481 482

483

484

485

486

487

488

489

490

491

492

493

494

495

497

498

499

502

505

509

To improve the performance of CodeLLM, we propose a supervised prompt tuning method to optimize prompt templates for code generation. As shown in Figure 3, the process begins by running Seper on a training dataset using initial prompt templates to generate and execute code. We construct a sample pool, initially populated with all training samples. After evaluating the results, we select representative examples, especially those that frequently cause errors, for prompt refinement. In an iterative process, error-prone samples are retained in the pool, while correct samples are gradually discarded, with only a small subset kept. These examples help the reviewer model identify failure patterns and refine the prompt templates. Details are provided in Appendix F.2.

5 Experiment

5.1 Experimental Setup

Benchmarks. We evaluate our method on the GQA (Hudson and Manning, 2019) and OK-VQA (Marino et al., 2019) datasets. GQA focuses on multi-hop questions using human-annotated scene graphs for complex compositional queries, while OK-VQA assesses answering image-related questions requiring external knowledge.
Implementation. We utilized ChatGPT(OpenAI,

2023) to serve as both code generator and three reviewers, each with distinct sessions to ensure independence. Choosing ChatGPT is to maintain consistency with the mainstream visual programming methods that have predominantly used ChatGPT, and also because, as previous research(Chen et al., 2023). We employ BLIP-2(Li et al., 2023), GroundingDINO(Liu et al., 2023), and XVLM(Zeng et al., 2021) as the pre-trained models underpinning our API. Detailed parameters are provided in the Appendix A.

5.2 Results

Table 1 presents results on the GQA test-dev and OK-VQA validation datasets, showing that Seper outperforms all other zero-shot methods. Despite utilizing large foundational models, programmingbased methods like ViperGPT and CodeVQA are Table 1: Result on the GQA and OK-VQA,*indicates an evaluation on a stratified sample of the test set, which may not be directly comparable to results on the full test set.

Model	GQA Acc.(%)	OK-VQA Acc.(%)
Supervised		
VisRep	59.2	-
LXMERT	60.0	-
NSM	63.0	-
CRF	72.1	-
KAT	-	54.4
RA-VQA	-	54.5
REVIVE	-	58.0
PromptCap	-	58.8
Zero-Shot		
BLIP-2	44.7	45.9
Flamingo	-	50.6
ViperGPT	48.1	51.9
CodeVQA	49.0	53.5
VisProg*	50.5	52.6
Seper(ours)	51.3	54.1

limited by their single-pass approach, resulting in lower accuracy. Seper, with its iterative framework for code generation and evaluation, achieves the highest accuracy of 51.3% on GQA and 54.1% on OK-VQA, outperforming the next best methods, VisProg, by 0.8% and 1.5%, respectively. This demonstrates the advantage of Seper's multi-stage review process in enhancing performance.

Figure 6 compares the accuracy of the baseline and Seper across various question types in GQA. With the integration of the Code Evaluation module and prompt tuning, Seper outperforms the baseline in several categories. Notably, the accuracy for "positionChoose" questions increased dramatically from 0.37 to 0.95, highlighting the effectiveness of the Code Evaluation module in handling complex positional relationships, as demonstrated in line 7 of Figure 2(a). Although some question types saw a slight decrease in accuracy compared to the baseline, a detailed analysis of these variations can be found in Appendix F.3.

5.3 Ablations

Table 2 presents an ablation study on the GQA dataset, evaluating the impact of various components on model performance. The baseline model achieves 48.1% accuracy. Incorporating code error

510



Figure 6: Accuracy distribution among different question types in GQA.

Table 2: Ablation results on the GQA and OK-VQA. \checkmark indicates the component is included in the configuration.

Reviewers				
Code Error	API	Logic	P-Tuning	Acc(%)
Baseline				48.1 (+0.0)
Two-compo	ation			
\checkmark	\checkmark			48.7 (+0.6)
\checkmark		\checkmark		49.3 (+1.2)
\checkmark			\checkmark	48.8 (+0.7)
	\checkmark	\checkmark		49.1 (+1.0)
	\checkmark		\checkmark	48.6 (+0.5)
		\checkmark	\checkmark	49.8 (+1.7)
Three-component Configuration				
\checkmark	\checkmark	\checkmark		50.4 (+2.3)
\checkmark	\checkmark		\checkmark	49.5 (+1.4)
\checkmark		\checkmark	\checkmark	50.0 (+1.9)
	\checkmark	\checkmark	\checkmark	50.9 (+2.8)
Full Configu	ıratio	n		
\checkmark	\checkmark	\checkmark	\checkmark	51.3 (+3.2)

and API reviewers improves accuracy to 48.7%, while adding logic reviewers boosts it further to 49.3%. Combining code error and logic reviewers results in 49.8% accuracy. The configuration with all three components—code error, API, and logic reviewers—achieves 50.4%, and adding prompt tuning further increases accuracy to 50.9%. The full configuration, including all components and P-Tuning, reaches the highest accuracy of 51.3%, representing a 3.2% improvement over the baseline. These results underscore the significance of multi-level reviews and prompt tuning in enhancing VQA model performance, with each additional component contributing to more refined and accu-

536

537

538

539

541

543

545

546

549

rate outcomes.

Error Analysis. For the GQA dataset, we manually inspected 100 samples to analyze error sources. As shown in Figure 7, our method outperforms Viper by eliminating basic code errors, significantly reducing logical errors, and effectively identifying API-related issues in models.

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

567

568

569

570

571

572

573

574

575



Figure 7: Error analysis of Seper and Viper.

6 Conclusion

We introduced the Seper framework, an iterative approach to visual question answering that integrates multi-level reviews and innovative evaluation stages with distinct reviewers. We also proposed a supervised prompt template tuning method to automatically optimize prompts for code reviewers. Our approach significantly outperforms existing programming-driven reasoning methods, addressing API limitations and enhancing logic, while showcasing the potential of LLMs as effective code reviewers. However, the reliance on a closed-source ChatGPT for reviews has increased costs, and we did not fully address the API limitations or enhance it to resolve associated issues. The logic review process also lacks the insights necessary for generating improved solutions. Future work will focus on reducing costs by leveraging open-source language models.

7 Limitations

In this study, our methodology did not show a significant advantage over existing programmingdriven reasoning approaches. Additionally, we ac-579 knowledge that the limitations of the API were not fully addressed, nor was the API enhanced to ef-582 fectively tackle the underlying issues. As a result, the logic review was not robust enough to provide deeper insights, which are essential for generating more refined code solutions. Furthermore, the reliance on a closed-source ChatGPT for multiple 586 review iterations led to relatively high costs, a factor that should be considered in future research. 588

References

- Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Dan Klein. 2016. Neural module networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 39–48.
- Johannes Bader, Andrew Scott, Michael Pradel, and Satish Chandra. 2019. Getafix: Learning to fix bugs automatically. *Proceedings of the ACM on Programming Languages*, 3(OOPSLA):1–27.
- Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. Teaching large language models to self-debug.
- Zehui Chen, Weihua Du, Wenwei Zhang, Kuikun Liu, Jiangning Liu, Miao Zheng, Jingming Zhuo, Songyang Zhang, Dahua Lin, Kai Chen, et al. 2023. T-eval: Evaluating the tool utilization capability step by step. *arXiv preprint arXiv:2312.14033*.
- Zi-Yi Dou, Cheng-Fu Yang, Xueqing Wu, Kai-Wei Chang, and Nanyun Peng. 2024. Re-rest: Reflectionreinforced self-training for language agents. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 15394– 15411.
- Tanmay Gupta and Aniruddha Kembhavi. 2022. Visual programming: Compositional visual reasoning without training.
- Yushi Hu, Hang Hua, Zhengyuan Yang, Weijia Shi, Noah A Smith, and Jiebo Luo. 2023. Promptcap: Prompt-guided image captioning for vqa with gpt-3. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2963–2975.
- Drew A Hudson and Christopher D Manning. 2019. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6700–6709.
- Zaid Khan, Vijay Kumar BG, Samuel Schulter, Yun Fu, and Manmohan Chandraker. 2024. Self-training

large language models for improved visual program synthesis with visual reinforcement. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 14344–14353.

- Claire Le Goues, Michael Pradel, Abhik Roychoudhury, and Satish Chandra. 2021. Automatic program repair. *IEEE Software*, 38(4):22–27.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. Blip-2: Bootstrapping language-image pretraining with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*.
- Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. 2023. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*.
- Fan Long and Martin Rinard. 2016. Automatic patch generation by learning correct code. In *Proceedings* of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, pages 298–312.
- Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. 2019. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/cvf conference on computer vision and pattern recognition*, pages 3195–3204.
- Theo X Olausson, Jeevana Priya Inala, Chenglong Wang, Jianfeng Gao, and Armando Solar-Lezama. 2023. Is self-repair a silver bullet for code generation. arXiv preprint arXiv:2306.09896.
- OpenAI. 2023. Chatgpt (march 14 version) [large language model]. https://chat.openai.com.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools.
- Zhenwei Shao, Zhou Yu, Meng Wang, and Jun Yu. 2023. Prompting large language models with answer heuristics for knowledge-based visual question answering. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 14974– 14983.
- Sanjay Subramanian, Medhini Narasimhan, Kushal Khangaonkar, Kevin Yang, Arsha Nagrani, Cordelia Schmid, Andy Zeng, Trevor Darrell, and Dan Klein. 2023. Modular visual question answering via code generation.

670

671

672

673

674

675

676

677

678

679

680

627

628

629

630

631

632

576

590

591

593

598

599

605

610

611

612

613

614

615

616

619

623

- 681 Dóidac Suróis, Sachit Menon, and Carl Vondrick. 2023.
 682 Vipergpt: Visual inference via python execution for reasoning.
- Yan Zeng, Xinsong Zhang, and Hang Li. 2021.
 Multi-grained vision language pre-training: Aligning texts with visual concepts. *arXiv preprint arXiv:2111.08276*.
- Kechi Zhang, Zhuo Li, Jia Li, Ge Li, and Zhi Jin. 2023.
 Self-edit: Fault-aware code editor for code generation. *arXiv preprint arXiv:2305.04087*.

Α **Implemetation Details**

We set the temperature parameter of codellm to 0, the temperature parameter of the reviewer to 0.25, and the rest of the parameters remain as usual. $\frac{1}{22}$ Besides, we capped the maximum number of itera-²³ tions at 3 to prevent infinite function iterations. For each dataset, we employed 16 in-context examples. 25 The pre-trained model and parameter settings we use are shown below:

19

20

34

36

39

44

45

46

47

48

49

50

image.

upper : int

GroundingDINO. We use the implementation from the official GitHub repository. The thresh-²⁶ 701 old of box is 0.30 and the threshold of text is 0.25. $_{27}$ 702 **XVLM**. We used the official implementation, 703 specifically the version finetuned for retrieval on MSCOCO. In addition to this, we set the similarity 28 705 threshold to 0.5. Meaning if the similarity between ²⁹ 706 the image and the text is more than 0.5, the two are considered to match. 30

> **BLIP-2**. We used the Flan-T5 XXL version from ³¹ Huggingface. We set the beams in the generation parameter to 5. 33

B **API Listings**

710

711

712

713

715

We provide different API definitions for different datasets.

B.1 API for GQA

```
38
          1 from PIL import Image
            from vision_functions import
718
                find_in_image, simple_qa,
719
                verify_property, best_text_match
721
           def bool_to_yesno(bool_answer: bool)->
          4
                str:
                                                         40
723
                return "yes" if bool_answer else "no
724
725
          6
726
          7
           class ImagePatch:
                """A Python class containing a crop
                                                         41
          8
                                                         42
                of an image centered around a
                                                         43
729
                particular object, as well as
730
                relevant information.
731
                Attributes
          9
732
         10
733
                cropped_image : array_like
734
                    An array-like of the cropped
735
                image taken from the original image.
736
                left : int
         13
737
                    An int describing the position
         14
                of the left border of the crop's
                bounding box in the original image.
                lower : int
         15
741
                    An int describing the position
         16
742
                of the bottom border of the crop's
743
                bounding box in the original image.
744
         17
                right : int
                    An int describing the position
745
         18
746
                of the right border of the crop's
747
                bounding box in the original image.
```

upper : int An int describing the position the top border of the crop's of bounding box in the original image. Methods find(object_name: str)->List[ImagePatch] Returns a list of new ImagePatch objects containing crops of the image centered around any objects found in the image matching the object_name. simple_query(question: str=None)-> str Returns the answer to a basic question asked about the image. If no question is provided, returns the answer to "What is this?" exists(object_name: str)->bool Returns True if the object specified by object_name is found in the image, and False otherwise. verify_property(property: str)->bool Returns True if the property is met, and False otherwise. best_text_match(string1: str, string2: str)->str Returns the string that best matches the image. crop(left: int, lower: int, right: int, upper: int)->ImagePatch Returns a new ImagePatch object containing a crop of the image at the given coordinates. def __init__(self, image, left: int= None, lower: int=None, right: int= None, upper: int=None): """Initializes an ImagePatch object by cropping the image at the given coordinates and stores the coordinates as attributes. If no coordinates are provided, the image is left unmodified, and the coordinates are set to the dimensions of the image. Parameters image : array_like An array-like of the original image. left : int An int describing the position of the left border of the crop's bounding box in the original image. lower : int An int describing the position of the bottom border of the crop's bounding box in the original image. right : int An int describing the position of the right border of the crop's bounding box in the original

```
52 An int describing the
818
                                                           95
                position of the top border of the
819
820
                crop's bounding box in the original
821
                image.
822
          53
                                                           97
                     .....
823
          54
824
                     if left is None and right is
          55
                                                           98
825
                None and upper is None and lower is
826
                None:
                                                           99
827
                          self.cropped_image = image
          56
828
          57
                          self.left = 0
829
                          self.lower = 0
          58
                          self.right = image.shape[2]
          59
831
                 # width
                                                          100
832
                          self.upper = image.shape[1]
          60
833
                 # height
                                                          101
834
                     else:
          61
835
                         self.cropped_image = image
          62
836
                [:, lower:upper, left:right]
                                                          102
                          self.left = left
837
          63
                                                          103
                          self.upper = upper
838
          64
839
                          self.right = right
          65
                                                           104
840
                          self.lower = lower
          66
841
          67
842
                     self.width = self.cropped_image.
          68
843
                shape[2]
                                                          105
844
                     self.height = self.cropped_image
          69
845
                .shape[1]
846
          70
                                                           106
847
                     self.horizontal_center = (self.
          71
                                                          107
                left + self.right) / 2
                                                          108
                     self.vertical_center = (self.
          72
850
                lower + self.upper) / 2
                                                           109
851
          73
                                                          110
852
                 def find(self, object_name: str)->
          74
853
                List["ImagePatch"]:
                     """Returns a new ImagePatch
854
          75
                object containing the crop of the
855
                image centered around the object
                                                          112
                specified by object_name.
858
                     Parameters
          76
                                                          114
859
          77
                                                          115
860
                     object_name : str
          78
861
                         A string describing the name
          79
                 of the object to be found in the
                                                          116
863
                image.
                                                          117
                     .....
864
          80
                                                          118
                     return find_in_image(self.
                                                          119
          81
                cropped_image, object_name)
867
          82
                                                          120
                 def simple_query(self, question: str
          83
                =None)->str:
                     """Returns the answer to a basic
870
          84
                 question asked about the image. If
871
                                                          122
872
                no question is provided, returns the
                 answer to "What is this?".
873
                                                          123
874
          85
                     Parameters
875
          86
                     _____
                                                          124
876
          87
                     question : str
                        A string describing the
877
          88
                question to be asked.
878
                                                          126
879
          89
                     Examples
          90
                                                          127
881
          91
                     _____
                                                          128
882
          92
883
          93
                     >>> # Which kind of animal is
                not eating?
884
                                                          129
885
          94
                     >>> def execute_command(image)->
                str:
```

```
>>> image_patch = ImagePatch
(image)
            animal_patches =
   >>>
image_patch.find("animal")
    >>>
            for animal_patch in
animal_patches:
                if not animal_patch.
   >>>
verify_property("animal", "eating"):
    >>>
                   return
animal_patch.simple_query("What kind
of animal is eating?") # crop would
include eating so keep it in the
query
           # If no animal is not
   >>>
eating, query the image directly
   >>> return image_patch.
simple_query("Which kind of animal
is not eating?")
    >>> # What is in front of the
horse?
    >>> # contains a relation (
around, next to, on, near, on top of
 in front of, behind, etc), so ask
directlv
    >>> return image_patch.
simple_query("What is in front of
the horse?")
    >>>
    .....
    return simple_qa(self.
cropped_image, question)
def exists(self, object_name: str)->
bool:
    """Returns True if the object
specified by object_name is found in
the image, and False otherwise.
    Parameters
    object_name : str
        A string describing the name
of the object to be found in the
image.
    Examples
    _____
    >>> # Are there both cakes and
gummy bears in the photo?
    >>> def execute_command(image)->
str:
    >>>
            image_patch = ImagePatch
(image)
   >>>
            is_cake = image_patch.
exists("cake")
          is_gummy_bear =
   >>>
image_patch.exists("gummy bear")
            return bool_to_yesno(
    >>>
is_cake and is_gummy_bear)
    return len(self.find(object_name
)) > 0
def verify_property(self,
object_name: str, property: str)->
bool:
    """Returns True if the object
possesses the property, and False
otherwise.
```

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

956	130	Differs from 'exists' in that it	170	
957		presupposes the existence of the		
958		object specified by object_name,		
959		instead checking whether the object	171	
960		possesses the property.		
961	131	Parameters	172	
962	132		173	
963	133	object_name : str		
964	134	A string describing the name	174	
965		of the object to be found in the	175	
966		image.		
967	135	property : str		
968	136	A string describing the	176	
969		property to be checked.		
970	137	Examples	177	
971	138	Examples	178	
972	139		179	
973	140	solor?	180	
974	1.4.1	COIDT :	101	
076	141	<pre>ctr:</pre>	181	
970	142	>>> image patch = ImagePatch	162	
978	142	(image)	183	
979	142	<pre>>>> letters natches =</pre>	103	
980	145	image patch find("letters")	104	
981	144	>>> # Question assumes only	185	
982	144	one letter patch	186	
983	145	>>> if len(letters patches)	100	
984		== 0:	187	
985	146	>>> # If no letters are	188	
986		found, query the image directly	189	
987	147	>>> return image_patch.		
988		simple_query("Do the letters have		
989		blue color?")	190	
990	148	<pre>>>> return bool_to_yesno(</pre>	191	
991		<pre>letters_patches[0].verify_property("</pre>		
992		letters", "blue"))	192	
993	149			
994	150	<pre>return verify_property(self.</pre>		
995		<pre>cropped_image, object_name, property</pre>	193	
996)	194	
997	151		195	
998	152	<pre>def best_text_match(self,</pre>	196	
999		option_list: List[str]) -> str:	197	
1000	153	"""Returns the string that best		
1001		matches the image.	198	
1002	154	Parameters	199	
1003	155	option list , str	200	
1004	156	A list with the names of the	200	
1005	157	different options	201	
1000	1.50	nrefiv · str	202	
1008	150	A string with the prefives	202	
1000	1.59	to append to the options	203	
1010	160		204	
1011	161	Examples	205	
1012	162		206	
1013	163	>>> # Is the cap gold or white?	207	
1014	164	>>> def execute command(image)->	208	
1015	101	str:		
1016	165	>>> image_patch = ImagePatch		
1017		(image)	209	
1018	166	>>> cap_patches =	210	
1019		<pre>image_patch.find("cap")</pre>		
1020	167	>>> # Question assumes one		
1021		cap patch		
1022	168	>>> if len(cap_patches) ==		
1023		0:		
1024	169	>>> # If no cap is found		B.
1025		, query the image directly		

```
>>> return image_patch.
                                                1026
simple_query("Is the cap gold or
                                                1027
white?")
                                                1028
    >>>
             return cap_patches[0].
                                                1029
best_text_match(["gold", "white"])
"""
                                                1030
                                                1031
    return best_text_match(self.
                                                1032
cropped_image, option_list)
                                                1033
                                                1034
def crop(self, left: int, lower: int
, right: int, upper: int)->"
                                                1035
                                                1036
ImagePatch":
                                                1037
    """Returns a new ImagePatch
                                                1038
cropped from the current ImagePatch.
                                                1039
    Parameters
                                                1040
    _____
                                                1041
    left : int
                                                1042
        The leftmost pixel of the
                                                1043
cropped image.
                                                1044
                                                1045
    lower : int
       The lowest pixel of the
                                                1046
cropped image.
                                                1047
                                                1048
    right : int
        The rightmost pixel of the
                                                1049
cropped image.
                                                1050
    upper : int
                                                1051
                                                1052
        The uppermost pixel of the
cropped image.
                                                1053
                                                1054
    _____
    .....
                                                1055
    return ImagePatch(self.
                                                1056
cropped_image, left, lower, right,
                                                1057
upper)
                                                1058
                                                1059
def overlaps_with(self, left, lower,
                                                1060
 right, upper):
    """Returns True if a crop with
                                                1061
                                                1062
the given coordinates overlaps with
                                                1063
this one,
                                                1064
    else False.
                                                1065
    Parameters
                                                1066
                                                1067
                                                1068
    left : int
        the left border of the crop
                                                1069
to be checked
                                                1070
    lower : int
                                                1071
        the lower border of the crop
                                                1072
 to be checked
                                                1073
    right : int
                                                1074
        the right border of the crop
                                                1075
 to be checked
                                                1076
    upper : int
                                                1077
        the upper border of the crop
                                                1078
 to be checked
                                                1079
                                                1080
    Returns
                                                1081
                                                1082
    _____
    bool
                                                1083
        True if a crop with the
                                                1084
given coordinates overlaps with this
 one, else False
                                                1086
                                                1087
    return self.left <= right and</pre>
                                                1088
self.right >= left and self.lower <=</pre>
                                                1089
                                                1090
 upper and self.upper >= lower
        Listing 2: API for GQA.
```

B.2 API for OK-VQA

```
1092
           1 from PIL import Image
1093
           2 from vision_functions import
1094
                obtain_query_response_from_image
1095
             from nlp_functions import llm_query,
1096
                 process_guesses
1097
           4
1098
           5 def llm_query(question: str)->str:
1099
           6
                 '''Answers a text question using GPT
                 -3. The input question is always a
1100
1101
                 formatted string with a variable in
1102
                 it.
1103
           7
1104
                 Parameters
           8
1105
           9
1106
                 question: str
          10
1107
                     the text question to ask. Must
1108
                 not contain any reference to 'the
1109
                 image' or 'the photo', etc.
1110
                 , , ,
1111
          13
                 return llm_query(question)
1112
          14
1113
          15 def process_guesses(question: str,
                 guesses: List[str])->str:
1114
1115
                  '''Processes a list of guesses for
          16
                 an answer to a question and returns
1116
                 the best answer.'''
1117
1118
          17
                 return process_guesses(question,
1119
                 guesses)
1120
          18
1121
          19 class ImagePatch:
1122
                 def __init__(self, image, left: int=
          20
                 None, lower: int=None, right: int=
1123
1124
                 None, upper: int=None):
                     if left is None and right is
1125
          21
1126
                 None and upper is None and lower is
1127
                 None:
1128
                          self.cropped_image = image
          22
1129
                          self.left = 0
          23
                          self.lower = 0
1130
          24
1131
          25
                          self.right = image.shape[2]
1132
                          self.upper = image.shape[1]
          26
1133
          27
                      else:
1134
                          self.cropped_image = image
          28
1135
                 [:, lower:upper, left:right]
                          self.left = left
1136
          29
1137
          30
                          self.lower = lower
1138
                          self.right = right
          31
                          self.upper = upper
1139
          32
1140
          33
1141
                      self.width = self.cropped_image.
          34
1142
                 shape[2]
1143
                     self.height = self.cropped_image
          35
1144
                 .shape[1]
1145
          36
1146
          37
                 def simple_query(self, query: str):
                      """Answer basic queries about
1147
          38
1148
                 the image patch.
1149
          30
                      Parameters
1150
          40
1151
          41
1152
                      _____
          42
1153
                      query: str
          43
1154
                         the simple query about the
          44
1155
                 image patch in the form of a
1156
                 question
1157
          45
1158
                      Returns
          46
1159
          47
                      _ _ _ _ _ _
1160
          48
                      str
1161
                          a guess for the answer to
          49
```

```
the question
                                                      1162
                                                       1163
50
           answer =
                                                      1164
51
      obtain_query_response_from_image(
                                                      1165
                                                      1166
      self.cropped_image, query)
                                                      1167
52
           return answer
            Listing 3: API for OK-VQA.
      Prompt For Reviewer
  С
                                                      1168
  C.1 Prompt for vga type APIs
                                                      1169
  Prompt for result asscessment
                                                      1170
prompt = f''You are a reviewer, your
                                                      1171
      task is to work hand in hand with
                                                      1172
      another AI to solve a question about
                                                       1173
       a certain picture,
                                                      1174
 2 Your task is to check whether the answer
                                                       1175
       meets the following two criteria
                                                      1176
      and give a judge result with
                                                       1177
      reasoning process:
                                                      1178
3 1. Is it possible that this answer is the
                                                      1179
       correct answer? Output Sure,
                                                      1180
      possible, or impossible.
                                                      1181
4 2.Is it this answer clear enough?Output
                                                      1182
      yes or no.
                                                      1183
                                                      1184
                                                      1185
  Examples as follows:
                                                       1186
7
  expected_result: material_type
                                                      1187
8
9 object: airplane
                                                      1188
10 query: What is the material?
                                                      1189
in answer: plastic
                                                      1190
12 judge:
                                                      1191
  - Is it possible that this answer is the
                                                      1192
13
       correct answer?:Possible. Plastic
                                                      1193
      is not a common material for
                                                      1194
      airplanes, unless it is a toy
                                                       1195
                                                      1196
      airplane.
  - Is it this answer clear enough?: Yes. It
                                                      1197
14
       is clear enough.
                                                      1198
                                                      1199
15
16 expected_result:skateboard_kind
                                                      1200
17 object:skateboard
                                                      1201
18 query: What is the this?
                                                      1202
                                                      1203
19 answer:kid
20 judge:
                                                       1204
  - Is it possible that this answer is the
                                                      1205
21
       correct answer?: Impossible.Kid is
                                                      1206
      not a kind of skateboard.
                                                      1207
  - Is it this answer clear enough?: Yes. It
                                                       1208
22
       is clear enough.
                                                      1209
                                                      1210
23
24 expected_result: horse_kind
                                                      1211
25 object:horse.
                                                      1212
26 query:What kind of horse is it?
                                                      1213
27 answer:horse.
                                                      1214
28 judge:
                                                      1215
29 - Is it possible that this answer is the
                                                      1216
       correct answer?:Sure.
                                                      1217
30 - Is it this answer clear enough?:No.It
                                                      1218
      is vague.
                                                      1219
                                                      1220
31
32 expected_result: person
                                                      1221
33 object:clothing
                                                      1222
34 query:Who will love this suit?
                                                      1223
                                                      1224
```

35 answer:No one.

36 judge:

```
1226
          37 - Is it possible that this answer is the 21 Judge from last step: Kid is not a kind
1227
                  correct answer?: Impossible. The
1228
                 answer should be occupation, gender,
1229
                  and everything else that is used to
1230
                  describe people.
1231
             - Is it this answer clear enough?:No.It
          38
1232
                 is vague.
1233
          39
1234
          40
1235
          41 (Examples finished)
1236
          42
1237
          43 expected_result:{expected_result}
1238
          44 object:{object}
1239
          45 query:{query}
1240
          46 answer:{model_output}
1241
          47 judge:
1242
          48
             , , ,
```

Listing 4: Prompt for vqa type APIs reviewer on resualt asscessment phase.

Prompt for review generation

1243

```
1 \text{ prompt} = f''
1244
1245
           2 You are a reviewer. You are reviewing a
1246
                code execution flow, and there is an
1247
                 error in an API call statement.
1248
                 Another large language model has
1249
                 already identified the error and its
                 cause. Please provide suggestions
1250
1251
                for modification.
1252
1253
           4 Here are some suggestion for you to
1254
                provide:
1255
           5 1. Replace the query parameter of
1256
                 original function. Use a more
1257
                informative question.
1258
           6 2.Use other function to answer the
                original visual question.
1259
1260
            For example, if there are options given
                or implied in the question, you can
1261
1262
                use best_text_match with giving
                options.
1264
          8
1265
          9
1266
          10 Some of the api's provided to you may be
1267
                 limited by the model resulting in
1268
                less than expected results, these
1269
                api's include:
          ii find(object_name): different names for
1270
1271
                the same object may lead to
                different results.
1272
1273
          12 simple_query(question): for different
1274
                 questions with the same semantics,
1275
                may lead to different results.
1276
          13 verify_property(object_name, property):
1277
                for different object_name or
1278
                property with same semantics, may
1279
                lead to different results.
1280
          14 best_text_match(option_list): Returns
                 the string that best matches the
1282
                 image.
1283
          15
1284
          16
1285
          17 Examples as follows:
1286
          18 (Example 1)
1287
          19 API statement: skateboard_kind =
1288
                 skateboard.simple_query('What is
1289
                 this?')
          20 Actual Result: Kid.
1290
```

```
1291
      of skateboard.
                                                       1292
  Review: The original query parameter
                                                       1293
      What is this?' is too simple.
                                                       1294
23 The code should use a more informative
                                                       1295
                                                       1296
      information.
                                                       1297
24
25 (Example 2)
                                                       1298
                                                       1299
26 API statement: horse_kind = horse_patch.
      simple_query("What kind of horse is
                                                       1300
      it?")
                                                       1301
27 Actual Result: horse
                                                       1302
28 Judge from last step: The answer is too
                                                       1303
                                                       1304
      vague.
29 Review: The original query parameter '
                                                       1305
      What kind of horse is it?' is
                                                       1306
      specific enough. But the answer is
                                                       1307
                                                       1308
      still too vague. Try to use "
      best_text_match" and giving some
                                                       1309
      horse kind as options.
                                                       1310
                                                       1311
30
31 (Example finished)
                                                       1312
                                                       1313
32
                                                       1314
33
  API statement:{code}
34 Actual Result:{result}
                                                       1315
35 Judge from last step:{judge}
                                                       1316
                                                       1317
36 Review:
37
  , , ,
                                                       1318
```

Listing 5: Prompt for vqa type APIs reviewer on review generation phase.

```
C.2 Prompt for location type API
                                                      1319
  Prompt for result asscessment
                                                      1320
prompt = f''You are an expert in
                                                      1321
      textual reasoning, and you are good
                                                      1322
      at reasoning from textual
                                                      1323
      information to get answers to
                                                      1324
                                                      1325
      questions.
2 Now here's a image with an AI-generated
                                                      1326
      description of that image(not 100%
                                                      1327
      correct), and the query designed for
                                                      1328
       that image.
                                                      1329
<sup>3</sup> You need to speculate on the answer to
                                                      1330
      another question based on the
                                                      1331
      original query and the description.
                                                      1332
4
                                                      1333
                                                      1334
5
6 Here are some examples:
                                                      1335
  (Example1)
                                                      1336
7
  ImageCaption: The image depict an aux
                                                      1337
8
      and bar.
                                                      1338
  OriginalQuery: The color of the top of
                                                      1339
      the aux is green or blue?
                                                      1340
10 Question1: Does foo exists in the image?
                                                      1341
11 Answer: Although the existence of foo is
                                                      1342
       not mentioned in the caption, its
                                                      1343
      presence is implied in the query,
                                                      1344
      because if foo does not exist, the
                                                      1345
      question becomes unanswerable.Yes, it
                                                      1346
       exists.
                                                      1347
                                                      1348
                                                      1349
13
14 (Example2)
                                                      1350
15 ImageCaption: There is an bar in the
                                                      1351
      left of the foo.
                                                      1352
16 OriginalQuery: The color of the aux of
                                                      1353
```

foo is green or blue?

```
1355
          17 Question: Does aux exists in the image?
1356
          18 Answer: If aux does not exist, there
1357
                will be no answers to the query.
1358
                 Althought it does not mention in the
1359
                  caption, the aux exists.
1360
          19
1361
          20
1362
          21
1363
          22 (Example3)
1364
          23 ImageCaption: There are an foo and an
1365
                 aux in the image.
1366
          24 OriginalQuery: Can you see bar or aux in
1367
                  the image?
1368
          25 Question: Does bar exists in the image?
1369
          26 Answer: No, it does not exist. Based on
                 the image caption and there is no
1370
                 affirmation of its existence in the
1371
1372
                 query.So bar does not exist in the
1373
                 image.
1374
          27
1375
          28 (Example finished)
1376
          29 ImageCaption:{(caption)}
          30 OriginalQuery:{(query)}
1377
          31 Question: Based on the the query "{(
                                                    "{(
                 query)}" and image caption,Does
1379
                 object)}" exists in the image?
1380
1381
          32
1382
          33 Remeber:
          34 - Learn from previous examples.
1384
          35 - The information in the query more
1385
                 important than the information in
                 the caption, because caption is
1386
1387
                 generated by ai that not 100 percent
1388
                 correct.
          36 Output:
1389
1390
          37 {
1391
          38
                 'existence':'Sure' or 'Impossible',
            }
1392
          39
            , , ,
1393
          40
```

2

3

4

5

6

7

8

Listing 6: Prompt for location type APIs reviewer on resualt asscessment phase.

Prompt for review generation

1394

```
prompt = f'''
1395
1396
            You are a reviewer. You are reviewing a
          2
                code execution flow, and there is an
1397
1398
                 error in an API call statement.
1399
                Another large language model has
                already identified the error and its
1400
1401
                 cause. Please provide suggestions
1402
                for modification.
1403
          3
1404
          4 Here are some suggestion for you to
1405
                provide:
1406
           5 1. Replace the object parameter in
                function. Try use more other synonym
1407
1408
                 object name replace the original
1409
                object name.
1410
          6 2.Use other function as error handler
1411
                while you can not find tagret object
1412
1413
1414
          8 Some of the api's provided to you may be
1415
                 limited by the model resulting in
1416
                less than expected results, these
1417
                api's include:
          9 find(object_name): different names for
1418
1419
                the same object may lead to
1420
                different results.
```

10	<pre>simple_query(question): for different</pre>	1421
	questions with the same semantics,	1422
	may lead to different results.	1423
11	verify property(object name.property):	1424
	for different object name or	1425
	property with same semantics may	1426
	lead to different results	1427
12	hest text match(option list): Returns	1428
12	the string that best matches the	1/120
	imago	1/20
10	Image.	1/01
13	Providuo codo fail to find the (abject)	1401
14	which chevild he in the image based	1432
	which should be in the image based	1433
	on the {query}. If y to give some	1434
	suggestions.	1430
15	- 1 1 1	1436
16	Examples as follows:	1437
17	(Example1)	1438
18	Query: What do the vase and the paper	1439
	have in common?	1440
19	Object should be found: paper	1441
20	Review: Try use more other synonym	1442
	object name replace 'paper' in the	1443
	paper_patches = image_patch.find("	1444
	paper") for increase chance to find	1445
	paper in the image.	1446
21		1447
22	(Example2)	1448
23	Query: How are the vehicles to the right	1449
	of the walking person that is	1450
	walking on the sidewalk called?	1451
24	Object should be found: vehicles	1452
25	Review: Explore more parameters.Replace	1453
	the vehicle with 'car' or other	1454
	words of the same semantics.	1455
26		1456
27	(Example3)	1457
28	Ouerv: What is the color of the carpet?	1458
29	Object should be found: carpet	1459
30	Review: Add simple query function for	1460
	error handler.	1461
31		1462
32	(Examples finished)	1463
32	Ouerv · { ouerv }	1464
24	Object should be found (object)	1/65
25	Review	1/66
20		1400
10		

Listing 7: Prompt for location type APIs reviewer on review generation phase.

Prompt for result asscessment of answer error1469f'''1470Suppose you are a checker, and you are working on a visual question task with another AI, who will give his approximation147114701471147114711472147114731473
f''' Suppose you are a checker, and you are working on a visual question task with another AI, who will give his approved to the visual question 1470 1470 1470 1470 1470 1470 1470 1470 1470 1470 1470 1470 1470 1470 1470 1471 1470 1470 1471 1470 1471 1470 1471 1470 1471 1470 1471 1472 1473 1473 1473 1473 1473 1473 1473 1474 1473 1473 1473 1474 1473 1473 1474 1473 1474 1473 1474 1473 1473 1473 1474 1473 1474 1473 1474 1473 1474
Suppose you are a checker, and you are working on a visual question task with another AI, who will give his appeared to the visual question
working on a visual question task 1472 with another AI, who will give his 1473
with another AI, who will give his 1473
answer to the viewal question 1474
allswer to the visual question, 1474
and your task is to determine whether 1475
the answer conforms to the form of 1476
the question 1477
Output 'Accept' or 'Refuse' with your 1478
thinking process. 1479
Learn from examples. 1480
Examples as follows: 1481
(Example1) 1482
Visual Question:What kind of horse is it 1483
? 1484

```
1485
          9 AI Answer:horse.
1486
          10 Judge: Refuse. The answer is too vague.
1487
1488
          12 (Example2)
1489
          13 Visual Question: Who will love this suit?
1490
          14 AI Answer:No one.
          15 Judge: Refuse. The answer should be
1491
1492
                occupation, gender, and everything
1493
                else that is used to describe people
1494
1495
          16 (Example finished)
1496
          17
1497
          18 Remeber:
1498
          19
1499
          20 - The AI answer does not need to be
1500
                extremely precise, but rather within
                 a certain range. For example, the
1501
                answer to the question Where
                                                 is
1503
                here?
                        can be "street" or "parking
                 lot",
1504
                       rather than a specific
1505
                location.
1506
          21 - Do not give 'correct answer' of the
                question. Becuase you can not give
1507
                correct answer for an unseen image.
1510
          23 Visual Question:{{(question)}}
1511
          24 AI Answer:{{(answer)}}
1512
          25 Judge:
1513
          26
```

23

24

25

26

27

28

29

30

31

33

34

37

30

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

Listing 8: Prompt for logic reviewer on result asscessment phase of answer error.

Prompt for review generation of answer error

```
1 \text{ Prompt} = f''
1515
1516
           2 I want you to serve as a senior code
1517
                engineer. I will provide you with
1518
                 the code, including its intermediate
1519
                states and return values.
           3 You should be able to trace back from
                the return values to identify the
1522
                 root cause statements that lead to
                 these return values.
1524
           4
1525
           5 Examples as follows:
1526
           6 (Example1)
1527
           7 code:
1528
           8 def execute_command(image)->str:
1529
           9
                 image_patch = ImagePatch(image)
                 person_patches = image_patch.find("
          10
1531
                person")
                 # Question assumes only one person
1533
                patch
1534
                 if len(person_patches) == 0:
          12
1535
                     # If no person is found, query
          13
1536
                 the image directly
1537
                     return image_patch.simple_query
          14
1538
                 ("What does the person to the left
                of the helmet hold?")
1539
1540
          15
                 person_patch = person_patches[0]
1541
          16
                 helmet_patches = image_patch.find("
1542
                 helmet")
                 # Question assumes only one helmet
1544
                patch
1545
                 if len(helmet_patches) == 0:
          18
1546
                     return "no"
          19
                    helmet_patch in helmet_patches:
1547
          20
                 for
1548
                     if helmet_patch.
1549
                 horizontal_center > person_patch.
1550
                horizontal_center:
```

```
return helmet_patch.
                                                      1551
      simple_query("What does the person
                                                      1552
      to the left of the helmet hold?")
                                                      1553
      return "no"
                                                      1554
  intermediate_result:
                                                      1555
       'image_patch': 'a baseball game',
                                                      1556
       'person_patches': 'lens:4'
                                                      1557
       'len(person_patches)': '4',
                                                      1558
      'if len(person_patches) == 0?': '
                                                      1559
      False',
                                                      1560
       'person_patch': 'a baseball player
                                                      1561
      holding a bat'
                                                      1562
       'helmet_patches': 'lens:3',
                                                      1563
       'len(helmet_patches)': '3',
                                                      1564
      'if len(helmet_patches) == 0?': '
                                                      1565
      False',
                                                      1566
      'helmet_patch.horizontal_center':
                                                      1567
      '333.5',
                                                      1568
       'person_patch.horizontal_center':
                                                      1569
      '230.0',
                                                      1570
      'if helmet_patch.horizontal_center >
                                                      1571
       person_patch.horizontal_center?':
                                                      1572
                                                      1573
      True'
36 return value: right
                                                      1574
                                                      1575
38 Return value traceability analysis The
                                                      1576
      root line is the "helmet_patch.
                                                      1577
      simple_query("What does the person
                                                      1578
      to the left of the helmet hold?")'
                                                      1579
      line. The reason is because
                                                      1580
      horizontal_center of one of the
                                                      1581
      helmet_patch bigger than that of
                                                      1582
                                                      1583
      person_patch.
                                                      1584
                                                      1585
40 (Example2)
41 # Which color is the device the monitor
                                                      1586
      is to the right of?
                                                      1587
  def execute_command(image)->str:
                                                      1588
      image_patch = ImagePatch(image)
                                                      1589
      monitor_patches = image_patch.find("
                                                      1590
                                                      1591
      monitor")
      # Question assumes only one monitor
                                                      1592
                                                      1593
      patch
      if len(monitor_patches) == 0:
                                                      1594
           # If no monitor is found, query
                                                      1595
      the image directly
                                                      1596
                                                      1597
           return image_patch.simple_query
      ("Which color is the device the
                                                      1598
      monitor is to the right of?")
                                                      1600
      for monitor_patch in monitor_patches
                                                      1601
           device_patches = image_patch.
                                                      1602
      find("device")
                                                      1603
                                                      1604
           for device_patch in
      device_patches:
              if device_patch.
                                                      1606
      horizontal_center < monitor_patch.</pre>
                                                      1607
      horizontal_cent er:
                   return device_patch.
                                                      1609
      simple_query("Which color is the
                                                      1610
      device?")
                                                      1611
      return "unknown"
                                                      1612
                                                      1613
                                                      1614
  intermediate_result:
       "image_patch": "a computer monitor
                                                      1615
      and a laptop computer on a desk",
                                                      1616
         "monitor_patches":"lens:2"
                                                      1617
         "len(monitor_patches)":"2",
                                                      1618
         "if len(monitor_patches) == 0?":"
                                                      1619
      False",
                                                      1620
```

```
1621
                   "device_patches":"lens:4",
          61
                   "device_patch.horizontal_center
1622
          62
1623
                 ":"285.0",
1624
                   "monitor_patch.horizontal_center":
          63
1625
                  "319.0",
                   "if device_patch.horizontal_center
1626
          64
                  < monitor_patch.horizontal_center
                 ?": "True"
1628
1629
          65 return value: right
1630
          66
1631
          67 Return value traceability analysis: The
                 root line is the "return
1632
1633
                 device_patch.simple_query("Which
1634
                 color is the device?")" line.The
1635
                 reason is
1636
            because horizontal_center of one of the
1637
                 device_patch smaller than that of
                 monitor_patch.
1639
            (Examples finished)
          69
1640
          70
1641
          71 The first code is:{code}
1642
          72 The intermediate result:{inter_res}
1643
          73 return_value:{return_value}
          74
            Return value traceability analysis:
          75
```

Listing 9: Prompt for logic reviewer on review generation phase of answer error.

20

21

23

24

25

26

27

28

29

30

31

32

33

34

37

38

39

42

43

44

45

47

48

49

50

51

52

53

54

55

56

57

Prompt for result asscessment and review generation phase of potential error

1646

```
1 \text{ Prompt} = f''
1648
1649
           2 I want you to serve as a senior code
1650
                 engineer. I will provide you with
1651
                 the code, including its intermediate
1652
                states and return values.
1653
           3 First.Summarize the content of the
                intermediate variables. The term '
1655
                 image_patch' represents the content
1656
                of the image; the number of patches
1657
                 indicates how many of a certain
                object are present in the image.
1659
           4 Second.Analyse from guestions code and
1660
                 intermediate results. Check if there
                 are errors in the loops and
1661
1662
                 selection structures in this code,
1663
                or if there is a lack of loops and
1664
                 selection structures.
           5
1666
           6 Examples as follows:
           7 (Example1)
            code:
1668
           8
1669
           9
            def execute_command(image)->str:
                 image_patch = ImagePatch(image)
1670
          10
1671
                 person_patches = image_patch.find("
          11
1672
                 person")
1673
          12
                 # Question assumes only one person
1674
                patch
1675
          13
                 if len(person_patches) == 0:
          14
                     # If no person is found, query
1677
                 the image directly
                     return image_patch.simple_query
          15
1679
                 ("What does the person to the left
1680
                 of the helmet hold?")
1681
                 person_patch = person_patches[0]
          16
1682
                 helmet_patches = image_patch.find("
          17
                 helmet")
1683
1684
                 # Question assumes only one helmet
          18
1685
                 patch
```

```
if len(helmet_patches) == 0:
                                                      1686
19
           return "no"
                                                      1687
                                                      1688
          helmet_patch in helmet_patches:
       for
           if helmet_patch.
                                                      1689
      horizontal_center > person_patch.
                                                      1690
      horizontal_center:
               return helmet_patch.
      simple_query("What does the person
      to the left of the helmet hold?")
                                                      1694
       return "no"
                                                      1695
  intermediate_result:
       'image_patch': 'a baseball game',
                                                      1697
       'person_patches': 'lens:4',
                                                      1698
       'len(person_patches)': '4'
                                                      1699
       'if len(person_patches) == 0?': '
                                                      1700
      False',
                                                      1701
       'person_patch': 'a baseball player
                                                      1702
      holding a bat',
                                                      1703
       'helmet_patches': 'lens:3'
                                                      1704
       'len(helmet_patches)': '3',
                                                      1705
       'if len(helmet_patches) == 0?': '
                                                      1706
      False',
                                                      1707
       'helmet_patch.horizontal_center':
                                                      1708
                                                      1709
       '333.5'.
       'person_patch.horizontal_center':
                                                      1710
       '230.0'.
                                                      1711
       'if helmet_patch.horizontal_center >
                                                      1712
       person_patch.horizontal_center?':
                                                      1713
                                                      1714
      True'
                                                      1715
  return value: right
                                                      1716
                                                      1717
                                                      1718
40
  Summary of intermediate variables: The
      image shows a baseball game, there
                                                      1719
                                                      1720
      are 4 persons in the image and 3
                                                      1721
      helmet.
41 Possible problems: The code assume that
                                                      1722
      there are only one person patch, but
      actually there are 4 persons. The
                                                      1724
      code is missing loop structure for
                                                      1725
      handling of different situations
                                                      1726
      when the number of people is greater
                                                      1727
       than 1.
                                                      1728
                                                      1729
                                                      1730
  (Example2)
                                                      1731
                                                      1732
  # Which color is the device the monitor
      is to the right of?
                                                      1733
46
  def execute_command(image)->str:
                                                      1734
       image_patch = ImagePatch(image)
                                                      1735
       monitor_patches = image_patch.find("
                                                      1736
      monitor")
                                                      1737
       # Question assumes only one monitor
                                                      1738
                                                      1739
      patch
                                                      1740
       if len(monitor_patches) == 0:
           # If no monitor is found, query
                                                      1741
      the image directly
                                                      1742
           return image_patch.simple_query
                                                      1743
      ("Which color is the device the
                                                      1744
      monitor is to the right of?")
                                                      1745
                                                      1746
      for monitor_patch in monitor_patches
                                                      1747
                                                      1748
           device_patches = image_patch.
                                                      1749
       find("device")
           for device_patch in
                                                      1750
      device_patches:
                                                      1751
                                                      1752
               if device_patch.
      horizontal_center < monitor_patch.</pre>
                                                      1753
      horizontal_cent er:
                                                      1755
          return device_patch.
```

```
1756
                 simple_query("Which color is the
1757
                 device?")
                 return "unknown"
1758
          58
1759
          59
1760
            intermediate_result:
    "image_patch": "a computer monitor
          60
          61
1762
                 and a laptop computer on a desk",
                    "monitor_patches":"lens:2",
1763
          62
                    "len(monitor_patches)":"2",
1764
          63
1765
                    "if len(monitor_patches) == 0?":"
          64
                 False"
1766
                                                           10
                    "device_patches":"lens:4",
1767
          65
                   "device_patch.horizontal_center
1768
          66
1769
                 ":"285.0"
                    "monitor_patch.horizontal_center":
1770
          67
                  "319.0",
1771
1772
                   "if device_patch.horizontal_center
          68
1773
                  < monitor_patch.horizontal_center
1774
                 ?": "True"
1775
          69
             return value: right
1776
          70
                                                           10
1777
          71
                                                           10
1778
          72 Summary of intermediate variables: The
                                                           10
                 image shows a computer monitor and a
                                                           10
1780
                  laptop computer on a desk, There are
                  2 monitor and 4 devices.
1781
          73 Possible problems: Although there are
1783
                 multiple Monitors in the problem,
1784
                 for the device, the judgment
1785
                 condition is already satisfied at
1786
                 the first device, and the answer is
                 returned without obtaining the
1787
1788
                 colors of other devices that meet
1789
                 the condition later. The code is
1790
                 lack of loop structure.
          74
1792
             (Example3)
          75
              In which part of the photo is the
1793
          76
             #
                 bridge, the bottom or the top?
1794
1795
          77
             def execute_command(image)->str:
                 image_patch = ImagePatch(image)
1796
          78
                 bridge_patches = image_patch.find("
1797
          79
                 bridge")
1798
1799
                 # Question assumes only one bridge
          80
1800
                 patch
                 if len(bridge_patches) == 0:
          81
                      # If no bridge is found, query
          82
1803
                 the image directly
                      return image_patch.simple_query
          83
1805
                 ("In which part of the photo is the
                                                            5
1806
                 bridge, the bottom or the top?")
                 if bridge_patches[0].vertical_center
1807
          84
1808
                  > image_patch.vertical_center:
1809
                      return "bottom"
          85
1810
          86
                 else:
                      return "top"
1811
          87
1812
          88
            intermediate_result:
          89
1814
             image_patch:A bridge cross the river.
          90
1815
          91
             bridge_patches:lens:1
             if bridge_patches[0].vertical_center >
1816
          92
1817
                 image_patch.vertical_center? True
1818
          93
1819
          94
             Summary of intermediate variables: The
1820
          95
1821
                 image shows 1 bridge cross the river
1822
1823
          96
             Possible problems: The judgment of up
          97
1825
                and down in the code is reversed.
```

	The patch with a larger vertical	1826
	center should be placed above, and	1827
	similarly the patch with a larger	1000
	Similarly, the patch with a larger	1020
	horizontal center should be placed	1829
	on the right side of the image. The	1830
	code snippet uses incorrect	1831
	attributes in the selection	1832
	structure.	1833
8		1834
0	(Examples finished)	1835
	(Examples (Infined)	1000
00		1030
)1	Your return format should be as follows:	1837
)2	Summary of intermediate variables:	1838
	Included length of list and text	1839
	description.	1840
13	Possible problems. The possible problems	1841
5	10331bie problems. The possible problems	10-11
	nide in code based on query and	1842
	summary of intermediate variables.	1843
)4		1844
)5	The code is:{code}	1845
)6	The intermediate result:{inter res}	1846
17	return value.{return value}	1847
	Cummony of intermedicto verichles.	10/10
18	Summary of Intermediate Variables:	1848
)9	, , ,	1849

Listing 10: Prompt for result asscessment and review generation phase of potential error.

C.4	Prompt For	Code Error	Reviewer	
------------	------------	-------------------	----------	--

f'''Error in execution of statement {	1851
<pre>code}:{error}.</pre>	1852
Here are 2 error types of code execution	1853
	1854
Syntax Error: Syntax errors are the most	1855
basic type of error, typically	1856
occurring when the written code does	1857
not adhere to the syntax rules of	1858
the programming language. For	1859
example missing semicolons	1860
mismatched parentheses. or spelling	1861
errors in variable names.	1862
Runtime Error: Runtime errors occur	1863
during program execution and are	1864
often caused by invalid input data	1865
memory overflow or division by zero	1866
Such errors can notentially lead	1867
to program crashes	1868
	1869
What is the type of {error}? If its a	1870
syntax error please provide the	1871
corrected code directly If its a	1872
runtime error add exception	1873
handling statements at the	1874
corresponding location in the	1875
original code	1975
original code.	1970
fixed code:'''	1077
TIXEU_COUC.	1070

Listing 11: Prompt for code error reviewer on review generation phase.

D Prompt For Rationale Generation 1879

```
f'''You are currently assisting another
    artificial intelligence in
    completing a Visual Question
    Answering (VQA) task. You need to
    break down the solution process for
    1884
```

```
1885
                the following question into 1-5 sub-
1886
                steps, allowing the AI to invoke
                tools for each sub-step and
1887
1888
                ultimately arrive at an answer.
                Example as follows:
1889
1890
          2 (Example1)
          3 Ouerv: Where is this?
          4 1. Identify prominent landmarks or
1892
1893
                features in the image that can help
1894
                determine the location.
1895
           5 2. Analyze the surroundings of the
                landmarks or features to gather
1896
1897
                contextual information that could
                indicate the location.
1898
           6 3. Based on contextual information,
1900
                infer the most likely location
1901
                depicted in the image.
          7
1903
          8 (Example2)
1904
          9 Query: What kind of car is near the pond
1905
1906
          10 1. Find the pond in the image.
1907
          11 2. Determine which cars are near the
                pond based on their positions
                relative to the pond's location.
1910
          12 3. Identify the type or model of the car
1911
1912
          13 (Examples finished)
1913
          14
1914
          15 Query:{query}
1915
          16
```

Listing 12: Prompt for rationale generation in code generation satge.

E Prompt For Supervised Prompt Tuning

1916

```
f'''As a professional prompt engineer
1917
1918
                specializing in prompt optimization
                based on review feedback, I'll help
1920
                you create a template for analyzing
1921
                reports and optimizing prompts:
1922
          2
1923
            As an experienced Prompt Optimization
           3
1924
                Specialist, please carefully analyze
1925
                 the following Review Report and
1926
                optimize the Prompt Template
1927
                accordingly.
          4
1929
          5 Review Report:
          6 {Review Report}
1931
1932
          8 Original Prompt Template:
1933
          9 {Prompt Template}
1934
          10
1935
          II Please follow these steps for analysis
1936
                and optimization:
1937
1938
          13 Analyze key issues and improvement
                suggestions mentioned in the Review
1940
                Report
          14 Based on this feedback, evaluate the
1942
                original prompt in the following
1943
                aspects:
1944
          15 Clarity and specificity of instructions
1945
          16 Accuracy of role definition
          17 Output format standardization
1946
1947
          18 Reasonableness of constraints
1948
          19 Helpfulness of examples
```

20	Propose specific optimization	1949
	suggestions	1950
21	Generate an optimized version of the	1951
	prompt	1952
22	Optimization focus:	1953
23		1954
24	Maintain the core functionality of the	1955
	prompt	1956
25	Enhance instruction executability	1957
26	Improve output consistency	1958
27	Supplement necessary context	1959
28	Perfect constraints	1960
29	Improve example explanations	1961
30	Please provide:	1962
31		1963
32	Problem analysis and optimization	1964
	approach	1965
33	List of recommended modifications	1966
34	Complete optimized prompt	1967
35	I will optimize according to the	1968
	original prompt's core intent while	1969
	addressing issues raised in the	1970
	review.',	1971

Listing 13: Prompt for supervised prompt tuning.

F Additional Results

F.1 GQA Validation Result

We also evluate Seper on GQA Validation set, which contains 2000 balanced question typed samples. As shown in 4, our method is two percent better than current visual programming method. 1973

1974

1975

1976

1978

1979

1980

1981

1982

1983

1984

1985

1986

1987

1988

1989

1990

1991

1992

1993

1994

1995

1996

F.2 Prompt Tuning Experiment

We selected the GQA training dataset to tune the prompt templates in Seper. During each tuning iteration, we categorize the samples into four types: 1) **Invalid**: Both the previous and current prompt templates result in incorrect answers for these samples. 2) **Valid**: The previous prompt templates result in incorrect answers, but the current iteration leads to correct answers. 3) **Toxic**: The previous prompt templates result in correct answers, but the current iteration leads to incorrect answers. 4) **Other**: Both the previous and current prompt templates result in correct answers.

As shown in Table 3, the prompt tuning process significantly reduces the number of toxic samples, with the 5-th iteration leading to a modest overall increase in accuracy from 55.3% to 57.2%. Consequently, we used the prompt templates in the 5-th iteration for subsequent experiments.

F.3 Accuracy Analysis across Question Types

8 As shown in Figure x, we categorized the samples in GQA according to the variation in accuracy1998between Seper and the baseline. The five categories2000are: 0 to 50% increase, 50 to 100% increase, over2001

Iteration	Other	Valid	Invalid	Toxic	Acc
0	-	-	-	-	55.3
1	44.2	10.9	33.8	11.1	55.1
3	44.8	10.7	34.0	10.5	55.5
5	46.2	11.0	33.7	9.1	57.2
10	46.0	10.8	33.9	9.3	56.8

Table 3: Accuracy results on the GQA dataset across iterations in supervised prompt tuning.

Table 4: Result On GQA Validation.

	Accuracy(%)
ViperGPT	55.2
CodeVQA	55.3
Seper	57.8

100% increase, 0 to 10% decrease, and 10 to 100% decrease. Nearly 70% of the samples achieved accuracy improvements, with 17.6% of the question categories seeing over 100% increases in accuracy, demonstrating the effectiveness of the Seper framework. However, 30% of the samples still saw accuracy declines, and around 3.9% of the question categories had accuracy drops of over 10%, which include the following three categories:

2003

2007

2008

2011 2012

2013

2014

2015

2016

2017 2018

2019

2020

2022

Туре	Baseline (%)	Seper (%)
activity	70.71	58.49 (-17.27%)
verifyAttrCThis	39.12	13.83 (-64.65%)
verifyAttrs	84.74	70.50 (-16.80%)

 Table 5: Accuracy Reduction for Problematic Categories

	Thr	ough	sam	pling	the e	examp	les for	these
t	hree	ques	tion	categ	gories,	, we	found	d that
t	hese	type	s of	que	stions	are	more	prone
t	o o'	verloo	king	glol	bal i	nform	ation	issues.
1	<pre>def execute_command(image)->str:</pre>							
3		retur	n im	age_p	atch.	simpl	e_quer	у("

Listing 14: Code generated by Viper.

For example, for the 'activity' category question "What is the man doing", Viper's generated code is as shown in Listing 14. Seper's generated code is as shown in Listing 15. The difference is that Seper generated more detailed steps - it first locates the specific person, and then judges their action. However, this also led to a lack of global



Figure 8: Accuracy Variation between Seper and Baseline in GQA Samples.



Figure 9: Accuracy Variation between Seper and Baseline in GQA Samples.

information in the judgment, as shown in Figure 9. In Seper's case, since it was unable to perceive the existence of the Frisbee, it was unable to answer the question correctly. 2023

2026

We considered emphasizing the role of global 2027 information in the review process, such as using an image-captioning model to generate descriptions of images and enforcing these descriptions as reference information for reviewers during their 2031 evaluation. However, we found that this approach 2032 could also introduce potential biases due to inaccuracies or distortions in the captions generated by the image-captioning model, thereby negatively impacting the review process. As a result, we ultimately decided to retain our current approach. In 2037 future work, we will explore more effective ways to incorporate global information into the review process while minimizing hallucinations or errors. 2040



Figure 10: (a) Proportion of samples with different iteration counts for Seper in GQA, we divided samples with different iteration counts into three groups: A, B, and C. These three groups underwent 1, 2, and 3 iterations respectively. (b) Accuracy of Seper and Viper in groups A, B, and C samples.

```
1 def execute_command(image)->str:

2     image_patch = ImagePatch(image)

3     # subtask1: find the man

4     man_patch = image_patch.find("man")

5     # subtask2: Identify the man's

activity

6     return man_patch.simple_query("What

     is the man doing?")
```

Listing 15: Code generated by Seper.

F.4 Group Analysis

Based on Seper's iteration counts, we categorized the GQA validation samples into three groups: A, B, and C, corresponding to 1, 2, and 3 iterations, respectively. Viper, in contrast, can be regarded as having executed only a single iteration across all sample types.

Figure 10(a) illustrates the distribution of samples by iteration count, with approximately onefourth undergoing a second iteration by Seper, and only about 5% necessitating a third. Figure 10(b) contrasts Seper and ViperGPT's performance, showing a declining accuracy trend across the groups. This trend highlights significant difficulty variation, with Group A being the easiest and Group C the most challenging. In Group A, where only one iteration was required, the self-correcting framework shows minimal deviation from ViperGPT. However, in the more complex Group B, the self-correcting framework surpasses ViperGPT by approximately 15% in accuracy. Even in Group C, the most difficult, Seper demonstrates a clear advantage, indicating that Seper's ability to address complex problems is substantially enhanced by self-correcting techniques.

As shown in Table 6, we conducted a comparative analysis of the average time consumption per sample between Seper and Viper. The results indicate that our method requires an average of 3.12

Metric	Viper	Seper
Avg Time	2.37s	3.12s
		71% × 1
Iteration Distribution	$100\% \times 1$	$24\% \times 2$
		$5\% \times 3$

 Table 6: Time Consume Comparison Between Viper and Seper.

seconds per sample, while Viper completes the task in 2.37 seconds, representing a 31.6% increase in processing time. This elevated time consumption is primarily attributed to the multiple rounds of interaction with the Large Language Model required for code generation in our approach. In future work, we aim to optimize these temporal costs through various strategies, such as minimizing prompt length or deploying local LLM implementations.

2071

2072

2074

2075

2076

2077

2079

2080

2041

2042