HumanEval-V: Benchmarking High-Level Visual Reasoning with Complex Diagrams in Coding Tasks

Anonymous ACL submission

Abstract

Understanding and reasoning over diagrams is a fundamental aspect of human intelligence. While Large Multimodal Models (LMMs) have demonstrated impressive capabilities across 005 various tasks, existing benchmarks lack comprehensive evaluation of their diagram interpretation and reasoning abilities, particularly in coding contexts. We present HumanEval-V, a rigorous benchmark of human-annotated coding tasks that spans six task types and evaluates 011 diverse visual reasoning capabilities. Each task features carefully crafted diagrams paired with 013 function signatures and test cases, employing novel code generation tasks to thoroughly as-015 sess models' diagram comprehension. Through extensive experiments with 22 LMMs, we find that even top-performing models achieve mod-017 est success rates, with Claude 3.5 Sonnet reach-019 ing only 36.8% pass@1, highlighting substantial room for improvement. Our analysis reveals that current LMMs struggle with spatial 021 transformations, topological relationships, and dynamic patterns that humans find intuitive. These findings provide valuable insights for advancing LMMs' visual reasoning abilities.¹

1 Introduction

027

034

035

040

High-level intelligence, whether in humans or advanced AI systems, requires the ability to understand and reason over visual information represented in diagrams. Diagrams are essential in many domains, including science, engineering, and mathematics, as they serve as a powerful medium for abstracting and communicating complex data, relationships, and processes, encoding rich information in a visual and structured format. The abilities required to comprehend diagrams extend beyond simple pattern recognition; they necessitate sophisticated cognitive capabilities, including interpreting transformation patterns, recognizing hierarchical structures, and integrating multiple visual cues

A B C D	<pre>Function Signature: def solution(input_matrix: List[List[str]]) -> List[List[str]]:</pre>
D C B A / * +	Transform the input matrix based on the shown pattern.
1, ' L	- Parameters:
D C B A / * - +	input_matrix: Input matrix as a 2d array. - Returns:
A B C D + - * /	output_matrix: Transformed matrix as a 2d array. """
Test Cases: assert solution([['''','a'], ['3','#']])==[['a','1'], ['#','3']]; assert solution ····
Ground Truth Cod	e Solution:
det solution(input	_matrix):
# Reverses ead	h row in the given matrix
return [row[::	-1] for row in input_matrix]

Figure 1: A task example from HumanEval-V. LMMs are required to figure out the facts and patterns in the diagram and complete the function body.

such as arrows, symbols, and their relative positions to perform spatial or logical reasoning. The rapid development of Large Multimodal Models (LMMs) has led to the creation of various benchmarks designed to assess the alignment between LMMs' capabilities and human intelligence. However, there remains a significant gap in benchmarks that specifically evaluate the ability to understand and reason over complex diagrams. 041

042

044

045

051

054

057

059

060

061

062

063

064

065

067

068

Popular multimodal benchmarks, such as MMMU (Yue et al., 2024), MathVista (Lu et al., 2023), and ChartQA (Masry et al., 2022), focus primarily on scientific, mathematical, and chartbased analytical questions over various domains of images, testing LMMs' multidiscipline knowledge rather than diagram understanding. While abstract visual reasoning tasks (Zhang et al., 2019; Nie et al., 2020) from IQ tests typically feature static patterns based on visual analogies or numerical inference over diagrams, they often lack the complexity and diversity in visual patterns and diagram types. This gap highlights the need for benchmarks that assess more intricate diagram reasoning abilities. The field of coding presents an underexplored opportunity, as developers frequently use various diagrams to illustrate data structures, algorithms, and problem constraints. A recent study, MMCode (Li et al., 2024b), evaluated LMMs on coding prob-

¹Data and code available at this anonymous repository



Figure 2: The six task categories in HumanEval-V, along with their quantitative distribution and representative diagram examples.

Figure 3: Core knowledge required for understanding diagrams in HumanEval-V.

104

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

130

131

132

133

134

135

136

138

lems with visual contexts directly crawled from competition platforms. However, competition coding problems often include comprehensive textual descriptions, making the visual information supplementary. Their results, which showed similar LMM performance with and without visual information, further underscore the gap in evaluating genuine diagram understanding abilities.

075

087

095

100

101

103

To address this gap, we introduce HumanEval-V, a novel benchmark designed to provide a focused evaluation of complex diagram understanding and reasoning abilities in programming contexts. Unlike MMCode, our benchmark is dedicated to assessing visual capabilities through a rigorous annotation pipeline that creates coding tasks capturing the essence of real-world problems. Each task features an indispensable, self-explanatory diagram with minimal textual clues, as demonstrated by our experiments where top LMMs failed all tasks without the provided diagrams. HumanEval-V consists of 253 human-annotated coding tasks. Each task features (1) a diagram encoding the problem context, (2) a function signature defining the task's input-output structures, and (3) test cases to verify solution correctness. Figure 1 provides an example task, where the diagram illustrates spatial transformation patterns, requiring the model to comprehend fine-grained visual elements such as matrices, arrow directions, and spatially ordered data points. This task aligns with the ARC-AGI (Chollet, 2019) benchmark in inferring transformation patterns from limited visual examples. However, unlike ARC's matrix-formatted diagrams, HumanEval-V offers a more diverse and complex set of diagrams spanning six task types (Figure 2), demanding versatile capabilities (Figure 3) for diagram understanding and reasoning. For a comparison of diagrams from existing benchmarks and ours, see Figures 21 and 22.

Another novelty of HumanEval-V lies in using code generation tasks for evaluation instead of the multiple-choice or short-answer questions commonly used in other multimodal benchmarks. This approach offers compelling benefits: code generation is more challenging, requiring comprehensive logical thinking and visual understanding with minimal chance of correct guesses, and test cases could rigorously verify whether the model captures all critical visual information, rather than relying on similarity matching with ground truth. Additionally, we utilize a two-stage evaluation pipeline that supports LMMs with limited coding abilities by first prompting them to generate a structured diagram description summarizing the visual context, then using a more capable coder model to implement the solution, ensuring the evaluation prioritizes visual understanding over coding proficiency.

Through extensive experiments with 22 LMMs, we observe the following key findings: (1) Our benchmark presents unique challenges not addressed by other multimodal benchmarks. The topperforming model, Claude 3.5 Sonnet, achieves 36.8% pass@1, while the best open-weight model, Pixtral 124B, reaches 21.3%. (2) Current LMMs exhibit stronger vision-to-language alignment than vision-to-code. Their best performance occurs when they serve as diagram describers, with GPT-40 acting as the coder model. (3) LMMs' performance can be further enhanced through sampling or iterative self-refinement. For instance, Claude 3.5 Sonnet achieves a 74.3% pass rate with 100
samples, and it can reach 55.3% pass@1 with four
self-refining iterations based on test case execution
feedback. (4) Current LMMs still have difficulty
understanding diagrams that are trivial for humans,
particularly understanding spatial transformations,
topological relationships, and dynamic patterns.

2 Benchmark Construction

146

147

148

149

150

151

152

153

154

155

156

157

158

159

161

162

163

164

165

166

167

168

170

171

172

173

174

175

176

178

179

183

184

187

Task Definition: As shown in Figure 1, each coding task in HumanEval-V includes: (1) a single diagram D providing visual context, (2) a Python function signature σ with input parameters, return type, and brief instructions, and (3) a set of test cases $T = t_1, t_2, \ldots, t_n$ to validate the correctness of the generated output O, a complete Python function produced by the LMM.

Annotation Standards: We establish rigorous standards to ensure high-quality coding tasks: (1) The visual context must be essential for solving the task, with all relevant information contained in a single image; (2) Tasks should be designed around the visual context with minimal textual description; (3) Unnecessary programming complexities, such as recursion, intricate constraints, or complex data structures, should be avoided.

Task Annotation Pipeline: We define the task annotation pipeline with *four steps* as in Figure 4. In the first step, we collect a large set of coding problems from prominent Q&A and coding challenge platforms that incorporate images, then screen them to exclude questions that (1) require specific programming frameworks or libraries, or contain images that (2) are not illustrative diagrams, (3) provide no useful information for solving the problem, or (4) require extensive textual context for interpretation. For the second step, we distill the screened problems to identify critical visual elements, along with the data structures, operations, transformations, or conditional rules involved, categorizing them into six task types (Figure 2) and outlining the key capabilities needed for diagram understanding. For the third step, we design new coding tasks based on these distilled ideas, eliminating unnecessary complexities, refining input/output structures and function signatures, and creating visual objects and layouts using basic shapes and a consistent color scheme in PowerPoint. We also generate tailored test cases, resulting in 100 newly crafted seed tasks after excluding tasks with de-



Figure 4: HumanEval-V task construction pipeline.

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

sign or formulation challenges. *For the fourth step*, we expand the task set by diversifying the seed tasks using a hybrid approach with GPT-40. GPT-40 identifies relevant capability aspects for each seed task and suggests modifications involving new spatial transformations, mathematical operations, dynamic patterns, or variations in data structures and object attributes. Human annotators then refine and annotate these new tasks and diagrams, creating 0 to 2 diversified versions per seed task based on complexity, culminating in *253* tasks in HumanEval-V and finalizing the capability aspects shown in Figure 3. Further details and examples of the data collection and diversification processes are provided in Appendix B.1 and Appendix B.2.

Quality Assurance: Our annotation team comprises four experienced programmers, each with over four years of Python programming experience. Initially, each annotator independently annotates their assigned tasks following pre-defined guidelines. Subsequently, all annotators review each other's work by annotating ground truth code solutions and diagram descriptions to ensure tasks are visually grounded, solvable with the provided information, and free of design or conceptual errors. Any identified issues are resolved collaboratively,



Figure 5: Distribution analysis of the benchmark data.

with tasks finalized only after consensus is reached.
Additionally, one annotator ensures consistent formatting and style across all visual representations
and coding tasks. Each annotator contributes over
200 hours to the annotation process.

219

220

224

231

232

236

Benchmark Statistics: To further demonstrate the quality of our benchmark, we conduct statistical analyses on several key aspects and present the distribution of these statistics in Figure 5. First, we strictly control diagram sizes, capping the maximum width or height at 1024 pixels to eliminate the need for high-resolution perception. Second, each task includes at least five test cases, with the majority containing ten, ensuring full statement and branch coverage over human-annotated code solutions. Third, the token length of humanannotated diagram descriptions is predominantly around 400 (measured using tiktoken (OpenAI, 2024c)), demonstrating that our diagrams encapsulate rich visual context. Lastly, our humanannotated code solutions exhibit cyclomatic complexity (Gill and Kemerer, 1991) levels comparable to HumanEval (Chen et al., 2021), a widely used coding benchmark designed for entry-level programming tasks.

3 Benchmarking Setup

Models: We evaluate 22 state-of-the-art LMMs,
including a representative mix of leading proprietary and open-weight models. Our evaluation covers five of the latest proprietary modclaude 3.5 Sonnet, GPT-4o, GPT-4o-mini,
Gemini 1.5 Pro, and Gemini 1.5 Flash. We



Figure 6: Evaluation pipelines employed in the experiments: (1) Direct translation of visual context into code (V2C), (2) with an optional Chain-of-Thought prompt (V2C w/ CoT); (3) Translation of visual context into a textual description, which is then processed to generate code (V2T2C); and (4) A variant of the third pipeline, where a stronger coder model is used to generate the code solution (V2T2C w/ SC).

also assess 17 top-performing open-weight models spanning various parameter sizes, including InternVL 2.5 (4/8/26/78B), Qwen2-VL (7/72B), Pixtral (12/124B), LLaVA-OV (7/72B), Llama-3.2-V (11/90B), Molmo-D (7/72B), Chameleon (7/30B), and Phi-3.5-V (4B). Further details are in Table 4.

247

248

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

Prompting: we employ multiple strategies for our evaluation pipelines as illustrated in Figure 6, where LMMs may encounter four different prompting scenarios: (1) Direct Code Generation. The model directly generates code based on the given diagram D and function signature σ , denoted as $P_{V2C}(D,\sigma)$; (2) Chain of Thought (CoT). This variant enhances the V2C pipeline by incorporating a zero-shot CoT instruction I_{CoT} (Wei et al., 2022), prompting the model to outline its reasoning process before generating the code. This is denoted as $P_{V2C}(D, \sigma, I_{CoT})$; (3) Intermediate Textual Representation. The model first produces a structured textual *problem specification* PS based on D and σ , denoted as $P_{V2T}(D, \sigma)$. The problem specification consists of three key sections: Problem Restatement, Visual Facts, and Visual Patterns. This structured representation is derived from our benchmark annotation process, which we found to be effective in capturing a comprehensive description of the problem context; (4) Code Generation from Text. The model generates code based on PS rather than the original diagram D, denoted as $P_{T2C}(PS, \sigma)$. The corresponding prompt templates are shown in Figure 27, with further details on prompt design available in Appendix C.2.

Models	V	2C	V2C	w/ CoT	V21	F2C	V2T2C	w/ GPT-4o
pass@k	k=1	<i>k</i> =3	k=1	<i>k</i> =3	k=1	k=3	k=1	k=3 🎖
			Propri	etary LMN	/Is			
Claude 3.5 Sonnet	28.1	37.9	36.8 ↑8.7	47.9 ↑10.0	33.2 ↑5.1	43.6 ^{†5.7}	31.6 ^{3.5}	43.7 ↑5.8
GPT-40	<u>24.1</u>	<u>33.8</u>	<u>27.7</u> ↑3.6	<u>40.0</u> ↑6.2	26.5 ^{+2.4}	<u>40.5</u> ↑6.7	26.5 ^{†2.4}	<u>40.5</u> ↑6.7
Gemini 1.5 Pro	23.3	26.9	22.9 _{10.4}	34.1 ^{†7.2}	<u>28.5</u> ↑5.2	36.4^9.5	<u>26.9</u> ↑3.6	37.3110.4
Gemini 1.5 Flash	15.4	20.5	17.4 ^{2.0}	$24.9_{\uparrow 4.4}$	15.8 _{10.4}	$22.0{\scriptstyle\uparrow1.5}$	18.6 ^{†3.2}	27.2↑6.7
GPT-4o-mini	9.90	16.0	15.8 ↑5.9	$21.1{\scriptstyle\uparrow5.1}$	14.2 ↑4.3	22.7↑6.7	18.2↑8.3	24.6
	Open	-weight	LMMs w	ith more th	an 70B pa	rameters		
Pixtral 124B	12.6	20.3	16.6 ↑4	28.1 ↑7.8	21.3 ↑8.7	29.9 ↑9.6	<u>21.3</u> ↑8.7	31.6 ^{11.3}
InternVL 2.5 78B	<u>12.3</u>	<u>19.7</u>	13.4 ^{1.1}	<u>27.3</u> ↑7.6	<u>17.8</u> ↑5.5	<u>25.7</u> ↑6	21.7 ↑9.4	<u>31.4</u> ↑11.7
Qwen2 VL 72B	9.10	15.7	<u>14.2</u> ↑5.1	19.4 ^{3.7}	10.7 ^{1.6}	19.1 ^{+3.4}	16.6↑7.5	25.1
LLaVA-OV 72B	6.70	7.70	6.30,0.4	11.4 ^{3.7}	10.7 ↑4	13.1	13.8 ↑7.1	19.7 ¹²
Molmo-D 72B	3.20	4.80	3.20	$8.80{\scriptstyle\uparrow4.0}$	1.60	7.00 ^{↑2.2}	5.10 ^{1.9}	14.219.4
Llama-3.2-V 90B	4.30	6.10	4.00	8.20 ^{+2.1}	5.90 ↑1.6	$10.9_{\uparrow 4.8}$	4.70 <u>↑0.4</u>	11.0 ^{+4.9}
	Open-	weight	LMMs w	ith fewer th	an 70B pa	arameters		
Pixtral 12B	4.0	<u>5.9</u>	6.3 ^{2.3}	12.2 ↑6.3	5.5 ^{1.5}	12.6 ^{+6.7}	13.8 ↑9.8	21.3 ↑15.4
InternVL 2.5 26B	<u>3.6</u>	6.6	<u>4.3</u> <u>10.7</u>	<u>6.8</u> 10.2	2.840.8	<u>6.7^{10.1}</u>	<u>8.3</u> ↑4.7	<u>16.7</u> ↑10.1
Qwen2 VL 7B	0.8	3.3	1.6	3.9 _{10.6}	2.4 ^{1.6}	6.3 ^{+3.0}	6.3 ↑5.5	14.7 ^{11.4}
InternVL 2.5 8B	0.8	2.0	0.8	3.7 ^{1.7}	1.2 ^{10.4}	3.3 ^{1.3}	5.1 ^{+4.3}	13.6 ^{11.6}
InternVL 2.5 4B	1.2	4.1	3.2 ↑2.0	3.4.10.7	<u>3.2</u> ↑2.0	5.8 ^{1.7}	5.9 ↑4.7	13.5
LLaVA-OV 7B	2.0	1.9	1.640.4	2.4 _{10.5}	2.0 ^{10.0}	3.2 ^{1.3}	5.1 ^{3.1}	10.2↑8.3
Phi-3.5-V 4B	0.0	0.0	0.0	0.0 <u>^0.0</u>	0.0 <u>^0.0</u>	0.0	5.9 ↑5.9	9.40
Llama-3.2-V 11B	2.0	2.0	1.640.4	3.9 ↑1.9	2.0	5.2 ^{+3.2}	4.0 ↑2.0	8.80^6.8
Molmo-D 7B	1.2	1.0	0.4.0.8	1.8 ^{10.8}	0.8	1.0 <u>^0.0</u>	2.8 ^{1.6}	8.40
Chameleon 7B	0.0	0.0	0.0	0.2 ^0.2	0.0 <u>^0.0</u>	0.0 <u>^0.0</u>	1.2 ^{1.2}	2.20 ^{2.2}
Chameleon 30B	0.0	0.0	0.010.0	0.2 ^(0.2)	0.4	0.0	0.0	1.90 ^{1.9}



Table 1: Performance of LMMs across different settings. Models are ranked based on the δ column. The **best** and <u>second-best</u> performances in each column are highlighted. The numerical values are color-coded to indicate performance changes relative to the corresponding pass@k values in the V2C column: green represents improvement, red indicates decline, and yellow denotes minimal change.

Figure 7: Comparison of LMM performance on HumanEval-V and other popular multimodal benchmarks.

Hyper-parameters & Post-processing: We apply two distinct decoding strategies for both code and description generation. First, we use greedy decoding to produce a single deterministic output, assessing model performance in a constrained setting. Additionally, we employ a sampling method with Top p = 0.95, Top k = 20, and a temperature of 0.8 to generate diverse outputs, allowing us to evaluate the models' ability to produce correct solutions when given multiple attempts. We set the maximum output length to 2048 tokens for both code and description generation. To facilitate extraction, we prompt the models to encapsulate their generated code within Markdown-style code blocks. We then apply an abstract syntax tree parser to detect and retrieve generated import statements, class definitions, and function definitions. These components are concatenated to form the final code solution. An additional ablation study on the temperature setting is presented in Appendix C.3.

278

279

283

284

291

292

295

296

297

Evaluation Metrics Following established practices in code generation evaluation (Chen et al., 2021, 2022), we use the pass@k metric to assess functional correctness. A task is considered solved if at least one of the k selected solutions passes all test cases, and pass@k is the percentage of solved tasks. We report results for k = 1, 3. In the V2C setting, we generate n code samples per task and randomly select k for evaluation. For greedy decoding, n = 1 for pass@1, while for sampling-based evaluation, n = 6 for pass@3. In the V2T2C setting, we first sample six problem specifications per task, then use greedy decoding to generate one code solution per PS, resulting in six solutions per task for pass@3 computation.

4 Benchmarking Results

Main Results: We present the benchmarking results of 22 LMMs in Table 1, covering the four evaluation pipelines introduced in Figure 6. Additionally, Figure 7 provides a correlation analy310

311

312

313

314

315

316

317

298

299



Figure 8: Iterative evaluation pipelines.

321

331

341

Figure 9: Performance of LMMs under the iterative evaluation settings.

sis to illustrate the performance gap between the 318 evaluated LMMs on HumanEval-V and other pop-319 320 ular benchmarks (more details on the correlation analysis are in Appendix A.1). Based on these results, we highlight the following key findings: (1) Our benchmark presents unique challenges not cap-323 tured by other benchmarks. As shown in Figure 7, most evaluated LMMs exhibit significantly larger 325 performance gaps on HumanEval-V compared to 326 other benchmarks. While MMMU demonstrates the highest correlation with our benchmark, its results still lack sufficient discrimination between models. (2) LMMs generally achieve their best performance under the V2T2C w/ GPT-40 setting. This is particularly evident for LMMs with fewer than 70B parameters, which struggle to complete tasks in the V2C setting. These findings validate the importance of decoupling visual understanding from coding abilities. Additionally, CoT prompting and the decoupled V2T2C pipeline show similar performance distributions, with more capable 338 LMMs benefiting more from these enhancements than smaller models. (3) Open-weight LMMs still lag behind top proprietary models. Although highcapacity open-weight LMMs (e.g., Pixtral 124B) 342 outperform the mini/flash versions of proprietary models, they still fall short of the most capable pro-344 prietary LMMs. For smaller-scale models, Pixtral 12B, Qwen2 VL 7B, and InternVL 2.5 4B demonstrate a high performance-to-size ratio. (4) Certain LMMs exhibit anomalously poor performance. Models such as Molmo-D, Llama-3.2-V, and the Chameleon series perform significantly worse than other LMMs of similar scale. Another case is Phi-3.5-V, which appears to lack coding ability, achieving a performance score of 0 in the V2C settings, compared to 9.4% pass@3 when assisted by GPT- 40 for code generation. (5) Additional Results of o1 and QVQ: We also evaluated reasoning-enhanced LMMs that leverage test-time scaling by generating long chain-of-thought (CoT) reasoning. Specifically, we assessed OpenAI o1 (OpenAI, 2024b) and QVQ-72B-Preview (Team, 2024b) under the V2C w/ CoT setting, achieving pass@1 scores of 40.6% and 19.0%, respectively. Our case study reveals that both models still struggles with visual understanding, often failing to identify rules or patterns in the diagrams. Meanwhile, QVQ primarily fails due to excessively long CoT reasoning, with 35% of cases unable to generate a valid code solution within the 20k token limit. These results underscore the complexity of the diagrams in our benchmark. Example cases for o1 and QVQ are shown in Figures 34, 39, 44 and Figures 35, 40, 45.

355

356

357

360

361

362

363

364

365

367

368

369

371

373

374

375

377

378

379

380

382

385

386

390

Iterative Benchmarking: We introduce an iterative benchmarking pipeline to evaluate LMMs' ability to reason over environmental feedback and perform self-refinement-an essential skill for realworld problem-solving. Figure 8 illustrates two types of iterative pipelines derived from the V2Cand V2T2C w/ SC settings. In these pipelines, LMMs must refine either their generated code solutions or textual descriptions based on feedback from the execution environment. To support this process, we design new prompt templates (Figure 28) that guide the refinement steps. Specifically, for each task, LMMs perform an additional iteration if the generated code contains syntax errors or fails to pass all test cases. The feedback includes detailed error messages or the failed test cases' inputs and expected outputs.

For the iterative evaluation, we select the most capable LMMs across different parameter scales,



Figure 10: Performance with increased sample size.

391

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

Figure 11: LMMs' performance F with human problem specifications.

nce Figure 12: LLM-as-Judge ratings for LMMs in the *V2T2C* setting.

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

using greedy decoding and the pass @1 metric. Figure 9 presents the results, where *iter 0* represents the first round of generation without feedback. We observe that LMMs generally improve across iterations, with more capable models achieving larger performance gains. Notably, some models, such as Claude and QVQ, exhibit stronger self-refinement capabilities. We also investigate the cases that are corrected after iterations and find that approximately 90% of these cases are corrected due to the models' improved understanding of the diagram and task. The remaining 10% are cases that fix edge conditions highlighted by the test case feedback. None of these corrections result from hard-coding the exposed test cases into the code solutions.

5 Experimental Analysis

This section presents our analysis of model performance under various settings, including increased sampling sizes, human-annotated problem specifications, the use of GPT-40 judge for rating LMMs' diagram descriptions, and error pattern analysis. Our goal is to examine both the potential and limitations of the LMMs on HumanEval-V. Further analysis is provided in Appendix D, where we explore the co-occurrence of capability aspects required in our benchmark tasks, the stability of using QwenCoder-32B as a strong coder instead of GPT-40, and experimental evidence supporting the value of tasks diversified from the seed tasks.

420Performance with Increased Sample Size: We421scale up the number of samples for five proprietary422LMMs to explore their potential performance. We423increase the sampling number n to 200, using the424same Top_p , Top_k , and max token limitations425outlined in Section 3 to calculate the pass@100426score under the V2C w/ CoT setting. As shown

in Figure 10, we observe a consistent performance improvement across all models with larger sample sizes. Notably, Claude 3.5 Sonnet achieves a significant improvement, reaching 74.3% pass@100, underscoring the strong potential for these models when scaling up sample sizes.

Coding Performance with Human-Annotated Problem Specifications: We evaluate all LMMs on a new task where they generate code based on human-annotated problem specifications, without direct access to the diagrams. This setup isolates their ability to perform visual reasoning and generate code. We also calculate the success parsing rate using Pylint (Wikipedia, 2024), which measures the syntactic correctness of the generated code, independent of its functional accuracy. The results, presented in Figure 11, show that most LMMs demonstrate strong coding capabilities, generally outperforming their best results from Table 1. Notably, GPT-40 achieves 96.5% pass@3, a significant improvement over its 40.5% pass@3 in the V2T2C setting. Smaller models, such as InternVL 2.5 4B, also show substantial improvement.

We also evaluate a setting where LMMs generate code based solely on the function signature, without access to diagrams or descriptions, and find that *none of the five proprietary models are able to pass any tasks*. This underscores the necessity of visual context in our benchmark. These results suggest that current LMMs face more challenges in visual reasoning than coding on HumanEval-V.

LLM-as-Judge Ratings: We evaluate the problem specifications (PS) generated by LMMs in the *V2T2C* setting using GPT-40 as the judge. GPT-40 rates the PS in three dimensions: Basic-Level Perception (identifying basic visual elements), High-Level Comprehension (understanding objects, pat-

terns, transformations, and operations), and Contex-464 tual Interpretation (clear description without vague-465 ness or hallucinations) as outlined in the prompt 466 template shown in Figure 29. Ratings are on a 1-3 467 scale, where 1 indicates severe errors and 3 reflects 468 near perfection in the capability dimension. The re-469 sults, shown in Figure 12, reveal that while LMMs 470 generally excel in basic perceptual abilities, they 471 struggle with high-level comprehension and clar-472 ity of expression. Notably, the performance gap 473 between models is small. We also find the differ-474 ence in ratings between passed and failed tasks is 475 minimal. For example, GPT-40 scores 2.9, 2.0, and 476 1.3 across the three dimensions on passed tasks, 477 compared to average ratings of 2.8, 1.6, and 1.2 478 across all tasks, highlighting the limitations of us-479 ing LLM-as-judge as an evaluation tool. This lack 480 of robustness may stem from rigid comparisons 481 to human-annotated PS, further emphasizing the 482 importance of pass rates as the evaluation metric. 483

Error Analysis: We conducted a comprehensive 484 error analysis to understand the limitations of cur-485 486 rent LMMs in HumanEval-V, as detailed in Appendix E. Our analysis examined correlations be-487 tween model performance and three key factors: 488 task types, general capability dimensions, and spe-489 cific capability aspects. The results reveal that 490 LMMs particularly struggle with tasks involving 491 492 Transformation and Iterative Calculation. And models show notable difficulties with specific ca-493 pabilities such as understanding dynamic patterns 494 (e.g., spirals, circular arrangements) and spatial 495 transformations (e.g., stacking, translation). Inter-496 estingly, our investigation of task difficulty metrics 497 shows that LMM performance correlates poorly 498 with human-perceived difficulty measures, includ-499 ing both programming complexity (measured by cyclomatic complexity) and visual comprehension 501 difficulty (measured by description length). This suggests a fundamental gap in LMMs' visual reasoning capabilities, where even tasks considered trivial by humans can prove challenging for state-505 of-the-art models. For concrete examples of these challenges, we present representative error cases in Figures 31 to 47.

6 Related Work

510

511

512

513

Benchmarks Involving Diagrams: Prior work on multimodal benchmarks can be categorized into several groups: (1) General-purpose multimodal evaluation benchmarks (Yue et al., 2024; Liu et al., 2023; Yu et al., 2023; Li et al., 2023; Ying et al., 2024; Chen et al., 2024a) that assess models' broad multidisciplinary capabilities; (2) Scientific diagram understanding (Lu et al., 2022; Kembhavi et al., 2016); (3) Mathematical visual reasoning (Lu et al., 2023; Wang et al., 2024a; Zhang et al., 2024); (4) Data visualization comprehension (Masry et al., 2022; Wang et al., 2024c; Chollet, 2019) that focus on plots and charts; (5) Abstract reasoning (Zhang et al., 2019; Jiang et al., 2024; Nie et al., 2020; Chia et al., 2024); and (6) Specialized diagram understanding including abstract symbol interpretation and geometric spatial reasoning (Lu et al., 2021; Rahmanzadehgervi et al., 2024). While these benchmarks cover various aspects of visual understanding, they do not address the complex diagrams in the coding context.

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

Multimodal Code Generation: Recent work in multimodal code generation has focused on two main Categories. In the first category, researchers have explored derendering web pages into functional code (Si et al., 2024; Laurençon et al., 2024) and converting scientific figures into their corresponding plotting code (Shi et al., 2024; Wu et al., 2024). The second category includes Programbased VQA approaches, where models leverage pre-defined modules to answer visual questions (Surís et al., 2023; Subramanian et al., 2023). MM-Code (Li et al., 2024b) is the most related coding benchmark to ours, evaluating LMMs' coding abilities using problems with visual demonstrations from competition platforms. However, our benchmark differs in its dedicated focus on assessing the visual capabilities of LMMs. We provide a detailed discussion in Appendix F, highlighting the differences between HumanEval-V and MMCode in terms of visual indispensability, task complexity, and evaluation design.

7 Conclusion

In this paper, we introduced HumanEval-V, a novel benchmark designed to evaluate LMMs' capabilities in understanding and reasoning over diagrams in programming contexts. Through comprehensive experiments, we demonstrated that current LMMs, while showing promising performance, still face significant challenges in complex diagram understanding and reasoning. Our extensive experimental results and analysis provide valuable insights for the future development of more sophisticated visual reasoning abilities in AI systems.

8 Limitations

564

567

Despite the valuable contributions of our benchmark, several limitations remain that we aim to 566 address in future work:

Limited Benchmark Size: The size of our 568 benchmark is constrained by the significant cost 569 of human annotation, as we prioritize highquality task design to ensure meaningful insights, with each annotator dedicating over 200 hours 572 to constructing HumanEval-V. Nevertheless, our benchmark includes 253 tasks, comparable to 574 many well-established human-annotated benchmarks in academia and industry, such as Hu-576 manEval (Chen et al., 2021) with 164 tasks, MM-Vet (Yu et al., 2023) with 218, and Vibe-Eval (Padlewski et al., 2024) with 269. Notably, none of the current popular multimodal benchmarks feature manually drawn diagrams, 581 further distinguishing HumanEval-V. Furthermore, HumanEval-V offers a diverse and balanced set of task types covering a wide range of capability aspects, enabling us to uncover unique insights into 585 the limitations of current LMMs.

Limited Model Coverage: While our experiments evaluate a representative set of top-588 589 performing LMMs, the rapid pace of model development means newly released models may not be covered in our current evaluation. To address this, we plan to publicly release our evaluation toolkit and dataset, along with an up-to-date leaderboard to track ongoing advancements. This will enable benchmarking of new models as they become avail-595 able, ensuring HumanEval-V remains relevant and continuously updated.

Limitations in Exploring Advanced Methods: 598 While our experiments cover various evaluation settings, including chain-of-thought (CoT), iterative refinement, and long-CoT-enhanced LMMs, our exploration of more advanced CoT techniques is limited. Methods such as supervised fine-tuning (Chen et al., 2024b), reinforcement learning (Snell et al., 2024), or more complex CoT approaches (Yao 605 et al., 2024; Mitra et al., 2024) could further enhance LMM reasoning capabilities. However, these techniques are challenging to apply to diagram reasoning due to the lack of high-quality training data in this domain. As our primary objective is to 610 bridge the gap in diagram reasoning benchmarks, 611 we leave the exploration of more sophisticated reasoning-enhancing methods to future work. 613

References

Pravesh Agrawal, Szymon Antoniak, Emma Bou Hanna, Baptiste Bout, Devendra Chaplot, Jessica Chudnovsky, Diogo Costa, Baudouin De Monicault, Saurabh Garg, Theophile Gervet, et al. 2024. Pixtral 12b. arXiv preprint arXiv:2410.07073.

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

- Anthropic. 2024. Claude 3.5 sonnet.
- Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu Chen. 2022. Codet: Code generation with generated tests. arXiv preprint arXiv:2207.10397.
- Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, et al. 2024a. Are we on the right way for evaluating large vision-language models? arXiv preprint arXiv:2403.20330.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. arXiv preprint arXiv:2107.03374.
- Qiguang Chen, Libo Qin, Jin Zhang, Zhi Chen, Xiao Xu, and Wanxiang Che. 2024b. M3cot: A novel benchmark for multi-domain multi-step multi-modal chain-of-thought. arXiv preprint arXiv:2405.16473.
- Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. 2024c. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. arXiv preprint arXiv:2412.05271.
- Yew Ken Chia, Vernon Toh Yan Han, Deepanway Ghosal, Lidong Bing, and Soujanya Poria. 2024. Puzzlevqa: Diagnosing multimodal reasoning challenges of language models with abstract visual patterns. arXiv preprint arXiv:2403.13315.
- François Chollet. 2019. On the measure of intelligence. arXiv preprint arXiv:1911.01547.
- Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, et al. 2024. Molmo and pixmo: Open weights and open data for state-of-the-art multimodal models. arXiv preprint arXiv:2409.17146.
- Haodong Duan, Junming Yang, Yuxuan Qiao, Xinyu Fang, Lin Chen, Yuan Liu, Xiaoyi Dong, Yuhang Zang, Pan Zhang, Jiaqi Wang, Dahua Lin, and Kai Chen. 2024. Vlmevalkit: An open-source toolkit for evaluating large multi-modality models. Preprint, arXiv:2407.11691.
- Geoffrey K Gill and Chris F Kemerer. 1991. Cyclomatic complexity density and software maintenance productivity. IEEE transactions on software engineering, 17(12):1284-1288.

776

Google. 2024a. Introducing gemini 1.5, google's next-generation ai model.Google. 2024b. Llama 3.2: Revolutionizing edge ai and

vision with open, customizable models.

670

673

675

680

706

707

710

711

713

714

715

716

717

718

719

720

721

722

- Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, et al. 2023. Hallusionbench: An advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models. *arXiv preprint arXiv:2310.14566*.
- Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Kai Dang, et al. 2024. Qwen2. 5-coder technical report. arXiv preprint arXiv:2409.12186.
- Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. 2024. Livecodebench: Holistic and contamination free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*.
- Yifan Jiang, Jiarui Zhang, Kexuan Sun, Zhivar Sourati, Kian Ahrabian, Kaixin Ma, Filip Ilievski, and Jay Pujara. 2024. Marvel: Multidimensional abstraction and reasoning through visual evaluation and learning. *arXiv preprint arXiv:2404.13591*.
 - Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi.
 2016. A diagram is worth a dozen images. In Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14, pages 235–251. Springer.
- Hugo Laurençon, Léo Tronchon, and Victor Sanh. 2024. Unlocking the conversion of web screenshots into html code with the websight dataset. *arXiv preprint arXiv:2403.09029*.
- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2024a. Llavaonevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*.
- Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. 2023. Seed-bench: Benchmarking multimodal llms with generative comprehension. *arXiv preprint arXiv:2307.16125*.
- Kaixin Li, Yuchen Tian, Qisheng Hu, Ziyang Luo, Zhiyong Huang, and Jing Ma. 2024b. Mmcode: Benchmarking multimodal large language models for code generation with visually rich programming problems. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 736–783.
- Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi

Wang, Conghui He, Ziwei Liu, et al. 2023. Mmbench: Is your multi-modal model an all-around player? *arXiv preprint arXiv:2307.06281*.

- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2023. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*.
- Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. Advances in Neural Information Processing Systems, 35:2507–2521.
- Pan Lu, Liang Qiu, Jiaqi Chen, Tony Xia, Yizhou Zhao, Wei Zhang, Zhou Yu, Xiaodan Liang, and Song-Chun Zhu. 2021. Iconqa: A new benchmark for abstract diagram understanding and visual language reasoning. *arXiv preprint arXiv:2110.13214*.
- Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. *arXiv preprint arXiv:2203.10244*.
- Microsoft. 2024. Discover the new multi-lingual, highquality phi-3.5 slms.
- Chancharik Mitra, Brandon Huang, Trevor Darrell, and Roei Herzig. 2024. Compositional chain-of-thought prompting for large multimodal models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14420–14431.
- Weili Nie, Zhiding Yu, Lei Mao, Ankit B Patel, Yuke Zhu, and Anima Anandkumar. 2020. Bongard-logo: A new benchmark for human-level concept learning and reasoning. *Advances in Neural Information Processing Systems*, 33:16468–16480.
- OpenAI. 2024a. Hello gpt-4o.
- OpenAI. 2024b. Introducing openai o1.
- OpenAI. 2024c. tiktoken.
- Piotr Padlewski, Max Bain, Matthew Henderson, Zhongkai Zhu, Nishant Relan, Hai Pham, Donovan Ong, Kaloyan Aleksiev, Aitor Ormazabal, Samuel Phua, et al. 2024. Vibe-eval: A hard evaluation suite for measuring progress of multimodal language models. *arXiv preprint arXiv:2405.02287*.
- Pooyan Rahmanzadehgervi, Logan Bolton, Mohammad Reza Taesiri, and Anh Totti Nguyen. 2024. Vision language models are blind. In *Proceedings of the Asian Conference on Computer Vision*, pages 18– 34.
- Chufan Shi, Cheng Yang, Yaxin Liu, Bo Shui, Junjie Wang, Mohan Jing, Linran Xu, Xinyu Zhu, Siheng Li, Yuxiang Zhang, et al. 2024. Chartmimic: Evaluating Imm's cross-modal reasoning capability via chart-tocode generation. *arXiv preprint arXiv:2406.09961*.

860

861

862

863

830

831

Chenglei Si, Yanzhe Zhang, Zhengyuan Yang, Ruibo Liu, and Diyi Yang. 2024. Design2code: How far are we from automating front-end engineering? *arXiv preprint arXiv:2403.03163*.

778

783

784

785

788

790

791

792

795

804

809

810

811

812

813

814

815

816

817

819

825

829

- Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. 2024. Scaling llm test-time compute optimally can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*.
- 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. arXiv preprint arXiv:2306.05392.
- Dídac Surís, Sachit Menon, and Carl Vondrick. 2023. Vipergpt: Visual inference via python execution for reasoning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11888– 11898.
 - Chameleon Team. 2024a. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint arXiv:2405.09818*.
- Qwen Team. 2024b. Qvq: To see the world with wisdom.
 - Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Mingjie Zhan, and Hongsheng Li. 2024a. Measuring multimodal mathematical reasoning with math-vision dataset. *arXiv preprint arXiv:2402.14804*.
 - Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. 2024b. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*.
 - Zirui Wang, Mengzhou Xia, Luxi He, Howard Chen, Yitao Liu, Richard Zhu, Kaiqu Liang, Xindi Wu, Haotian Liu, Sadhika Malladi, et al. 2024c. Charxiv: Charting gaps in realistic chart understanding in multimodal llms. *arXiv preprint arXiv:2406.18521*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- Wikipedia. 2024. Pylint Wikipedia, the free encyclopedia. http://en.wikipedia.org/w/index. php?title=Pylint&oldid=1191495734. [Online; accessed 24-October-2024].
- Chengyue Wu, Yixiao Ge, Qiushan Guo, Jiahao Wang, Zhixuan Liang, Zeyu Lu, Ying Shan, and Ping Luo. 2024. Plot2code: A comprehensive benchmark for evaluating multi-modal large language models in code generation from scientific plots. *arXiv preprint arXiv:2405.07990*.

- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2024. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36.
- Kaining Ying, Fanqing Meng, Jin Wang, Zhiqian Li, Han Lin, Yue Yang, Hao Zhang, Wenbo Zhang, Yuqi Lin, Shuo Liu, et al. 2024. Mmt-bench: A comprehensive multimodal benchmark for evaluating large vision-language models towards multitask agi. *arXiv preprint arXiv:2404.16006*.
- Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. 2023. Mm-vet: Evaluating large multimodal models for integrated capabilities. arXiv preprint arXiv:2308.02490.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. 2024. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567.
- Chi Zhang, Feng Gao, Baoxiong Jia, Yixin Zhu, and Song-Chun Zhu. 2019. Raven: A dataset for relational and analogical visual reasoning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5317–5327.
- Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan Lu, Kai-Wei Chang, Yu Qiao, et al. 2024. Mathverse: Does your multi-modal llm truly see the diagrams in visual math problems? In *European Conference on Computer Vision*, pages 169–186. Springer.

Appendix

865	Α	Com	parison with Other Benchmarks	12
866		A.1	Correlation Analysis	12
867		A.2	Diagrams in Other Benchmarks	12
868		A.3	Diagrams in HumanEval-V	13
869	В	Mor	e Details on Data Annotation	14
870		B .1	Data Collection and Screening	14
871		B.2	Recreation and Diversification	15
872	С	Mor	e details on Experimental Setup	15
873		C .1	Evaluated Models	15
874		C.2	Prompt Templates	15
875		C.3	Ablation on Temperature	16
876	D	Deep	per Analysis on HumanEval-V	16
877		D .1	Co-occurrence of Capability Aspects	16
878		D.2	Comparison of Strong Coder Models	16
879		D.3	The Effect of Diversified Tasks	17
880	E	Mor	e Detailed Error Analysis	19
881		E .1	Error Patterns and Taxonomy	19
882		E.2	Error Analysis by Task Difficulty .	19
883	F	Mor	e Discussion on MMCode	19
884	G	Othe	er Considerations	20

A Comparison with Other Benchmarks

A.1 Correlation Analysis

To assess whether HumanEval-V identifies specific weaknesses not captured by existing benchmarks, we select seven widely used multimodal benchmarks that cover a range of multidisciplinary abilities. These include AI2D (Kembhavi et al., 2016), MMVet (Yu et al., 2023), MMBench (Liu et al., 2023), MMStar (Lu et al., 2023), MMMU (Yue et al., 2024), MMStar (Chen et al., 2024a), and HallusionBench (Guan et al., 2023). Performance results for the 22 LMMs evaluated in this paper are collected from the OpenVLM Leaderboard (Duan et al., 2024), as well as corresponding papers and reports. These results are shown alongside the pass@3 scores for HumanEval-V under the V2T2C w/ GPT-4o setting in Table 2.

From the analysis, we observe that open-weight LMMs with more than 70B parameters generally perform well on the selected benchmarks, with models like Pixtral, InternVL 2.5, and Qwen2 VL even outperforming proprietary models such as GPT-40 and Claude 3.5 Sonnet in several cases. Llama-3.2-V also shows competitive performance. However, open-weight LMMs exhibit significantly lower performance on HumanEval-V, suggesting that our benchmark uncovers model weaknesses that may not be apparent in other benchmarks.

To quantify the relationship between HumanEval-V and the other five benchmarks, we visualize the performance of the 22 LMMs across all benchmarks using regression plots for each benchmark pair in Figure 20. The plots reveal low correlations between HumanEval-V and the other benchmarks, with notable differences in performance across models. Overall, the performance of all models remains lower on HumanEval-V compared to the other benchmarks.

A.2 Diagrams in Other Benchmarks

Figure 21 presents a comprehensive comparison of five distinct categories of diagrams commonly used in various benchmarks and coding platforms, showcasing the diverse range of visual reasoning challenges in the open world. The first category consists of real-world images from benchmarks such as MMMU, MMBench, and MM-Vet, encompassing everyday photographs of food, sports, architecture, art, and wildlife in both color and monochrome formats. These images test general visual recognition and understanding capabilities, contrasting sharply

Models	AI2D	MM-Vet	MMBench	MathVista	MMMU	MMStar	HallusionBench	HumanEval-V 🎖
			Prop	rietary LM	IMs			
Claude 3.5 Sonnet	81.2	66.0	81.7	67.7	65.9	65.1	55.1	43.7
GPT-40	84.9	69.1	84.3	61.3	69.2	65.1	56.2	40.5
Gemini 1.5 Pro	79.1	64.0	82.8	57.5	60.6	59.1	45.6	37.3
Gemini 1.5 Flash	78.5	63.2	76.9	51.2	58.2	55.8	48.5	27.2
GPT-4o-mini	77.8	66.9	76.0	52.4	60.0	54.8	46.1	24.6
			Open-	Weight Ll	MMs			
Pixtral 124B	93.8	-	-	69.4	64.0	-	-	31.6
InternVL 2.5 78B	89.2	64.4	87.7	65.6	58.3	72.1	57.4	31.4
Qwen2 VL 72B	83.0	74.0	81.0	70.5	64.5	25.9	58.7	25.1
LLaVA-OV 72B	86.2	63.7	82.6	67.5	56.6	65.8	47.9	19.7
Molmo-D 72B	83.4	61.1	79.5	55.2	52.8	63.3	46.4	14.2
Llama-3.2-V 90B	92.3	64.1	77.3	57.3	60.3	55.3	44.1	11.0
Pixtral 12B	77.4	58.5	72.7	56.3	44.1	54.5	47.0	21.3
InternVL 2.5 26B	86.2	60.0	84.6	59.4	50.7	66.5	55.8	16.7
Qwen2 VL 7B	88.3	62.0	85.9	58.2	54.1	16.3	50.4	14.7
InternVL 2.5 8B	84.6	54.3	82.5	58.3	51.2	63.2	49.0	13.6
InternVL 2.5 4B	81.4	50.9	78.2	58.1	48.3	58.7	46.6	13.5
LLaVA-OV 7B	82.8	57.5	80.9	63.2	46.8	61.9	31.6	10.2
Phi-3.5-V 4B	77.8	43.2	67.4	43.2	44.6	47.5	40.5	9.4
Llama-3.2-V 11B	91.1	57.6	65.8	51.5	50.7	49.8	40.3	8.8
Molmo-D 7B	79.6	53.3	76.5	46.9	48.7	54.4	47.4	8.4
Chameleon 7B	46.0	8.3	19.8	22.5	22.4	31.1	17.1	2.2
Chameleon 30B	53.7	9.7	32.7	23.8	38.8	32.7	18.6	1.9

Table 2: A performance comparison of 22 LMMs across HumanEval-V and seven popular multimodal benchmarks. Models are ranked according to the \checkmark column. Results for HumanEval-V correspond to the *V2T2C w/ GPT-4o* setting from Table 1. The top two results for each column are highlighted in **bold**.

with the more structured representations found in other categories.

936

937

939

940

943

944

947

949

951

952

954

955

957

958

The second and third categories focus on analytical and scientific visualization. Analytical tables and charts, evaluated through benchmarks like ChartQA (Masry et al., 2022) and Charxiv (Wang et al., 2024c), comprise business and scientific data visualizations including bar charts, line graphs, and frequency tables. Scientific diagrams featured in MMMU (Yue et al., 2024), MMBench (Liu et al., 2023), and ScienceQA (Lu et al., 2022) present technical illustrations of molecular structures, particle dynamics, and ecosystem relationships. While both categories deal with data representation, they differ in their approach: analytical charts emphasize quantitative interpretation, whereas scientific diagrams focus on conceptual understanding.

Mathematical diagrams, assessed through benchmarks such as MathVista (Lu et al., 2023) and Math-Vision (Wang et al., 2024a), represent another crucial category that bridges pure mathematics with practical applications. These include function graphs, geometric constructions, and physics diagrams, demonstrating complex mathematical concepts through visual means. This category shares some common ground with programmingrelated diagrams, particularly in their emphasis on logical relationships and systematic thinking. 959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

The fifth category encompasses visual abstract reasoning, evaluated through benchmarks like ARC-AGI (Chollet, 2019), RAVEN (Zhang et al., 2019), and Bongard (Nie et al., 2020). These tests feature grid-based patterns and geometric transformations that assess abstract thinking and pattern recognition skills. This category bears the closest resemblance to programming-related diagrams in terms of logical abstraction and systematic problem-solving approaches.

A.3 Diagrams in HumanEval-V

Figure 22 presents six fundamental task types in the HumanEval-V benchmark, each representing distinct cognitive challenges in visual reasoning. Our benchmark employs a rich variety of visual elements including geometric shapes, symbolic notations, matrices, and directed graphs. These representations are enhanced through connecting lines, arrows, color-coding, and numerical annotations



Figure 13: Sources of the screened tasks for annotation.

to effectively capture relationships and transformations between components. The visual representations maintain clarity across all categories while scaling in complexity to accommodate different difficulty levels. Through careful design of visual elements and systematic progression of patterns, each task type provides a clear framework for evaluating specific aspects of visual reasoning and problemsolving abilities.

The six task categories demonstrate diverse problem-solving requirements: Aggregation tasks (18% of the benchmark) introduce new grouping and aggregation rules for input data; Validation tasks (17%) define conditional rules to verify or classify input data; Expansion tasks (16%) focus on defining new patterns that evolve or extend data elements; Rearrangement tasks (17%) establish new traversal patterns to reorganize data; Iteration tasks (17%) define new iterative operations applied to input data; and Computation tasks (20%) introduce new computational rules and operations.

What distinguishes our benchmark is not only its balanced distribution across task types but also the wide interconnection between categories. While each category emphasizes specific problem-solving skills, real-world scenarios often require combining multiple approaches. For instance, computation tasks may incorporate iterative processes, while aggregation problems might require validation steps. This interconnected design reflects the complexity of practical problem-solving scenarios where multiple cognitive skills must be applied simultaneously.

B More Details on Data Annotation

B.1 Data Collection and Screening

Our data collection process involves two primary sources: coding challenge platforms, such as Code-Forces, and the Q&A platform Stack Overflow (SO). Each coding problem undergoes a rigorous screening process to ensure it aligns with the standards of HumanEval-V. Annotators are instructed to exclude problems that: (1) require knowledge of specific programming frameworks or libraries, (2) contain images that are not abstract diagrams, (3) provide no useful information for solving the problem, or (4) require excessive textual context for interpretation. 1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

The majority of our tasks are sourced from coding challenge platforms, especially CodeForces, as shown in Figure 13, where we display the distribution of screened tasks by platform. For coding challenge platforms, we use the open-source MM-Code dataset (Li et al., 2024b), which includes coding problems from various platforms with visual elements in the problem descriptions. However, we find that most of these problems are unsuitable for HumanEval-V. Many images are non-essential, as they can be inferred from the textual problem descriptions. Some problems, though containing relevant visual information, are overly complex and require lengthy textual descriptions to interpret, violating our requirement for self-explanatory visual content. After careful screening, less than 5% of the viewed problems pass our standards.

We select SO for its extensive repository of realworld programming problems. To identify relevant posts, we first filter questions from 2020 that have non-negative votes and accepted answers. Then, we focus on posts that include images in the question body and code blocks in the corresponding answers, narrowing down further to those tagged with Python. After this automated filtering, we manually review the remaining posts, excluding topics related to front-end, mobile, or UI development, as these often require external frameworks and libraries that do not align with the goals of our benchmark. We also exclude posts where the images provide information in textual nature, such as code snippets, error messages, or execution outputs. Ultimately, we identified suitable questions primarily covering topics like geometry, plotting, and image processing.

To further illustrate our screening process, we present two negative examples that do not meet our standards in Figure 23: (1) The first example is a coding problem from CodeForces, where the task is to determine an optimal stacking method for a set of books with identical heights, given their thickness and width, in order to minimize the total thickness. While the provided image shows a possible stacking configuration, it lacks critical information, such as constraints on the stacking method and precise book dimensions. Moreover, the core problemsolving details are conveyed primarily through text,

1020

1022

982

983

making the image non-essential for understanding 1075 the solution. (2) The second example is a cod-1076 ing problem from GeeksForGeeks, which involves 1077 traversing a 2D matrix according to a specified 1078 pattern, starting from the top-left corner and identifying the traversal endpoint. Although the im-1080 age offers a basic representation of the matrix, the 1081 traversal pattern is too complex to be effectively 1082 captured visually and requires significant textual 1083 explanation. As a result, the textual description 1084 carries more problem-solving information than the 1085 image itself, violating our requirement for the vi-1086 sual context to be self-explanatory and serve as the 1087 primary source of information. 1088

B.2 Recreation and Diversification

1089

1090

1091

1092

1093

1094

1096

1097

1098

1100

1101

1102

1103

1104

1105

1106

1107 1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

We present three examples in Figure 24, Figure 25, and Figure 26 to demonstrate our recreation and diversification process. Each figure is divided into three parts: the original problem that meets our screening criteria (top), the recreated coding task based on the distilled ideas (middle), and the diversified variant (bottom). Below are detailed explanations of each example:

Figure 24 showcases a Stack Overflow problem where a developer needs to draw a parallelogram using four specified points. The image illustrates the connection between these points, providing the essential information needed to solve the task. Since the text merely restates the geometric properties shown in the image, we significantly reduce the textual content without losing crucial details. For recreation, we transform this into a five-pointed star problem, enriching the visual information with four examples showing different point connection patterns. The new function signature clearly defines the implementation requirements, including objectives, input parameters, and return value constraints. Instead of generating a parallelogram image, our task focuses on determining whether two specific points should be connected, simplifying the implementation while maintaining emphasis on visual reasoning. For diversification, we modify the visual pattern from a five-pointed to a six-pointed star while maintaining the same function signature.

Figure 25 presents a CodeForces problem involving polygon folding and area calculation. The image demonstrates the folding process along dashed lines, showing both initial and final states. For recreation, we simplify this into a matrix folding task where overlapping sections produce color changes. The input matrix uses two initial colors (white and light blue), which can result in three distinct outcomes after folding (white, light blue, and dark blue). Three illustrative examples clarify the folding mechanics. For diversification, we replace the color addition rule with numeric addition, requiring models to process numerical changes before and after folding. 1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

Figure 26, also from CodeForces, involves grid reduction following a specific pattern. The image effectively communicates the step-by-step transformation process. For recreation, we enhance the complexity by removing the reduction factor k as a parameter, requiring models to deduce that k = 2from the provided examples. We transform the original binary scaling operation into a statistical pooling operation (e.g., minimum value computation), demanding both OCR capabilities and advanced visual reasoning. For diversification, we increase the pooling stride from 2 to 3, requiring models to analyze larger matrices. Test cases are adjusted accordingly to maintain consistency with the modified patterns.

In addition to the three examples above, we provide further examples of how we perform diversification across specific capability aspects in Table 3.

C More details on Experimental Setup

C.1 Evaluated Models

In Table 4, we provide a detailed list of Large Multimodal Models (LMMs) used in our experiments. For each model, we specify the number of parameters and include direct links to relevant reports or Huggingface repositories for further reference.

C.2 Prompt Templates

We designed three main sets of prompts for the 1159 experiments. The first set is used for the evaluation 1160 pipelines in the main benchmarking experiments, 1161 covering scenarios such as Vision-to-Code, Vision-1162 to-Code with Chain-of-Thought (CoT), Vision-to-1163 Text, and Text-to-Code, as described in Section 3. 1164 The corresponding prompts for these scenarios are 1165 listed in Figure 27. The second set of prompts is 1166 used in the iterative refinement experiments, intro-1167 duced in Section 4. These prompts address sce-1168 narios where code or previously generated textual 1169 problem specifications are refined based on feed-1170 back from the execution environment. The rele-1171 vant prompts for this scenario are provided in Fig-1172 ure 28. The third scenario involves using GPT-40 1173 as a judge to rate the diagram descriptions gen-1174

Examples
Adjust concatenation, swapping, grouping, or stacking order; modify the direction of translation, flipping, rotation, or folding. Alter arithmetic operations, sorting order, or aggregation rules.
Reverse alternation sequences; switch increments to decrements; change the direction of spirals or zigzags; adjust layer layouts.
Change the color, size, shape, angle or position of objects.
Reverse nesting or overlap order; modify connection, intersection, or adjacency rules; adjust boundary interaction conditions.
Modify graph direction, data type in array, matrix dimensions.

Table 3: Examples of modifications applied to diversify seed tasks. Modifications of a task may span many aspects.

erated by LMMs. The prompt used is shown inFigure 29.

Our prompt design has undergone multiple 1177 rounds of optimization to address specific issues we 1178 encountered. For instance, we instruct the model 1179 to follow a markdown format code block for gen-1180 eration and avoid generating multiple code blocks 1181 to streamline post-processing and improve parsing 1182 1183 success rates. We also provide detailed instructions for prompting LMMs to generate diagram descrip-1184 tions in a structured problem specification format, 1185 including Problem Restatement, Visual Facts, and 1186 1187 Visual Patterns, ensuring a comprehensive capture and expression of the visual context. In the iter-1188 ative refinement setting, we specifically instruct 1189 LMMs not to hardcode test cases into their gener-1190 ated code to ensure that improvements stem from 1191 an enhanced understanding of the problem. Ad-1192 ditionally, for LLM-as-Judge experiments, we list 1193 clear steps for rating the LMMs' outputs, promot-1194 ing more robust and reliable rating results. 1195

C.3 Ablation on Temperature

1196

1197

1198

1199

1200

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

As described in Section 3, we set the sampling temperature to T = 0.8 for generating multiple predictions, following established practices in code generation benchmarking (Chen et al., 2021, 2022). Given that LMMs may exhibit varying performance at different temperatures, we conduct an ablation study to assess the rationale behind this choice. Specifically, we evaluated all 22 LMMs under the V2T2C w/ GPT-40 setting across a range of temperatures from 0.4 to 1.0. The results are presented in Table 5, which indicate that LMMs generally demonstrate consistent performance across these settings, with a few models showing slight variations, further validating the rationale for our chosen temperature setting.

D Deeper Analysis on HumanEval-V

1213 D.1 Co-occurrence of Capability Aspects

As illustrated in Figure 3, the diagrams in HumanEval-V encompass a wide range of capabil-



Figure 14: Performance comparison between GPT-40 and QwenCoder-32B as strong coders under the *V2T2C w/SC* setting.

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1235

ity aspects that require human-level intelligence for interpretation. Each diagram typically involves multiple capability aspects to be fully understood. To explore the relationships between these capability aspects, we conduct a co-occurrence analysis based on the aspect labels assigned by human annotators during task annotation. The results, presented in Figure 30, show a heatmap where each value represents the number of tasks in HumanEval-V that involve both corresponding aspect labels. Our analysis reveals that the diagrams in HumanEval-V exhibit a diverse distribution of capability aspects. Among them, adjacency, grid, matrix, and sequence are the most frequently occurring labels. Common co-occurrences include matrix-adjacency, grid-boundary, grid-path, and sequence-linear increment, highlighting the fundamental spatial and structural relationships embedded in these diagrams.

D.2 Comparison of Strong Coder Models

We use GPT-40 as the primary strong coder for our1236benchmarking experiments, leveraging its superior1237coding capabilities to translate problem specifica-1238tions generated by LMMs into code. This allows us1239

Models	Params	Links
Pr	oprietary L	MMs
OpenAI o1 (OpenAI, 2024b)	-	OpenAI o1
GPT-40 (0806) (OpenAI, 2024a)	-	OpenAI GPT-40
GPT-4o-mini (0718) (OpenAI, 2024a)	-	OpenAI GPT-40-mini
Claude 3.5 Sonnet (1022) (Anthropic, 2024)	-	Anthropic Claude
Gemini 1.5 Pro (002) (Google, 2024a)	-	Google Gemini 1.5 Pro
Gemini 1.5 Flash (002) (Google, 2024a)	-	Google Gemini 1.5 Flash
Open-weight LMM	s with mor	e than 70B parameters
Pixtral (Agrawal et al., 2024)	124B	mistralai/Pixtral-Large-Instruct-2411
Llama-3.2-V 90B (Google, 2024b)	88.8B	meta-llama/Llama-3.2-90B-Vision-Instruct
InternVL 2.5 78B (Chen et al., 2024c)	78.4B	OpenGVLab/InternVL2-5-78B
Owen2 VL 72B (Wang et al., 2024b)	73.4B	Qwen/Qwen2-VL-72B-Instruct
QVQ-72B-Preview (Team, 2024b)	73.4B	Qwen/QVQ-72B-Preview
Molmo-D 72B (Deitke et al., 2024)	73.3B	allenai/Molmo-72B-0924
LLaVA-OV 72B (Li et al., 2024a)	73.2B	llava-hf/llava-onevision-qwen2-72b-ov-chat-hf
Open-weight LMM	s with fewe	er than 70B parameters
Chameleon 30B (Team, 2024a)	34.3B	facebook/chameleon-30b
InternVL 2.5 26B (Chen et al., 2024c)	25.5B	OpenGVLab/InternVL2-5-26B
Pixtral 12B (Agrawal et al., 2024)	12.0B	mistralai/Pixtral-12B-2409
Llama-3.2-V 11B (Google, 2024b)	10.7B	meta-llama/Llama-3.2-11B-Vision-Instruct
Qwen2 VL 7B (Wang et al., 2024b)	8.3B	Qwen/Qwen2-VL-7B-Instruct
InternVl 2.5 8B (Chen et al., 2024c)	8.1B	OpenGVLab/InternVL2-5-8B
LLaVA-OV 7B (Li et al., 2024a)	8.03B	llava-hf/llava-onevision-qwen2-7b-ov-chat-hf
Molmo-D 7B (Deitke et al., 2024)	8.02B	allenai/Molmo-7B-D-0924
Chameleon 7B (Team, 2024a)	7.04B	facebook/chameleon-7b
Phi-3.5-V 4B (Microsoft, 2024)	4.2B	microsoft/Phi-3.5-vision-instruct
InternVL 2.5 4B (Chen et al., 2024c)	3.7B	OpenGVLab/InternVL2-5-4B
Open	-Weight Co	ode LLM
Qwen2.5 Coder 32B (Hui et al., 2024)	32.8B	Qwen/Qwen2.5-Coder-32B-Instruct

Table 4: List of LMMs with their parameter sizes and links to the official reports or Huggingface repositories.



Figure 15: Effect of task diversification on performance variation across models in the *V2T2C w/ SC* setting.

1240to focus on the evaluation of LMMs' visual under-1241standing ability in a more controllable manner. Our1242evaluation pipeline is designed to be robust, accom-1243modating different strong coders. To test the sta-1244bility of our results, we perform an ablation study1245by replacing GPT-40 with an open-weight LLM,

Qwen2.5-Coder-32B-Instruct (Hui et al., 2024), referred to as QwenCoder-32B. QwenCoder-32B demonstrates comparable coding performance to GPT-40, as evidenced by LiveCodeBench (Jain et al., 2024). This ablation allows us to explore whether switching to a different strong coder leads to deviations in our findings. 1246

1247

1248

1249

1250

1251

1252

1253

1254

1255

1256

1257

1258

1259

1260

The ablation is conducted under the *V2T2C w/ SC* setting, where QwenCoder-32B replaces GPT-40 to generate code based on the problem specifications provided by LMMs. The results, shown in Figure 14, reveal that GPT-40 and QwenCoder-32B exhibit near-perfect correlations, demonstrating the strong stability of our evaluation methodology.

D.3 The Effect of Diversified Tasks

As outlined in Section 2 and Section B.2, our task annotation pipeline includes a crucial step to create diversified versions of the seed tasks, expanding both the volume and variety of tasks in HumanEval-V. To evaluate whether these diversi-



Figure 16: Pass rates of LMMs across different task types.

Figure 17: Pass rates of LMMs l across main capability dimensions.

Models	T=0.4	T=0.6	T=0.8	T=1
Pro	oprietary	LMMs		
Claude 3.5 Sonnet	48.3	46.9	48.1	47.9
GPT-40	44.9	42.2	43.6	44.9
Gemini 1.5 Pro	41.1	39.2	39.4	41.6
Gemini 1.5 Flash	27.1	30	28.4	29.1
GPT-4o-mini	27.9	29	29.9	32.1
Ope	en-weight	t LMMs		
Pixtral 124B	35.6	39.8	34.2	37
InternVL 2.5 78B	37.2	36.1	35.9	36.2
Qwen2 VL 72B	31.1	29.1	28.7	31.5
LLaVA-OV 72B	25	22.7	24.1	22.2
Molmo-D 72B	16.9	13.2	14.4	14.5
Llama-3.2-V 90B	12.5	10.7	12.4	13.7
Pixtral 12B	21.2	23.4	23.2	21.1
InternVL 2.5 26B	20.9	21	20.7	21.8
Qwen2 VL 7B	13.3	17.2	16.6	18.5
InternVL 2.5 8B	13.6	16.6	17.3	16.3
InternVL 2.5 4B	16.9	17.6	15.2	20.5
LLaVA-OV 7B	15.1	15.2	13	15.2
Phi-3.5-V 4B	5.6	9.3	9.1	8.7
Llama-3.2-V 11B	7.5	10.2	9.6	9.9
Molmo-D 7B	9.2	10	11.2	11.2
Chameleon 7B	1	2.5	2.5	2.1
Chameleon 30B	2	0.5	3.3	2

Table 5: Ablation on LMMs' sampling temperature for the pass@3 results under the *V2T2C w/ GPT-4o* settings.

fied tasks introduce different challenges compared to the original seed tasks, we analyze the standard deviation of pass rates across the seed tasks and their diversified versions.

1266

1267

1268

1269

1270

1272

1273

1274

1275

1276

Specifically, we use the results of the 22 LMMs under the V2T2C w/ SC setting. We group tasks based on whether they are seed tasks or their variants, resulting in 100 task groups (since we have 100 seed tasks). We then calculate the standard deviation of the pass@3 results within each group, excluding groups where all tasks have a pass@3

Ms Figure 18: Pass rates of LMMs ns. across specific capability aspects.



Figure 19: Correlation Between LMM Pass Rates and Task Difficulty in Coding and Diagram Descriptions.

rate of 0. The standard deviation for each group is computed as:

$$SD_{group} = \sqrt{\frac{1}{N} \sum_{i=1}^{N} (pass@3_i - \mu_{group})^2}$$

where N is the number of tasks in the group, pass@ 3_i is the pass rate of task *i*, and μ_{group} is the mean pass@3 for the group. The resulting distribution of pass@3 standard deviations (SD) is shown in Figure 15 as box plots.

The results reveal that the 25th percentile of most models has an SD greater than 0.2, and the median SD is around 0.4. This demonstrates a notable performance gap between tasks within the same group, validating the effectiveness of our task diversification process.

1277 1278

1279

1280

1281

1283

1284

1285

1286

1287

1289

1291 1292 1293 1294 1295 1296

1298

1299

1300

1301

1303

1304

1305

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

1325

1326

1327

1328

1329

1331

1332

1333

1334

1335

1336

1337

1338

1339

1290

E More Detailed Error Analysis

E.1 Error Patterns and Taxonomy

To better understand where current LMMs fall short in solving the coding tasks in HumanEval-V, we conduct a statistical analysis examining the correlation between pass rates and three key factors: task type, general capability dimensions, and specific capability aspects (illustrated in Figure 3) required for understanding diagrams in HumanEval-V. The results are presented in Figure 16, Figure 17, and Figure 18. The pass rate in our analysis is the averaged pass@3 for the five proprietary models under the V2T2C w/ SC setting.

The analysis reveals that LMMs perform particularly poorly on tasks involving *Transformation* and *Iterative Calculation*, both achieving a pass@3 of approximately 32%. This suggests that these models struggle with understanding spatial transformations and tracking state changes over iterative steps. In terms of general capability dimensions, the difference in pass rates across various categories is minimal. Specifically, *Spatial Transformation, Topological Relations*, and *Dynamic Patterns* all yield an average pass@3 of 38%.

When examining specific capability aspects, we find that LMMs exhibit notable difficulty with diagrams involving dynamic patterns such as Spirals, Circular Arrangements, and Zigzags. Additionally, tasks requiring spatial transformations like Stacking, Translation, and Splitting, as well as mathematical operations such as Sorting and Absolute Value computations, pose significant challenges. We illustrate concrete error cases that highlight these challenges in Figure 31 to 45.

E.2 Error Analysis by Task Difficulty

We also investigate the correlation between LMM performance and task difficulty, using two key metrics. The first metric is the cyclomatic complexity (Gill and Kemerer, 1991) of the humanannotated solution code for each task, which reflects the complexity of programming logic. The second metric is the token length of the humanannotated diagram descriptions, which indicates the difficulty of understanding the diagram from a textual perspective. These two metrics represent human-perceived difficulty in both visual comprehension and programmatic reasoning.

For measuring LMM performance, we use the averaged pass@3 score of the top five proprietary models under the *V2T2C w/SC* setting, with corre-

lation results presented in Figure 19. Interestingly,1340the results suggest that LMM performance has little1341correlation with human-perceived difficulty, either1342in coding complexity or in visual description length,1343except for tasks with very high programming com-1344plexity or exceptionally long diagram descriptions.1345

1346

1347

1348

1349

1350

1351

1352

1353

1354

1355

1356

1357

1358

1359

Through a detailed case study, we find that many tasks in HumanEval-V are relatively easy for humans but remain challenging for LMMs, primarily because these models struggle to comprehend diagrams at a fundamental level. This limitation stems from their lack of basic visual perception and reasoning abilities, making it difficult to develop a precise metric that accurately captures LMMperceived difficulty in our benchmark.

The error cases in Figures 46 and 47 further illustrate this challenge. Tasks that appear trivial to humans often prove insurmountable for even the top-performing LMMs, highlighting the fundamental gap in their ability to interpret diagrams and reason visually.

F More Discussion on MMCode

MMCode (Li et al., 2024b) presents a coding 1362 dataset for evaluating LMMs' algorithmic problem-1363 solving capabilities in visual contexts, compris-1364 ing 3.5k questions crawled from competitive pro-1365 gramming platforms. However, as explained in 1366 Appendix B.1, the visual content in most coding 1367 challenges is redundant, with image information 1368 largely inferrable from textual descriptions. This 1369 redundancy is evident in MMCode's reported re-1370 sults, where performance on "language-only" in-1371 puts closely matches that of "vision + language" 1372 inputs. In contrast, HumanEval-V is specifically 1373 designed to evaluate visual understanding and rea-1374 soning capabilities rather than general coding pro-1375 ficiency. Our benchmark ensures that visual con-1376 text is integral to problem-solving. Experiments 1377 with five proprietary LMMs demonstrate a striking 1378 contrast: while providing only function signatures 1379 without diagrams or diagram descriptions results in 0% pass rates across all models, the same models 1381 achieve over 90% pass rates (Figure 11) when given 1382 human-annotated diagram descriptions. This dra-1383 matic performance difference confirms the essen-1384 tial role of visual information in HumanEval-V. Fur-1385 thermore, our difficulty analysis (Figure 5) shows 1386 that the coding tasks maintain moderate complex-1387 ity, enabling a focused assessment of visual reasoning abilities. Our evaluation pipeline also in-1389 1390troduces a two-stage code generation process, al-1391lowing LMMs with lower coding proficiency to1392generate diagram descriptions while delegating the1393coding implementation to more capable models.1394These deliberate design choices clearly distinguish1395HumanEval-V from MMCode by placing visual rea-1396soning at the forefront of evaluation.

G Other Considerations

Environmental Considerations: Our bench-1398 1399 mark's challenging tasks typically require largersized multimodal models, which raises environmen-1400 tal concerns regarding computational costs. How-1401 ever, we believe the solution lies in improving train-1402 ing efficiency rather than simply scaling up model 1403 size. Our future work will focus on developing 1404 resource-efficient training methods while maintain-1405 ing performance on our benchmark. 1406

Language Coverage: Currently, our benchmark 1407 primarily focuses on English and Python, which 1408 may appear limiting. This choice was deliberate, 1409 as these languages are most prevalent in LMMs' 1410 training data and best demonstrate their capabilities. 1411 While this focus allows for deeper analysis, we 1412 acknowledge the importance of linguistic diversity. 1413 Our annotation pipeline is language-agnostic and 1414 can be extended to other programming languages 1415 in future iterations. 1416

License and Distribution: Our benchmark con-1417 sists of manually created tasks, drawing inspiration 1418 from Stack Overflow discussions and the MMCode 1419 dataset (which includes problems from platforms 1420 like Codeforces and LeetCode). We intend to dis-1421 tribute our code and data under a research-only 1422 license Creative Commons Non-Commercial (CC 1423 BY-NC) to promote academic advancement. 1424

1425Data Privacy and Protection:We can conclu-1426sively confirm that our dataset contains no person-1427ally identifiable information. All tasks were created1428from scratch by our team, with careful attention1429to privacy considerations. We maintained strict1430protocols during the creation process to ensure no1431sensitive information was included.

1432Computing Infrastructure:Our experimental1433setup utilized a computing node equipped with 81434NVIDIA A800 GPUs, primarily for LMM infer-1435ence.1436(1 greedy decode and 6 repeated sampling), the

computational overhead remained manageable due1437to our focused dataset of 253 high-quality tasks.1438

Demographic of Annotators: The annotation 1439 team consisted of four highly qualified individuals 1440 - postgraduate and doctoral students specializing 1441 in computer science, each with over four years of 1442 Python programming experience. Their participa-1443 tion is entirely voluntary and research-motivated, 1444 with no monetary compensation involved. All anno-1445 tators explicitly consented to participate in this aca-1446 demic endeavor. While this arrangement worked 1447 well for our academic setting, we acknowledge that 1448 paid annotation might be necessary for larger-scale 1449 or commercial projects. 1450



Figure 20: Correlations between eight multimodal benchmarks, including HumanEval-V. Each subplot displays the relationship between two benchmarks, while the diagonal subplots show the performance distribution for the corresponding benchmark.



Figure 21: A comparison between diagrams covered in popular multimodal benchmarks.



Figure 22: A curated selection of diagrams representing the six task types in HumanEval-V.

Problem Description:		Image:
Shaass has n books. He wants to make a bookshelf bookshelf's dimensions to be as small as possible. The and its pages' width is equal to w_l . The thickness of books have the same parabolishes	<pre>for all his books. He wants e thickness of the <i>i</i>-th book : f each book is either 1 or 2.</pre>	the This image lacks is t_i essential information.
Shaass puts the books on the bookshelf in the following	g way. First he selects some of	F the
books and put them vertically. Then he puts the rest o vertical books. The sum of the widths of the horizont total thickness of the vertical books. A sample arrang the image. Help Shaass to find the minimum total thickness of the	f the books horizontally above al books must be no more than gement of the books is depicte vertical books that we can ach	e the h the d in hieve.
	Cutput:	
The first line of the input contains an integer n , (1 $\leq n$ the next n lines contains two integers t_i and w_i denoti	$n \leq 100$). Each of On the only of the thickness is minimum tot	line of the output print the al thickness of the vertical books
and width of the i -th book correspondingly, (1 \leq $t_i \leq$ 2,	, $1 \le w_i \le 100$). That we can	achieve.
and width of the <i>i</i> -th book correspondingly, $(1 \le t_i \le 2, \frac{Problem Description:}{1 \le t_i \le 2})$	$1 \le w_l \le 100$). That we can	mage:
and width of the <i>i</i> -th book correspondingly, $(1 \le t_l \le 2, \frac{Problem Description:}{0}$ Given a binary matrix of dimensions with R rows and C of 0), moving in the right direction. Perform the following	$1 \le w_l \le 100$). that we can columns. Start from cell(0, g operations:	achieve. mage: This image requires substantial textual explanation to understand
and width of the <i>i</i> -th book correspondingly, $(1 \le t_l \le 2, \frac{Problem Description:}{0}$ Given a binary matrix of dimensions with R rows and C of 0), moving in the right direction. Perform the following • If the value of matrix[i][j] is 0, then traverse in the the next value.	$1 \le w_l \le 100$). that we can columns. Start from cell(0, g operations: the same direction and check	achieve. mage: This image requires substantial textual explanation to understand $\rightarrow 0$ 1 $0 \rightarrow 1$
<pre>and width of the <i>i</i>-th book correspondingly, (1 ≤ t_l ≤ 2, <u>Problem Description:</u> Given a binary matrix of dimensions with R rows and C c 0), moving in the right direction. Perform the following • If the value of matrix[i][j] is 0, then traverse in the the next value. • If the value of matrix[i][j] is 1, then update matrix current direction clockwise. ie - up, right, down, or interval.</pre>	<pre>, 1 ≤ w_l ≤ 100). that we can columns. Start from cell(0, g operations: he same direction and check [i][j] to 0 and change the left directions change</pre>	This image requires substantial textual explanation to understand $0 \rightarrow 1$ $1 \ 0$ $1 \ 1 \ 0$
 and width of the <i>i</i>-th book correspondingly, (1 ≤ t_i ≤ 2, <u>Problem Description</u>: Given a binary matrix of dimensions with R rows and C of 0), moving in the right direction. Perform the following If the value of matrix[i][j] is 0, then traverse in the next value. If the value of matrix[i][j] is 1, then update matrix current direction clockwise. ie - up, right, down, or to right, down, left, and up respectively. 	<pre>1 ≤ w_l ≤ 100). that we can columns. Start from cell(0, g operations: he same direction and check [i][j] to 0 and change the left directions change</pre>	mage: This image requires substantial textual explanation to understand $\begin{array}{c} 0 & 1 \\ 1 & 0 \\ 0 & 0 \end{array}$
<pre>and width of the <i>i</i>-th book correspondingly, (1 ≤ t_l ≤ 2, <u>Problem Description:</u> Given a binary matrix of dimensions with R rows and C c 0), moving in the right direction. Perform the following • If the value of matrix[i][j] is 0, then traverse in the the next value. • If the value of matrix[i][j] is 1, then update matrix current direction clockwise. ie - up, right, down, or to right, down, left, and up respectively. Find the index of the cell where you will be forced to experiorming the given traversal.</pre>	<pre>1 ≤ w_i ≤ 100). that we can columns. Start from cell(0, g operations: he same direction and check [i][j] to 0 and change the left directions change exit the matrix while</pre>	mage: This image requires substantial textual explanation to understand $\begin{array}{c} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ \end{array}$
<pre>and width of the <i>i</i>-th book correspondingly, (1 ≤ t_l ≤ 2, <u>Problem Description:</u> Given a binary matrix of dimensions with R rows and C c Ø), moving in the right direction. Perform the following • If the value of matrix[i][j] is 0, then traverse in the the next value. • If the value of matrix[i][j] is 1, then update matrix current direction clockwise. ie - up, right, down, or to right, down, left, and up respectively. Find the index of the cell where you will be forced to experiment of the given traversal. • In the second se</pre>	<pre>1 ≤ w_i ≤ 100). that we can columns. Start from cell(0, g operations: he same direction and check [i][j] to 0 and change the left directions change exit the matrix while</pre>	mage: This image requires substantial textual explanation to understand $\begin{array}{c} & 0 & 1 \\ \hline 1 & 0 \\ \hline 1 & 0 \\ \hline 1 & 0 \\ \hline \end{array}$ $\begin{array}{c} 0 & 0 \\ \hline 1 & 0 \\ \hline \end{array}$
<pre>and width of the <i>i</i>-th book correspondingly, (1 ≤ t_i ≤ 2, <u>Problem Description:</u> Given a binary matrix of dimensions with R rows and C c 0), moving in the right direction. Perform the following • If the value of matrix[i][j] is 0, then traverse in the the next value. • If the value of matrix[i][j] is 1, then update matrix current direction clockwise. ie - up, right, down, or to right, down, left, and up respectively. Find the index of the cell where you will be forced to e performing the given traversal. <u>Input:</u> A two dimensional matrix matrix[I][] and the number of</pre>	<pre>1 ≤ w_i ≤ 100). that we can columns. Start from cell(0, g operations: he same direction and check [i][j] to 0 and change the left directions change exit the matrix while <u>Output:</u></pre>	mage: This image requires substantial textual explanation to understand $\begin{array}{c} & 0 \\ \hline \\ 1 \\ \hline \\ 1 \\ \hline \end{array}$ $\begin{array}{c} 0 \\ \hline \\ 1 \\ \hline \end{array}$ $\begin{array}{c} 0 \\ \hline \\ 1 \\ \hline \end{array}$ $\begin{array}{c} 0 \\ \hline \\ 1 \\ \hline \end{array}$

Figure 23: Two negative examples in our data screening process: the first example is sourced from CodeForces (https://codeforces.com/problemset/problem/294/B), and the second from GeeksforGeeks (https://www.geeksforgeeks.org/problems/last-cell-in-a-matrix/1).

↓ Original Task ↓



↓ Diversified Version ↓



Figure 24: Task annotation examples illustrating the recreation and diversification applied to the screened coding problem. The original problem is sourced from Stack Overflow (https://stackoverflow.com/questions/ 69163515).

↓ Original Task ↓



\downarrow Our Recreated Version \downarrow



↓ Diversified Version ↓



Figure 25: Task annotation examples illustrating the recreation and diversification applied to the screened coding problem. The original problem is sourced from CodeForces (https://codeforces.com/problemset/problem/1381/E).

↓ Original Task ↓



↓ Our Recreated Version ↓

Function Signature:	Imag	e:						
<pre>def solution(matrix: list[list[int]]) -> list[list[int]]: """</pre>	1	2			7			
Refer to the example cases illustrated in the figure, identify and implement the pooling operation on the matrix.	3	4						
Parameters:	1	3	4	6				
matrix: A 2d list representing the initial matrix.	5	3	8	7		1	4	
Returns: list[list[int]]: A 2d list representing the resulting matrix after the	6	2	9	0		2	0	
pooling operation.	8	2	5	1				
Test Cases: assert solution([[1, 3, 4, 2], [2, 1, 1, 3], [1, 2, 2, 4], [3, 2, 1, 0]]) == [[1, 2, 2, 4], [3, 2, 1, 0]])	1],	1, 0			solut	ion(•	

↓ Diversified Version ↓

Function Signature:	-								
Function Signature:	1100	ige:		1	2	6			
<pre>def solution(matrix: list[list[int]]) -> list[list[int]]: """</pre>				3	4	3	\implies	1	
Refer to the example cases illustrated in the figure, identify and implement				8	7	9			
matrix.	2	4	2	7	9	0			
Deservations	1	2	9	7	5	3			
matrix: A 2d list representing the initial matrix.	4	6	7	3	7	2		1	0
Defense	3	8	9	6	9	3	-	2	3
list[list[int]]: A 2d list representing the resulting matrix after the	4	8	5	7	5	4			
pooling operation.	8	9	2	4	9	8			
Test Cases: assert solution([[1, 3, 4, 2, 0, 3], [2, 1, 1, 3, 2, 6], [1, 2, 2, 4, 4, 7], [3, [1, 7, 5, 2, 2, 0], [2, 9, 1, 2, 3, 1]]) == [[1, 0], [1, 0]]; assert solution(2, 1	, 0 ,	1,	0],	[3	, 2,	1, 0,	1, 0],

Figure 26: Task annotation examples illustrating the recreation and diversification applied to the screened coding problem. The original problem is sourced from CodeForces (https://codeforces.com/problemset/problem/ 1996/B).

You are an exceptionally intelligent coding assistant with a deep understanding of Python programming and a keen ability to interpret visual data. Your responses are consistently accurate, reliable, and thoughtful.
<pre>**Objective:** You will be presented with a Python programming problem and an accompanying image. Please complete the function based on the provided image and code context.</pre>
<pre>**Note** - Remember, the signature by itself does not contain the entire problem; the image provides critical details Observe the image closely and determine how its visual elements correspond to the problem's inputs, outputs, operations, calculations, patterns (static/dynamic), and conditions. Place generate the complete code colution including its function signature and body. formatted in a single Putter code.</pre>
block, **without any additional text or explanation**.
(function_signature) Scenario: P _{V2C} (D, σ)
You are an exceptionally intelligent coding assistant with a deep understanding of Python programming and a keen ability to interpret visual data. Your responses are consistently accurate, reliable, and thoughtful.
Objective: You will be presented with a Python programming problem and an accompanying image. Please complete the function based on the provided image and code context.
<pre>**Note** - {problem_category_specification} - Remember, the signature by itself does not contain the entire problem; the image provides critical details Observe the image closely and determine how its visual elements correspond to the problem's inputs, outputs, operations, calculations, patterns (static/dynamic), and conditions First summarize the important clues or findings and write a step-by-step analysis Then generate the complete code solution, including the function signature and body, formatted in a single Python code block</pre>
Code Context: ```python
{function_signature} Scenario: $P_{V2C}(D, \sigma, I_{CoT})$
Instructions: You will receive a Python programming problem and an accompanying image for analysis:
<pre>**Code Context:** ```python {function_signature}</pre>
 Analyze the Function Signature <pre>Examine the provided function signature (its input, output, and goal) and identify any missing context. Remember, the signature by itself does not contain the entire problem; the image provides critical details.</pre>
Observe the image closely and determine how its visual elements correspond to the problem's inputs, outputs, operations, calculations, patterns (static/dynamic), and conditions. - First, describe the visual elements you see.
- Next, list the important facts from the image that are relevant for understanding the problem. - Finally, deduce any missing information from the problem based on the image.
Response Format: Scenario: $P_{V2T}(D, \sigma)$ Please structure your response in three main sections (use Markdown H1 headers):
 # Problem Restatement Provide a concise restatement of the problem, including relevant background and requirements. **# Visual Facts**
List the facts directly observed from the image that are necessary for interpreting or solving the problem. 3. **# Visual Patterns**
<pre>summarize any objects, operations, transformations, patterns, conditions, and relationships interred from these facts. **Important Note:** - Clearly separate facts (what you directly see in the image) from patterns (what you infer based on those facts) If complex visual information is difficult to express in plain language, use formal notation (mathematical or pseudo-code) - State only what you are sure of; do not introduce assumptions not supported by the image or give vague conclusions **Do not** include any code implementation in your response.</pre>
<pre>**Instructions:** You are an exceptionally intelligent coding assistant that consistently delivers accurate and reliable responses to user instructions. Please complete the function based on the provided problem specification, code context, and accompanying image (if provided). Return the complete solution, including the function signature, in a single response, formatted within a Python code block.</pre>
<pre>**Problem Specification:** ``markdown {problem_specification}</pre>
Code Context:
$Scenario: P_{T2C}(PS, \sigma)$

Figure 27: Prompting templates used for the four scenarios introduced in Section 3. {function_signature} and {problem_specification} serve as placeholders for the respective content.

interpret visual data. Four responses are consistently accurate, relia	rstanding of Python programming and a keen ability to ble, and thoughtful.
<pre>**Objective:** You will be presented with a Python programming problem, an accompany! previously generated. Your task is to refine both the **problem analy! feedback from the test cases.</pre>	ng image, and the problem analysis and code you is** and the **code solution** based on execution
**Code Contavt **	Scenario: Iterative V2C
<pre>```python {function_signature}</pre>	
Previous Problem Analysis and Solution: ```markdown	
{previous_prediction}	
Execution Feedback:	
{execution_feedback}	
<pre>**Note** - Remember, the signature by itself does not contain the entire proble - Observe the image closely and determine how its visual elements corr calculations, patterns (static/dynamic), and conditions Carefully review the execution feedback and analyze any errors or is - Based on the feedback, refine your understanding of the problem and previously neglected aspects from the image or problem analysis !!You must NOT directly include the test cases from the feedback in invalidates the solution.!! Instead, improve the logic to handle all p</pre>	m; the image provides critical details. respond to the problem's inputs, outputs, operations, usues that arose during testing. make necessary corrections. Ensure you revisit the your code. Doing so is considered cheating and rotential scenarios correctly.
Your task is to generate: 1. A revised version of the step-by-step problem analysis with an impr and conditions.	roved understanding of the visual details, operations,
2. A retined Python code solution, formatted in a single code block, e passes all test cases **without hardcoding specific values from the fe	nsuring that it addresses the identified issues and edback**.
You are an exceptionally intelligent coding assistant with a deep under interpret visual data. Your responses are consistently accurate, relia	rstanding of Python programming and a keen ability to ble, and thoughtful.
### Objective: You will be presented with a Python programming problem, an accompany: previously generated. Your task is to refine and generate a **new vers feedback from test cases.	ng image, and the **problem specification** you ion of the problem specification** based on execution
	Scenario: Iterative V2T2C
### Code Context:	
<pre>**# Code Context: ``python {function_signature}</pre>	
<pre>### Code Context: ```python {finction_signature} ``` #### Your Previous Version Problem Specification:</pre>	
<pre>### Code Context: ```python {function_signature} ### Your Previous Version Problem Specification: ```markdown {previous_prediction}</pre>	
<pre>### Code Context: ```python {function_signature} ### Your Previous Version Problem Specification: ```markdown {previous_prediction} ### Execution Feedback: ````````````````````````````````````</pre>	
<pre>### Code Context: ```python {function_signature} ```markdown {previous_prediction} ### Execution Feedback: ```` {execution_feedback} ````</pre>	
<pre>### Code Context: ```python {function_signature} ### Your Previous Version Problem Specification: ```markdown {previous_prediction} **** ### Execution Feedback: ```` {execution_feedback} **** ### Instruction for Refining Problem Specification:</pre>	
<pre>### Code Context: ```python {function_signature} ### Your Previous Version Problem Specification: ```markdown {previous_prediction} ### Execution Feedback: ``` {execution_feedback} *** ### Instruction for Refining Problem Specification: 1. **Analyze the Function Signature** - Examine the provided function signature (its input, output, and g - Remember, the signature by itself does not contain the entire pro- ************************************</pre>	coal) and identify any missing context. blem; the image provides critical details.
<pre>### Code Context: ```pton {function_signature} *** ### Your Previous Version Problem Specification: ```mankdown {previous_prediction} *** ### Execution Feedback: ```` {execution_feedback} *** ### Instruction for Refining Problem Specification: 1. **Analyze the Function Signature** - Examine the provided function signature (its input, output, and g - Remember, the signature by itself does not contain the entire pro 2. **Examine the Image** Observe the image closely and determine how its visual elements con calculations, patterns (static/dynamic), and conditions. - First, describe the visual elements you see.</pre>	coal) and identify any missing context. blem; the image provides critical details. respond to the problem's inputs, outputs, operations,
<pre>### Code Context: ```python {function_signature} ### Your Previous Version Problem Specification: ```markdown {previous_prediction} ### Execution Feedback: ``` {execution_feedback} *** ### Instruction for Refining Problem Specification: 1. **Analyze the Function Signature** - Examine the provided function signature (its input, output, and g - Remember, the signature by itself does not contain the entire pro 2. **Examine the Image** Observe the image closely and determine how its visual elements con calculations, patterns (static/dynamic), and conditions. - First, describe the visual elements you see. - Next, list the important facts from the image that are relevant f - Finally, deduce any missing information from the problem based or a **Eventian Enedback Analysis**</pre>	coal) and identify any missing context. blem; the image provides critical details. respond to the problem's inputs, outputs, operations, for understanding the problem.
<pre>### Code Context: ```python {function_signature} ### Your Previous Version Problem Specification: ```markdown {previous_prediction} ### Execution Feedback: ``` {execution_feedback} *** * Examine the provided function signature (its input, output, and g - Remember, the signature by itself does not contain the entire pro 2. **Examine the provided function signature (its input, output, and g - Remember, the signature by itself does not contain the entire pro 2. **Examine the Image** Observe the image closely and determine how its visual elements con calculations, patterns (static/dynamic), and conditions. - First, describe the visual elements you see. - Next, list the important facts from the image that are relevant f - Finally, deduce any missing information from the problem based or 3. **Execution Feedback Analysis** - Carefully review the execution feedback, especially error message - Analyze the issues that arose during testing, and consider how th that were previously unclear. overlooked. or incorrectly defined.</pre>	coal) and identify any missing context. blem; the image provides critical details. respond to the problem's inputs, outputs, operations, or understanding the problem. the image. es, or unexpected outputs with the expected results. ley may relate to aspects of the problem specification
<pre>### Code Context: '' python {function_signature} ### Your Previous Version Problem Specification: ''markdown {previous_prediction} ### Execution Feedback: ''' ''' ### Instruction for Refining Problem Specification: 1. **Analyze the Function Signature** - Examine the provided function signature (its input, output, and g - Remember, the signature by itself does not contain the entire pro 2. **Examine the provided function signature(its input, output, and g - Remember, the signature by itself does not contain the entire pro 2. **Examine the Image** Observe the image closely and determine how its visual elements con calculations, patterns (static/dynamic), and conditions. - First, describe the visual elements you see. Next, list the important facts from the image that are relevant f - Finally, deduce any missing information from the problem based or 3. **Execution Feedback Analysis** - Carefully review the execution feedback, especially error message - Analyze the issues that arose during testing, and consider how th that were previously unclear, overlooked, or incorrectly defined. - !!You must NOT directly include the test cases from the feedback considered cheating and invalidates the refinement process.!! - Instead, generalize your understanding to address all nossible co</pre>	coal) and identify any missing context. blem; the image provides critical details. respond to the problem's inputs, outputs, operations, or understanding the problem. I the image. s, or unexpected outputs with the expected results. ley may relate to aspects of the problem specification into the refined problem specification. Doing so is uses comprehensively.
<pre>### Code Context: '' python {function_signature} '''''''''''''''''''''''''''''''''</pre>	coal) and identify any missing context. bblem; the image provides critical details. brespond to the problem's inputs, outputs, operations, for understanding the problem. the image. is, or unexpected outputs with the expected results. ley may relate to aspects of the problem specification into the refined problem specification. Doing so is uses comprehensively. problem.
<pre>### Code Context: python {function_signature} ### Your Previous Version Problem Specification: "markdown {previous_prediction} "" ### Execution Feedback: "" texamine the provided function signature* - Ramine the provided function signature (its input, output, and g - Remember, the signature by itself does not contain the entire pro - Remember, the signature by itself does not contain the entire pro - Remember, the signature by itself does not contain the entire pro - *Examine the Image** Observe the image closely and determine how its visual elements con calculations, patterns (static/dynamic), and conditions. - First, describe the visual elements you see. Next, list the important facts from the image that are relevant f - Finally, deduce any missing information from the problem based or **Execution Feedback Analysis** - Carefully review the execution feedback, especially error message - Analyze the issues that arose during testing, and consider how the that were previously unclear, overlooked, or incorrectly defined. - !!You must NOT directly include the test cases from the feedback considered cheating and invalidates the refinement process.!! - Instead, generalize your understanding to address all possible ca - *#Refine the Problem Specification** Based on the execution feedback, revise your understanding of the p - Clarify or update any ambiguous parts of the specification. - Address missing or incorrect details in the initial problem specification. - Address missing or incorrect details in the initial problem specification. - Address missing or incorrect details in the initial problem specification. - Address missing or incorrect details in the initial problem specification. - Address missing or incorrect details in the initial problem specification.</pre>	<pre>goal) and identify any missing context. bblem; the image provides critical details. rrespond to the problem's inputs, outputs, operations, for understanding the problem. In the image. es, or unexpected outputs with the expected results. ley may relate to aspects of the problem specification into the refined problem specification. Doing so is uses comprehensively. problem. fication that were revealed by the test cases.</pre>
<pre>### Code Context: '' python {function_signature} ### Your Previous Version Problem Specification: ''markdown {previous_prediction} ### Execution Feedback: (execution_feedback) ### Instruction for Refining Problem Specification: 1. **Analyze the Function Signature** Examine the provided function signature (its input, output, and g Remeber, the signature by itself does not contain the entire provides for the image closely and determine how its visual elements contail conditions. First, describe the visual elements you see. Next, list the important facts from the image that are relevant the Finally, deduce any missing information from the problem based or 3. **Execution Feedback Analysis** Carefully review the execution feedback, especially error message Analyze the issues that arose during testing, and consider how the that were previously unclear, overlooked, or incorrectly defined. - !!You must NOT directly include the test cases from the feedback considered cheating and invalidates the refinement process.!! Instead, generalize your understanding to address all possible cast. **Refine the Problem Specification** Based on the execution feedback, revise your understanding of the g - Clarify or update any ambiguous parts of the specification. - Address missing or incorrect details in the initial problem specification. - Address missing or incorrect details in the initial problem specification. - Problem Restatement. Yisual Facts and Yisual Patterns - Problem Restatement. Yisual Facts and Yisual Patterns - Problem Restatement. Yisual Facts and Yisual</pre>	<pre>coal) and identify any missing context. bblem; the image provides critical details. respond to the problem's inputs, outputs, operations, for understanding the problem. 1 the image. 1 the image. 1 s, or unexpected outputs with the expected results. 1 ey may relate to aspects of the problem specification 1 into the refined problem specification. Doing so is 1 ses comprehensively. 1 oroblem. 1 fication that were revealed by the test cases. 1 headers):</pre>
<pre>### tode totext: '' python {function_signature} ### Your Previous Version Problem Specification: ''markdown {previous_prediction} ### Execution Feedback: ''' ''''''''''''''''''''''''</pre>	<pre>tool) and identify any missing context. bblem; the image provides critical details. rrespond to the problem's inputs, outputs, operations, for understanding the problem. the image. ts, or unexpected outputs with the expected results. ley may relate to aspects of the problem specification into the refined problem specification. Doing so is uses comprehensively. roblem. fication that were revealed by the test cases. headers):</pre>

Figure 28: Prompting templates used for the iterative benchmarking scenarios introduced in Section 4. {function_signature}, {previous_prediction}, and {execution_feedback} serve as placeholders for the respective content.

```
**[Task Definition]**
Your task is to evaluate and identify the **Root-Cause Error** in a Model-Generated Problem Specification (Gen-PS) by
comparing it to a reference Human-Written Problem Specification (Ref-PS). Specifically, you will:
1. Identify the fundamental error in the Gen-PS that prevents it from accurately reflecting the visual context or programming
requirements.
2. Attribute the error to one or more of the following **capability dimensions**: `Basic-Level Perception`, `High-Level Comprehension`, or `Contextual Interpretation`. Rate each dimension based on the severity of the error.
3. Label the specific **capability aspects** associated with the error.
4. Provide a clear explanation to justify your evaluation.
[Task Background & Key Terms]
   **Problem Specification**: A detailed analysis of an image's visual content and its alignment with a given Python function
signature to provide problem-solving instructions.
  Root-Cause Error: The most fundamental flaw in the Gen-PS that disrupts its ability to fully capture the problem's context.
- Capability Dimensions: The key areas for assessing the model's performance. Identify any conflicts between the Ref-PS and
Gen-PS in the following dimensions:

    **Basic-Level Perceptin**: Identifying basic visual elements like shapes, colors, text, numbers, layout, positions, etc.
    **High-Level Comprehension**: Understanding the objects, relationships, constraints, patterns, and operations depicted in

the visual context.
   **Contextual Interpretation**: The output contains three sections ("Problem Restatement", "Visual Facts", "Visual
Patterns"). The visual patterns are clearly described without vagueness or hallucinated details.
[Rating Scale]
Use a scale from **1** to **3**, where:
- **1**: Severe error in the specific capability dimension.
- **2**: Moderate error in the specific capability dimension.
- **3**: Nearly perfect in the specific capability dimension.
[Output Format]
Provide your evaluation in the following structure:
1. **Label**: A JSON object indicating the identified error and scores of three capability dimensions.
   json
{
     "ratings": {
         {
              "Basic-Level Perception": 1/2/3,
              "High-Level Comprehension": 1/2/3,
"Contextual Interpretation": 1/2/3
         }
    }
}__
2. **Explanation**: A concise justification for the identified root-cause error. Link it to the identified capability
dimensions and aspects.
[Inputs for Evaluation]
Reference Problem Specification (Ref-PS):
   markdown
{human_annotated_problem_specification}
Model-Generated Problem Specification (Gen-PS):
   markdown
{model_generated_problem_specification}
[Evaluation Notes]
1. Start by carefully comparing the Ref-PS and Gen-PS to identify discrepancies.
2. Determine the Root-Cause Error and assess which capability dimensions were not effectively addressed.
3. Rate each dimension from 1 to 3 based on the severity of the error.
```

Figure 29: Prompting templates used for the LLM-as-Judge rating experiment introduced in Section 5. {hu-man_annotated_problem_specification} and {model_generated_problem_specification} serve as placeholders for the respective content.



Figure 30: Analysis of the capability aspect co-occurrences in HumanEval-V tasks.

```
Function Signature:
Imaae:
                e = (-1)^2 + (0)^2 = 1
                                      def solution(nums: list[float], idx1: int, idx2: int) -> float:
                -1
                   0
                      2
                         -9
                            -7
             3
                                          Calculate the distance between two elements in a list of numbers.
                                          Parameters:
                                              nums (list[float]): A list of numbers.
                   -2 7 -3
             0
                2
       3
          -1
                            1
                                              idx1 (int): The first index position. (0-indexed)
                                              idx2 (int): The second index position. (0-indexed)
                              (-2)^2 = 4
                                          Returns:
   -4
                  -1.2 -2.5
3
                                              float: The calculated distance between the elements at the two index
                                      positions. (round two decimal places)
Example 3: distance = (-6)^2 + (-5)^2 + (-1.2)^2 + (-2.5)^2 = 68.69
                               distar
                                          ......
Test Cases:
assert solution([3,-1,0,2,-9,-7], 0, 3)==1; assert solution([2,-3,-1,-1,-2,-2,3,-3], 0, 6)==19; assert ···
Incorrect Diagram Description Generated by Claude 3.5 Sonnet:
# Problem Restatement
The problem requires implementing a function that calculates a "distance" between two elements in a list
of numbers based on their index positions. The distance calculation follows specific rules shown in the
examples, where elements between the two selected indices contribute to the total distance value. The
result should be rounded to 2 decimal places.
# Visual Facts
1. Three examples are shown with arrays of different lengths:
   - Example 1: [3, -1, 0, 2, -9, -7]
   - Example 2: [3, -1, 0, 2, -2, 7, -3, 1]
   - Example 3: [3, -4, -2, 4, -6, -5, -1.2, -2.5, 3, -2, 7, 3]
2. Each example shows distance calculations between different pairs of indices:
   - Distance annotations show formulas like (-1)^2 + (0)^2 = 1
   - Multiple distances can be calculated within the same array
   - Some calculations involve multiple terms
3. Distance values shown include:
   - Single-step distances = 1 (from (-1)^2+(0)^2)
   - Larger distances = 9 (from (-3)^2)
   - Complex distances = 68.69 (from (-6)^2 + (-5)^2 + (-1.2)^2 + (-2.5)^2)
   - Zero distance in some cases
# Visual Patterns
1. Distance Calculation Rules:
   - When elements are adjacent, distance is calculated using the differences squared
   - For elements with other numbers between them, all elements in the range contribute
   - The calculation uses the square of each element in the range between (and including) the selected
indices
                                              Incorrect: The elements shouldn't include both ends
   - Multiple terms are summed when there are multiple elements in the range
2. Mathematical Pattern:
   - Distance = \Sigma(x\_i)^2 where x\_i are the elements between and including the indexed positions
   - For adjacent elements or single elements, it simplifies to the square of the difference
   - When indices are the same, the distance is 0
3. Range Consideration:
   - The calculation includes both start and end indices
   - The order of indices doesn't matter (the distance is the same regardless of which index is first)
   - All elements between the two indices contribute to the final distance value
```

Figure 31: Example error case demonstrating LMMs' challenges with data structure manipulation and mathematical operations. The case shown for Claude 3.5 Sonnet is under the *V2T2C w/SC* setting.

```
Imaae:
                                      Function Signature:
               mce = (-1)^2 + (0)^2 = 1
                                      def solution(nums: list[float], idx1: int, idx2: int) -> float:
             3 -1 0 2 -9 -7
                                          Calculate the distance between two elements in a list of numbers.
            (-1)^2 + (0)^2 = 1
                                          Parameters:
                     dis
                        ce = (-3)^2 = 9
                                              nums (list[float]): A list of numbers.
Example 2: 3 -1 0 2 -2 7 -3 1
                                              idx1 (int): The first index position. (0-indexed)
                                              idx2 (int): The second index position. (0-indexed)
                              (-2)^2 = 4
                                          Returns:
   -4
      -2
                -5 -1.2 -2.5
                            -2
                                 3
                                             float: The calculated distance between the elements at the two index
                                      positions. (round two decimal places)
                      +(-2.5)^2 = 68.69
                  +(-1.2)^{\frac{1}{2}}
Test Cases:
assert solution([3,-1,0,2,-9,-7], 0, 3)==1; assert solution([2,-3,-1,-1,-2,-2,3,-3], 0, 6)==19; assert ···
Incorrect Diagram Description Generated by Pixtral 124B:
## Problem Restatement
The problem requires calculating a "distance" between two elements in a list of floating-point numbers.
The function `solution` accepts a list of numbers and two index positions as inputs. The goal is to
compute and return this distance, rounded to two decimal places.
## Visual Facts
1. The image displays multiple examples of numerical lists with annotations explaining the calculations
of "distance."
2. Each example calculates the distance based on selected values from the list.
3. The distance formula used in the examples is the squared sum of selected elements from the list.
4. The index positions chosen for the examples are noted to be zero-indexed.
5. The results of the distance calculations are annotated and matched with specific parts of the list in
each example.
### Example 1:
- List: [3, -1, 0, 2, -9, -7]
 - Selected indices: 1 and 2 (values: -1 and 0)
- Distance = (-1)^2 + (0)^2 = 1
### Example 2:
- List: [3, -1, 0, 2, -2, 7, -3, 1]
- Selected indices: 1 and 2 (values: -1 and 0)
   - Distance = (-1)^2 + (0)^2 = 1
- Selected indices: 4 and 5 (values: -2 and 7)
                                                    Hallucination
  - Distance = (-2)^2 + (7)^2 = 53
### Example 3:
- List: [3, -4, -2, 4, -6, -5, -1.2, -2.5, 3, -2,
- Selected indices: 1 and 2 (values: -4 and -2)
   Distance = (-4)^2 + (-2)^2 = 20
 - Selected indices: 9 and 10 (values: -2 and 7)
  - Distance = (-2)^2 + (7)^2 = 53
- Selected indices: 4 to 7 (values: -6, -5, -1.2, -1.5)
  - Distance = (-6)^2 + (-5)^2 + (-1.2)^2 + (-2.5)^2 = 68.69
- Selected indices: 11 and 11 (value: 3)
  - Distance = (3)^2 = 9
## Visual Patterns
1. The distance is calculated as the sum of the squares of selected elements from the list.
2. The selection of elements to calculate the distance depends on the provided indices.
3. The formula for distance between two indices `idx1` and `idx2` is the sum of squares of elements from
 <code>`idx1` to `idx2` (inclusive).' Incorrect: The elements shouldn't include both ends</code>

    If `idx1` equals `idx2`, the distance is simply the square of the element at that index.

5. The result should be formatted to two decimal places for precision.
The pattern is:
distance = sum(nums[i]^2 for i in range(idx1, idx2 + 1))
```

Figure 32: Example case on the same task in Figure 31, demonstrated by Pixtral 124B under the V2T2C w/ SC setting.

```
Imaae:
                                      Function Sianature:
                                      def solution(nums: list[float], idx1: int, idx2: int) -> float:
                 (-1)^{2+}(0)^{2} = 1
                -1
                   0 2 -9
             3
                                          Calculate the distance between two elements in a list of numbers.
                                          Parameters:
                                             nums (list[float]): A list of numbers.
       3 -1
             0 2 -2 7 -3
                                             idx1 (int): The first index position. (0-indexed)
                                              idx2 (int): The second index position. (0-indexed)
    e = (-4)^2 + (-2)^2 = 20
                             (-2)^2 = 4
                                          Returns:
   -4 -2 4
3
                -5
                  -1.2 -2.5
                         3 -2
                                              float: The calculated distance between the elements at the two index
                                      positions. (round two decimal places)
 Example 3:
                  +(-1.2)^2 + (-2.5)^2 = 68.69
                                          .....
Test Cases:
assert solution([3,-1,0,2,-9,-7], 0, 3)==1; assert solution([2,-3,-1,-1,-2,-2,3,-3], 0, 6)==19; assert ···
Incorrect Diagram Description Generated by Pixtral 12B:
# Problem Restatement
The problem involves calculating a "distance" between two elements in a list of numbers. This distance
is determined using specific mathematical operations defined in the image. The function `solution` takes
three arguments:- `nums`: A list of floating-point numbers.
- `idx1`: The first index position (0-indexed).
- `idx2`: The second index position (0-indexed).
The function returns the calculated distance between the elements at the specified index positions,
rounded to two decimal places.
# Visual Facts
From the image, the following visual elements and facts are observed:
1. Each example consists of a list of numbers and a distance calculation.
2. The distance is calculated using squared differences between the elements at specific indices.
3. Each example provides a sequence of numbers and corresponding distances calculated using specific
indices.
4. The calculations are shown in the image as:
   - 'Example 1: Distance between elements at indices 1 and 3.
                                                                            Hallucination
    - Example 2: Distance between elements at indices 1, 3, and 6.
    - Example 3: Distance between elements at indices 1, 3, 5, 7, and 9.
# Visual Patterns
From the visual facts, the following patterns and operations are inferred:
                                              Vague Expression: Does not specify how to calculate the distance.
1. The distance is calculated using the sum of squared differences between the elements at the specified
indices.
2. Each example shows the sum of squares for the differences:
   - Example 1: `distance = (-1)^2 + (0)^2 = 1`
   - Example 2: `distance = (-1)^2 + (0)^2 + (-3)^2 = 1 + 1 + 9 = 11`
    - Example 3: `distance = (-4)^2 + (-2)^2 = 16 + 4 = 20` and `distance = (-2)^2 = 4`
3. These distances appear to be cumulative sums of squared differences of elements at the specified
indices.
4. For multiple indices, each squared difference is summed to produce the final distance.
This understanding of the operations and calculations helps in implementing the function `solution` as
described in the provided function signature.
```

Figure 33: Example case on the same task in Figure 31, demonstrated by Pixtral 12B under the V2T2C w/ SC setting.

```
Image:
                                         Function Signature:
                                         def solution(nums: list[float], idx1: int, idx2: int) -> float:
               ance = (-1)^2 + (0)^2 = 1
     Example 1:
              3 -1
                    0
                       2 -9 -7
                                             Calculate the distance between two elements in a list of numbers.
                                             Parameters:
        dista
             (-1)^{2}+(0)
                                                 nums (list[float]): A list of numbers.
Example 2: 3 -1
              0
                 2
                    -2
                        7 -3 1
                                                 idx1 (int): The first index position. (0-indexed)
                                                 idx2 (int): The second index position. (0-indexed)
distance = (-4)^2 + (-2)^2 = 20
                                             Returns:
 3
    -4
       -2
              -6
                 -5
                    -1.2
                       -2.5
                              -2
                                 7
                                    3
           4
                           3
                                                 float: The calculated distance between the elements at the two index
                                         positions. (round two decimal places)
 Example 3: dista
             (-6)^{2}+(-5)^{2}+(-1.2)^{2}+(-2.5)^{2}=68.69
                                     = 0
                                             .....
Test Cases:
assert solution([3,-1,0,2,-9,-7], 0, 3)==1; assert solution([2,-3,-1,-1,-2,-2,3,-3], 0, 6)==19; assert ···
<u>Incorrect</u> Code Solution Generated by OpenAI o1:
**Analysis**
From the examples in the image:
                                             Incorrect: Fail to understand the aggregation rules.
- They consistently show something like "distance = (x)^2 + (y)^2."
- There is no visible square root; the result is simply the *sum of the squares* of the two chosen
elements.
- Finally, we round to two decimal places.
Hence, given a list `nums` and two indices `idx1`, `idx2`, we compute
1/
\text{distance} = \bigl(\text{nums[idx1]}\bigr)^2 + \bigl(\text{nums[idx2]}\bigr)^2
\mathbf{N}
and round it to two decimals.
```python
def solution(nums: list[float], idx1: int, idx2: int) -> float:
 Calculate the 'distance' between two elements in a list of numbers,
 defined here as the sum of their squares (no square root).
 The result is rounded to 2 decimal places.
 Parameters:
 nums (list[float]): A list of numbers.
 idx1 (int): The first index position (0-based).
 idx2 (int): The second index position (0-based).
 Returns:
 float: The calculated distance, rounded to two decimals.

 Incorrect: Fail to understand the aggregation rules.
 x = nums[idx1]
 y = nums[idx2]
 dist = x^{**2} + y^{**2}
 return round(dist, 2)
```

Figure 34: Example case on the same task in Figure 31, demonstrated by OpenAI o1 under the textitV2C w/ CoT setting.



Figure 35: Example case on the same task in Figure 31, demonstrated by QVQ 72B Preview under the textitV2C w/ CoT setting.



Figure 36: Error case highlighting LMMs' limitations in understanding dynamic patterns with alternating elements and linear increments. The case shown for Claude 3.5 Sonnet is under the *V2T2C w/ SC* setting.

```
Imaae:
 Function Signature:
 Example 2
 def solution(nums: List[int], T: int) -> List[int]:
 Example 1
 T=0: 0 2
 T=0: 5 -1
 1
 T=1: 2 0 2 0
 Returns the array at a specified time T.
 T=1: -1 5 -1 5
 Parameters:
 0 2 0 2 0 2
 1
 nums (List[int]): The initial array of integers
 1
 T=2: 5 -1 5 -1 5 -1
 T (int): The target time point
 :
 Returns:
T=3: -1 5 -1 5 -1 5 -1 5
 List[int]: The array at time T

Test Cases:
assert solution([5,-1], 1) == [-1, 5,-1, 5]; assert solution([0,2], 3) == [2,0,2,0,2,0,2,0]; assert ···
Incorrect Diagram Description Generated by Pixtral 124B:
Problem Restatement
Given an initial array of integers and a time point T, the function is required to return the state of
the array at that specific time point T. The time evolution of the array follows a specific
transformation pattern demonstrated in the image, which needs to be identified and implemented.
Visual Facts
1. **Initial States:**
 - Example 1 starts with the array \([5, -1]\) at T=0.
 - Example 2 starts with the array ([0, 2]) at T=0.
2. **Transformation at T=1:**
 - In Example 1, the array \([5, -1]\) evolves into \([-1, 5, -1, 5]\).
 - In Example 2, the array \([0, 2]\) evolves into \([2, 0, 2, 0]\).
3. **Transformation at T=2:**
 - In Example 1, the array \([-1, 5, -1, 5]\) evolves into \([5, -1, 5, -1, 5, -1]\).
 - In Example 2, the array ([2, 0, 2, 0]) evolves into ([0, 2, 0, 2, 0, 2]).
4. **Continuation Pattern:**
 - The transformations continue for each subsequent time step, producing longer arrays.
Visual Patterns
1. **Transformation Rule:**
 - At each time step, the array is transformed by duplicating the array and reversing the order of the
elements for the second half of the newly formed array.
 - Mathematically, for an array (A_t) at time (t), the array at time (t+1) is formed by
concatenating (A_t) with a reversed copy of (A_t).
2. **Repeated Growth:**
 Incorrect: Fail to understand the transformation rules.
 - The length of the array doubles at each transformation step.
 - For an initial array (A_0) of length (n), the length of the array at time (t) is (n \times n)
2^t\).
 Incorrect: Length(T) should be 2*(T+1)
3. **General Transformation Formula:**
 - If (A_t) is the array at time (t), then (A_{t+1} = A_t + text{reverse}(A_t)).
 Example: For (A_0 = [a, b]), (A_1 = [a, b, b, a]) and (A_2 = [a, b, b, a, a, b, b, a]).
4. **Binary Pattern Expansion:**
 - The expansion pattern observed is similar to creating a sequence where each element's position
follows a binary-like doubling pattern in its progression.
By carefully examining the image and understanding the transformation pattern, the problem can be solved
by implementing the transformation logic iteratively or recursively until the desired time point T is
reached.
```

Figure 37: Example case on the same task in Figure 36, demonstrated by Pixtral 124B under the V2T2C w/SC setting.



Figure 38: Example case on the same task in Figure 36, demonstrated by Pixtral 12B under the V2T2C w/ SC setting.

```
Function Signature:
Image:
 Example 2
 def solution(nums: List[int], T: int) -> List[int]:
 Example 1
 0 2

 T=0: 5 -1
 1
 2 0 2 0
 Returns the array at a specified time T.
 T=1: -1 5 -1 5
 Parameters:
 0 2 0 2 0 2
 1
 nums (List[int]): The initial array of integers
 T=2: 5 -1 5 -1 5 -1
 T (int): The target time point
 ÷
 Returns:
T=3: -1 5 -1 5 -1 5 -1 5
 List[int]: The array at time T
Test Cases:
assert solution([5,-1], 1) == [-1, 5,-1, 5]; assert solution([0,2], 3) == [2,0,2,0,2,0,2,0]; assert ···
Correct Code Solution Generated by OpenAI o1:
Analvsis
From the examples, each time step (T) corresponds to taking a "window" of length (2 \times T+1)
from an infinite repetition of the original two-element array, starting at index (T). Concretely:
1. We have an initial array of length 2, say \(\bigl[x_0,x_1\bigr]\).

 Imagine an infinite repeating sequence: \(\dots,x_0,x_1,x_0,x_1,x_0,x_1,\dots\).

3. At time (T), we skip the first (T) elements of this infinite sequence and then take the next
(2\times (T+1)) elements.
Visually, in **Example 1** with `nums=[5, -1]`:
- **\(T=0\)**: offset \(=0\), length \(=2\). We take the first 2 elements from \([5, -1, 5, -1,\dots]\),
yielding ([5, -1]).
- **\(T=1\)**: offset \(=1\), length \(=4\). We start from index 1 (\(-1\)) and take 4 elements: \([-1,
5, -1, 5]\).
- **\(T=2\)**: offset \(=2\), length \(=6\). From index 2 (which is 5 again), the next 6 elements are
\([5, -1, 5, -1, 5, -1]\). (The figure shows a truncated 5-element snippet, but logically it continues.)
Similarly for **Example 2** (`nums=[0, 2]`), the pattern toggles every time step, with lengths 2, 4, 6,
8, etc.
```python
from typing import List
def solution(nums: List[int], T: int) -> List[int]:
    .....
    Returns the array at a specified time T by viewing the original two-element
    array `nums` as an infinite repeating sequence and taking 2*(T+1) elements
    starting at index T.
    Parameters:
        nums (List[int]): The initial array of exactly two integers.
        T (int): The target time point.
    Returns:
        List[int]: The array of length 2*(T+1), starting from offset T in the
                   infinite repetition of nums.
    .....
    # The length of the resulting array
    length = 2 * (T + 1)
    # Build the result by indexing into the repeating [nums[0], nums[1]] pattern
    result = [nums[(T + i) % 2] for i in range(length)]
    return result
```

Figure 39: Example case on the same task in Figure 36, demonstrated by OpenAI o1 under the V2C w/ CoT setting.



Figure 40: Example case on the same task in Figure 36, demonstrated by QVQ 72B Preview under the V2C w/ CoT setting.

Image:	Function Signature:	
	<pre>def solution(matrix: List[List[int]]) -> List[int]:</pre>	
$1 \leftarrow 2 \leftarrow 3 \leftarrow 4 \qquad 1 \rightarrow 2 \rightarrow 3 \rightarrow 4$	Г ини Г	
$5 \rightarrow 6 \rightarrow 7 \rightarrow 8$	Given an M x N 2D matrix, traverse the matrix according to the spiral pattern	
	shown in the image	
1 + 2 + 3 + 4 $5 + 7 + 8$	Parameters:	
	- matrix (List[List[int]]): A 2D list of integers representing the matrix.	
$5 \rightarrow 6 \rightarrow 7 \qquad 8 \qquad 9 \rightarrow 10 \rightarrow 11 \qquad 12 \qquad \bullet$	Returns:	
$9 \rightarrow 10 \rightarrow 11 \rightarrow 12$ $13 \rightarrow 14 \rightarrow 15 \rightarrow 16$	- A list of integers representing the elements of the matrix in the order as	
	shown in the image	
	! ини 	
<u>Test Cases:</u>		
assert solution([[1,2]])==[1,2]; assert solution([[1,2],[3,4],[5,6],[7,8]])==[7,8,6,4,2,1,3,5]; assert ···		
Incorrect Diagram Description Generated by CLaude 3.5 Sonnet:		
# Problem Restatement		
The problem requires implement:	ing a function that traverses a matrix following specific arrow patterns	
shown in the image. The function takes a 2D matrix as input and should return a list of integers		
representing the elements in the	ne order they are visited according to the arrow patterns demonstrated in	
the image.		
# Visual Facts		
1. The image shows four differe	ent matrix examples with arrows indicating traversal paths:	
- Top left: 2x4 matrix		
- Top right: 1x4 matrix		
- Bottom left: 3x4 matrix		
- Bottom right: 4x4 matrix		
2. Each cell contains a sequen	tial number starting from 1	
3. Arrows connect adjacent cells indicating the direction of traversal		
4. The patterns show both horizontal and vertical movements		
5. All examples demonstrate con	nsistent numbering (sequential integers)	
# Visual Datterns		
1 Traversal Pattern Rules		
- Horizontal traversals alte	arnate between left-to-right and right-to-left	
- Vertical traversals are consistently unward at the nightmost column		
- When multiple rows exist	the leftmost column shows downward traversal	
- The nattern creates a sni	ral-like nath through the matrix	
2 Common Characteristics: Inc	correct: The traversal starts from the bottom-left cell	
The traversal always star	ts from the ton-left cell (1)	
- The nattern follows a snal	ke-like movement horizontally	
- Vertical movements occur	at the edges of the matrix	
- The traversal continues u	at the edges of the matrix	
3 Pattern Variations. Incorrect	t: Vertical movements may occur in any column	
- The nattern adapts to dife	Forent matrix dimensions (1v1 2v1 3v1 4v1)	
- The basic pules remain co	ncistant regardlass of matrix size	
- The number of vertical tr	ansitions depends on the matrix height	
- The number of horizontal traversals depends on the matrix width		
- The number of norizonial traversals depends on the matrix width		
- Each number is visited exactly once		
- Lach humber is visited exactly once		
- The crave sat creates a continuous path through the matrix		
- The order of visited cells follows the arrow directions strictly		
The order of visited cells	s rollows the allow directions structly	

Figure 41: Error case illustrating LMMs' difficulties in recognizing and reasoning about complex geometric arrangements, particularly spiral and circular patterns. The case shown for Claude 3.5 Sonnet is under the V2T2C w/SC setting.

<u>Image:</u>	Function Signature:	
$1 \leftarrow 2 \leftarrow 3 \leftarrow 4$ $1 \rightarrow 2 \rightarrow 3 \rightarrow 4$	<pre>def solution(matrix: List[List[int]]) -> List[int]:</pre>	
	Given an M x N 2D matrix traverse the matrix according to the sniral nattern	
$5 \div 6 \div 7 \div 8$ $1 \div 2 \leftarrow 3 \leftarrow 4$	shown in the image	
	Parameters:	
	- matrix (List[List[int]]): A 2D list of integers representing the matrix.	
$5 \rightarrow 6 \rightarrow 7 \qquad 8 \qquad 9 \rightarrow 10 \rightarrow 11 \qquad 12 \qquad 1$	Returns:	
9 + 10 + 11 + 12 13 + 14 + 15 + 16	- A list of integers representing the elements of the matrix in the order as	
	Shown in the image	
Test Cases:	<u>+</u>	
assert solution([[1,2]])==[1,2]; assert solution([[1,2],[3,4],[5,6],[7,8]])==[7,8,6,4,2,1,3,5]; assert ···		
Incorrect Diagram Description Generated by Pixtral 124B:		
# Problem Restatement		
Given a 2D matrix (a list of lists of integers), the problem asks us to traverse this matrix in a		
specific spiral pattern. The objective is to return a list of integers that represent the elements of		
	sed by this spiral pattern as shown in the image.	
# Visual Facts		
The image provides several illustrations of matrices with arrows indicating the order of traversal in a		
spiral pattern. Here are the observed facts:		
1. **Matrix Dimensions and Ele	ments**:	
- The matrices shown are 2x4, 2x4, 3x4, and 4x4 in size.		
2. **Traversal Directions**:	Incorrect: The traversal starts from the bottom-left cell	
- The traversal starts from	the top-left corner of the matrix.	
- Arrows indicate the follo	wing movements: right, down, left, and up.	
- The direction changes eac	h time the edge of the matrix or the previously traversed path is reached.	
3. **Traversal Paths**:		
- For the 2x4 matrix, the traversal first goes right across the top row, then down the last column;		
- For the 3x4 matrix the traversal includes the rightward movement across the top row downward		
movement in the last column, leftward movement across the second row, upward movement in the first		
column, and so on.		
- For the 4x4 matrix, the pattern continues similarly, with additional inner rows and columns being		
traversed in a spiral manner.	*	
# Visual Patterns	Hallucination	
Based on the visual facts obse	rved:	
1. **Spiral Traversal Algorith	m**: //	
- Start from the top-left e	lement.	
- Move right across the top	row until the last element of the row is reached.	
- Move down the last column until the bottom-most element is reached.		
- Move up the first column until the ton-most element of the next inner row is reached		
- Repeat this cycle for the inner sub-matrix until all elements are traversed.		
2. **Dynamic Boundaries**:		
- The boundaries (top, bottom, left, right) of the current sub-matrix are reduced after each complete		
traversal of a layer.		
- The right boundary decreases after a rightward traversal.		
- The portom boundary decreases after a downward traversal.		
- The top boundary increases after an upward traversal.		
3. **Edge Conditions**:		
- After traversing the oute	rmost layer, the process continues with the next inner layer, reducing the	
matrix's effective size progre	ssively.	
- The process ends when all elements have been visited, effectively reducing the matrix boundaries to		
the matrix dimensions are even).		
the matrix dimensions are even	··	

Figure 42: Example case on the same task in Figure 41, demonstrated by Pixtral 124B under the V2T2C w/SC setting.

<u>Image:</u>	Function Signature:	
	<pre>def solution(matrix: List[List[int]]) -> List[int]:</pre>	
$ \begin{array}{c} 1 + 2 + 3 + 4 \\ 5 + 6 + 7 + 8 \\ \hline 1 + 2 + 3 + 4 \\ 5 + 6 + 7 8 \\ 9 + 10 + 11 + 12 \\ \hline 1 + 2 + 3 + 4 \\ 5 - 6 + 7 8 \\ 9 + 10 + 11 + 12 \\ \hline 1 + 2 + 3 + 4 \\ 5 - 6 + 7 8 \\ 9 + 10 + 11 + 12 \\ \hline 1 + 2 + 3 + 4 \\ \hline 1 + 1 + 12 \\ \hline 1 + 1 + 1$	<pre>Given an M x N 2D matrix, traverse the matrix according to the spiral pattern shown in the image Parameters: matrix (List[List[int]]): A 2D list of integers representing the matrix. Returns: A list of integers representing the elements of the matrix in the order as shown in the image """</pre>	
<pre><u>Test Cases:</u> assert solution([[1,2]])==[1,2]; assert solution([[1,2],[3,4],[5,6],[7,8]])==[7,8,6,4,2,1,3,5]; assert ···</pre>		
Incorrect Diagram Description Generated by Pixtral 12B:		
<pre># Problem Restatement The problem involves traversing a given M x N 2D matrix according to a spiral pattern. The function `solution` takes a 2D list of integers (matrix) and returns a list of integers in the order specified by the spiral traversal as shown in the image.</pre>		
 # Visual Facts 1. The image shows a matrix divided into a grid of cells with numbers inside each cell. 2. Arrows within each cell indicate the direction of traversal. 3. Each traversal moves from one cell to the next in a spiral pattern: starting from the top-left corner and moving outward. 4. The numbers inside the cells indicate the sequence in which they will be visited. 5. There are multiple matrices of different sizes shown: - [A 2x2 matrix.] - [A 3x3 matrix.] - [A 4x4 matrix.] 5. The traversal pattern is consistent across all matrices: top-right, right-bottom, bottom-left, and left-top. 		
# Visual Patterns Incorrect: The traversal starts from the bottom-left cell in counter-clock order		
1. The traversal pattern follows a clockwise spiral starting from the top-left corner.		
2. The spiral pain alternates between norizontal and vertical movements. 3. Each laver of the spiral is traversed in a specific order:		
- Top row from left to right.		
- Rightmost column from top to bottom.		
- Bottom row from right to left.		
4. This pattern continues until all elements in the matrix have been visited.		
5. The sequence of numbers collected from the traversal will follow the spiral pattern as depicted.		
By understanding these visual facts and patterns, one can implement a function to traverse the matrix in the specified spiral order.		

Figure 43: Example case on the same task in Figure 41, demonstrated by Pixtral 12B under the V2T2C w/SC setting.

```
Image:
                                    Function Signature:
                                     def solution(matrix: List[List[int]]) -> List[int]:
 1 + 2 + 3 + 4
                    1 \rightarrow 2 \rightarrow 3 \rightarrow 4
                                         .....
                                         Given an M x N 2D matrix, traverse the matrix according to the spiral pattern
    \rightarrow 6 \rightarrow 7
           + 8
 5
                       2 + 3 + 4
                                     shown in the image
                                        Parameters:
   +2+3
             Λ
                        6
                          - 7
                    5
                                8
                                        - matrix (List[List[int]]): A 2D list of integers representing the matrix.
                            +
                                ŧ
             t
        → 7
    ♦ 6 •
             8
                    9.
                       10 + 11
                               12
 5
                                        Returns:
                                ŧ
                                         - A list of integers representing the elements of the matrix in the order as
 9 +10 +11 -
                    13 + 14 + 15 + 16
            +12
                                     shown in the image
Test Cases:
assert solution([[1,2]])==[1,2]; assert solution([[1,2],[3,4],[5,6],[7,8]])==[7,8,6,4,2,1,3,5]; assert ···
Incorrect Code Solution Generated by OpenAI o1:
### Step-by-Step Analysis
From the image, the task is to traverse a matrix in a spiral order. The order in which the elements are traversed
can be described as follows: Incorrect: The traversal starts from the bottom-left cell
1. **Start from the top-left corner of the matrix** (row 0, column 0).
2. **Move right** across the topmost row.
3. **Move down** along the rightmost column.

    **Move left** across the bottommost row (if more rows are left).

5. **Move up** along the leftmost column (if more columns are left).
6. **Repeat** this pattern, shrinking the boundary of traversal by excluding the outermost row and column until
all elements are visited.
This spiral order traversal involves the following boundaries:
- `top` (initially 0),
- `bottom` (initially the last row index),
- `left` (initially 0),
- `right` (initially the last column index).
These boundaries are adjusted after each movement:
- After moving right, increment the `top` boundary.
- After moving down, decrement the `right` boundary.
- After moving left, decrement the `bottom` boundary.
- After moving up, increment the `left` boundary.
This loop continues until all elements are visited.
### Code Implementation
```python
def solution(matrix: List[List[int]]) -> List[int]:
 if not matrix:
 return []
 result = []
 top, bottom, left, right = 0, len(matrix) - 1, 0, len(matrix[0]) - 1
 while top <= bottom and left <= right:</pre>
 for i in range(left, right + 1):
 result.append(matrix[top][i])
 top += 1
 for i in range(top, bottom + 1):
 result.append(matrix[i][right])
 right -= 1
 if top <= bottom:</pre>
 for i in range(right, left - 1, -1):
 result.append(matrix[bottom][i])
 bottom -= 1
 if left <= right:</pre>
 for i in range(bottom, top - 1, -1):
 result.append(matrix[i][left])
 left += 1
 return result
```

Figure 44: Example case on the same task in Figure 41, demonstrated by OpenAI o1 under the V2C w/ CoT setting.

```
Image:
 Function Signature:
 def solution(matrix: List[List[int]]) -> List[int]:
 1 \rightarrow 2 \rightarrow 3 \rightarrow 4
 2 🕂 3 🕇
 4
 1 •

 Given an M \times N 2D matrix, traverse the matrix according to the spiral pattern
 ► 7
 5
 →6
 8
 2 + 3 + 4
 1
 shown in the image
 Parameters:
 1 + 2 + 3 -
 7
 4
 5
 6
 8
 - matrix (List[List[int]]): A 2D list of integers representing the matrix.
 +
 +
 47
 5
 + 6
 8
 9
 10
 →11
 12
 Returns:
 +
 ŧ
 - A list of integers representing the elements of the matrix in the order as
 9 \rightarrow 10 \rightarrow 11
 12
 13 \rightarrow 14 \rightarrow 15 -
 +16
 shown in the image
Test Cases:
assert solution([[1,2]])==[1,2]; assert solution([[1,2],[3,4],[5,6],[7,8]])==[7,8,6,4,2,1,3,5]; assert ···
Unfinished Code Solution Generated by QVQ 72B Preview:
[...above is a very long CoT generation process...]
Alternatively, maybe the traversal is row by row for matrices where the number of rows is greater than or equal to
the number of columns, and spiral otherwise.
 Incorrect: The traversal pattern is irrelevant to row/column numbers
But that also seems arbitrary.
Given the confusion, perhaps the best approach is to implement a general spiral traversal algorithm, as initially
thought.
Here's a more detailed implementation:
 ```python
def solution(matrix: List[List[int]]) -> List[int]:
    if not matrix or not matrix[0]:
        return []
    rows, cols = len(matrix), len(matrix[0])
    result = []
    top, bottom = 0, rows - 1
    left, right = 0, cols - 1
    while top <= bottom and left <= right:</pre>
          for col in range(left, right + 1):
            result.append(matrix[top][col])
        top += 1
        for row in range(top, bottom + 1):
            result.append(matrix[row][right])
        right -= 1
        if top <= bottom:</pre>
             for col in range(right, left - 1, -1):
                 result.append(matrix[bottom][col])
            bottom -= 1
        if left <= right:</pre>
               for row in range(bottom, top - 1, -1):
                 result.append(matrix[row][left])
            left += 1
    return result
But given that the examples show row by row traversal, perhaps the problem expects a row by row traversal.
Alternatively, maybe the spiral traversal is supposed to alternate directions based on the current layer.
I need to make a decision here.
Given the confusion between row by row and spiral traversal, and considering that the examples show row by row
traversal, I will proceed with implementing a row by row traversal.
Here's the final code:
```python
def solution(matrix: List[List[int]]) -> List[int]:
 return [num for row in matrix for num in row]
. . .
This should satisfy the problem's requirements based on the examples provided.
[...unfinished CoT generation...]
```

Figure 45: Example case on the same task in Figure 41, demonstrated by QVQ 72B Preview under the V2C w/ CoT setting.

```
Image:
 Function Signature:
 Example 2 Operations = [1]
 Evem
 Block
 def solution(initial_blocks: List[str], operations: List[int]) -> List[str]:
 С

 в
 return the transformed blocks after applying the operations.
 Args:
 - initial_blocks: List of strings representing the initial blocks
 - operations: List of integers representing the operations to be applied
 D
 С
 Returns:
 в
 - List of strings representing the blocks after all operations
 А

Test Cases:
assert solution(["C","B","D","A"],[1,3,2])==["C","A","B","D"];assert solution(["B","C"],[1])==["C","B"];assert ···
Incorrect Diagram Description Generated by Claude 3.5 Sonnet:
Problem Restatement
The problem involves transforming a vertical arrangement of blocks (represented by letters) through a
series of operations. Each operation is specified by an index number that triggers a specific
transformation pattern. The
function needs to calculate the final arrangement of blocks after applying all specified operations in
sequence.
Visual Facts
1. Each example shows blocks arranged vertically with index positions (0 to n-1)
2. Each block contains a single letter (A, B, C, D)
3. Example 1 shows:
 - Initial state: [B, A]
 - Operations = [1]
 - Final state: [A, B]
4. Example 2 shows:
 - Initial state: [C, B, A]
 - Operations = [<mark>1</mark>]
 - Final state: [B, C, A]
5. Example 3 shows:
 - Initial state: [D, C, B, A]
 - Operations = [3, 1, 3]
 - Final state: [B, D, A, C]
6. Each transformation is shown with:
 - Dotted boxes indicating moving blocks
 - Blue arrows showing movement direction
 - White arrows showing sequence of steps
Visual Patterns
1. Operation Pattern:
 - When an index i is given in operations:
 - The block at index i moves up
 Incorrect: All blocks with index < i moves down,
 - The block at index i-1 moves down
 not just swapping i and i-1

 They effectively swap positions

2. Sequential Processing:
 - Operations are applied one at a time
 - Each operation affects exactly two adjacent blocks
 - The result of one operation becomes the input for the next
3. Movement Rules:
 - Only adjacent blocks can be swapped
 The operation index refers to the lower block in the swap
 - Blocks maintain their integrity (don't split or merge)
4. Transformation Properties:
 - Each operation preserves all blocks (no deletion/addition)
 - Multiple operations can affect the same blocks repeatedly
 - The final arrangement depends on the sequence of operations
```

Figure 46: Error case highlighting the gap between human intuition and LMM performance on seemingly straightforward visual reasoning tasks. The case shown for Claude 3.5 Sonnet is under the *V2T2C w/ SC* setting.



Figure 47: Another case highlighting the gap between human intuition and LMM performance on seemingly straightforward visual reasoning tasks. The case shown for Claude 3.5 Sonnet is under the *V2T2C w/SC* setting.