# Planning-Driven Programming: A Large Language Model Programming Workflow

**Anonymous ACL submission** 

#### Abstract

The strong performance of large language models (LLMs) raises extensive discussion on their application to code generation. Recent research suggests continuous program refinements through visible tests to improve code generation accuracy in LLMs. However, these methods suffer from LLMs' inefficiency and limited reasoning capacity. In this work, we propose an LLM programming workflow (LPW) designed to improve both initial code generation and subsequent refinements within a structured two-phase workflow. Specifically, the solution generation phase formulates a solution plan, which is then verified through visible tests to specify the intended natural language solution. Subsequently, the code implementation phase drafts an initial code according to the solution plan and its verification. If the generated code fails the visible tests, the plan verification serves as the intended solution to consistently inform the refinement process for correcting bugs. Compared to state-of-the-art methods across various existing LLMs, LPW significantly improves the Pass@1 accuracy by up to 16.4% on well-established text-to-code generation benchmarks. LPW also sets new state-ofthe-art Pass@1 accuracy, achieving 98.2% on HumanEval, 84.8% on MBPP, 59.3% on Live-Code, 62.6% on APPS, and 34.7% on Code-Contest, using GPT-40 as the backbone.

## 1 Introduction

011

017

018

019

026

027

031

034

042

Code generation, also known as *program synthesis*, studies the automatic construction of a program that satisfies a specified high-level input requirement (Gulwani et al., 2017). Recently, large language models (LLMs) pre-trained on extensive code-related datasets (Brown et al., 2020; Meta, 2024; Li et al., 2023; Roziere et al., 2023; Achiam et al., 2023; Muennighoff et al., 2023) have shown success in code-related tasks, such as code generation from natural language descriptions, also

named as text-to-code generation (Chen et al., 2021; Austin et al., 2021; Li et al., 2022), code translation (Pan et al., 2024; Yang et al., 2024), and code completion (Izadi et al., 2024). However, LLM-based code generation remains challenging due to stringent lexical, grammatical, and semantic constraints (Scholak et al., 2021).

043

045

047

049

051

054

055

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

079

083

Code generation substantially benefits from the empirical insights of human programmers. In practice, human programmers develop high-quality code by consistently rectifying errors through the analysis of test case executions (Huang et al., 2023c; Chen et al., 2023b). Different studies have refined programs based on execution results and LLM-generated information such as code and error explanation (Tang et al., 2023; Shinn et al., 2023; Madaan et al., 2023). Recent work further optimizes refinement (debugging) by performing rubber duck debugging processes (Chen et al., 2023b) and leveraging control flow graph information to assist LLMs in locating bugs (Zhong et al., 2024). However, the absence of precise correction instructions in feedback messages results in numerous refinements that deviate from the intended solution. Additionally, refining programs that significantly diverge from the problem description remains an open challenge (Tian and Chen, 2023).

To replicate each phase of program development, several studies (Lin et al., 2024; Qian et al., 2024; Dong et al., 2023b) have employed LLM instances as customized agents, assigning them diverse roles and facilitating their collaboration. Recent work incorporates extra visible tests (Huang et al., 2023a) and solution plans (Islam et al., 2024) generated by designated agents to improve the code refinements in multi-agent collaborations. However, the absence of methodologies for generating reliable visible tests and plans in these studies undermines their credibility, as incorrect visible tests and plans can lead to erroneous codes. Besides, with an increased number of agents, multi-agent collabora-



Figure 1: The pipeline of LPW, where the components highlighted in red are exclusive to LPW.

tions consume significant token resources for communication (Huang et al., 2023a). The detailed related work is discussed in Appendix A.

In this work, we propose LPW, a *large language* model programming workflow, specifically for textto-code generation, addressing the aforementioned limitations. LPW involves two phases for code generation: the solution generation phase for plan creation and plan verification, and the code implementation phase for initial code development and subsequent refinements. The pipeline of LPW is depicted in Figure 1. LPW leverages various information, including LLM-generated solution plan (Jiang et al., 2023) (block (b)), LLM-generated code explanation (Chen et al., 2023b) (block (g)), and runtime information from program execution (Zhong et al., 2024) (block (h)) to boost the code generation performance, and effectively incorporates them into an end-to-end framework. In LPW, aside from runtime information, all other messages are autonomously generated by LLMs using fewshot prompting, without the need for annotated corpora or additional training.

A unique feature of LPW is incorporating plan verification (block (c)) as the natural language intended solution for visible tests. LPW initially produces a solution plan that decomposes a complex problem into several tractable sub-problems (intermediate steps) (block (b)). LPW then verifies the solution plan against visible tests to assess its correctness, known as plan verification. For a visible test, the LLM-generated verification includes a text-based step-by-step analysis to derive the output for each intermediate step and the final output, ensuring that the final output is consistent with the visible test result. Subsequently, each inferred intermediate output is reviewed by LLMs (block (d)) to maintain logical consistency and mitigate hallucination throughout the verification. The plan verification encompasses comprehensive conditions and logical specifications for solving visible tests, eliminating potential misunderstandings before code generation. This is akin to *Test-Driven Development*, where human developers validate the intended solution with test cases (Beck, 2022).

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

159

160

161

162

163

164

165

167

LPW uses the plan and its verification in the initial code development (block (e)) to ensure that the initial code closely aligns with the problem description. Furthermore, LPW incorporates plan verification in the subsequent refinements. By comparing discrepancies between the expected output of each intermediate step, as recorded in the plan verification, against the execution trace on the failed visible test (block (h)), LPW accurately locates bugs, identifies logic flaws in the code implementation, and further generates detailed refinement suggestions, as documented in the error analysis (block (j)). Then, the error analysis when integrated with the code explanation (block (g)) serves as feedback to refine the code, surpassing conventional scalar rewards or deduced error analysis (Chen et al., 2023b; Zhong et al., 2024; Shinn et al., 2023) and thereby improving refinement efficiency and accuracy.

We first evaluate LPW on four text-to-code generation benchmarks: HumanEval (Chen et al., 2021), MBPP (Austin et al., 2021), and their extended test case variants, HumanEval-ET and MBPP-ET (Dong et al., 2023a). We conduct experiments on the proprietary LLM GPT-3.5 (Achiam et al., 2023), and open-source LLMs, Llama-3 (Meta, 2024) and Phi-3 (Abdin et al., 2024). The Pass@1 accuracy (Chen et al., 2021) is reported. Compared to the state-of-the-art LLM debugger, LDB (Zhong et al., 2024), LPW improves Pass@1 accuracy by 2% to 6.1% across all benchmarks with GPT-3.5 and achieves up to 16.4% improvement on MBPP with Llama-3. When evaluated on additional benchmarks using the advanced GPT-40 (OpenAI, 2024), LPW maintains its advantages, and achieves new state-of-the-art performance across evaluated benchmarks. Notably, on the contamination-free benchmark, LiveCode (Jain et al., 2024), and challenging benchmarks, APPS (Hendrycks et al., 2021) and CodeContests (Li

120

121

123

124

125

et al., 2022), LPW improves Pass@1 accuracy by
around 5%, 10%, and 5%, respectively, compared
to LDB. We outline our contributions as follows:

171

172

173

174

175

176

177

178

179

180

181

183

184

185

188

189

190

191

192

193

194

195

196

198

199

201

207

209

210

211

212

213

214

- We introduce an end-to-end large language model programming workflow, LPW, which significantly improves the code generation accuracy over the state-of-the-art methods.
- We derive the intended solution for visible tests, represented as the plan verification. The plan verification clarifies logic specifications required to solve the visible tests for the given problem, thereby increasing the LLMs' confidence during both the initial program generation and subsequent debugging processes.
  - We conduct extensive experiments across seven text-to-code generation benchmarks to validate the performance of LPW with various LLM backbones, provide a comprehensive analysis of their performance, token usage, and failure cases, and highlight the existing challenges.

#### **2 Problem Formulation**

We follow the problem formulation for text-to-code generation as outlined in Jiang et al. (2023), Chen et al. (2023b), and Zhong et al. (2024). The text-tocode generation problem is formulated as a triple  $\mathcal{P} = \langle Q, T_v, T_h \rangle$ , where Q represents the natural language problem specifications, and  $T_v$  and  $T_h$ are sets of visible and hidden tests, each containing input-output pairs  $(t^i, t^o) \in T = T_v \cup T_h$ . The goal is to leverage the LLM  $\mathcal{M}$  to generate a program function  $f, \mathcal{M} \to f$ , that maps each input  $t^i$ to its output  $t^o$  for all pairs in T, i.e.,  $f(t^i) = t^o$ , for  $(t^i, t^o) \in T$ . We note that  $T_h$  remains hidden during both solution generation and code implementation phases and only becomes visible if the generated f passes  $T_v$ . In LPW, for all components shown in Figure 1, the problem description Q is, by default, concatenated with task-specific prompts to produce the desired response from LLMs.

## **3** Workflow Structure

In this section, we first detail the two phases of LPW separately and then elaborate on the iterative update strategies used in each phase.

**Solution Generation**. Figure 2 displays the overall workflow of the solution generation phase in LPW (part (a)), with an example programming

	Problem Description	(1) def maximum(arr, k):
	$\overline{}$	Sorted list of length k with the maximum k numbers in arr."""
	Visible Tests	$ \begin{bmatrix} (1) \\ ($
	Solution Plan	(3) (3) (1. Sort the input array in descending order. (2. Return the first k elements of the sorted array in reverse order.
		(Dlan Varification for 1 (4)
—	Varification	1. Sort [-3, -4, 5] in descending order, which is [5, -3, -4].
		2. Return the first 3 elements in reverse order, which are [-4, -3, 5].
		[Results Compare] The correct output is [-4, -3, 5]. The analysis output is [-4, -3, 5]. [-4, -3, 5] = [-4, -3, 5]. So the plan is correct. [Correct Plan].
	Verification	[Verification Check] (5)
	Check	- Sort is correct.
		- Returning which are [-4, -5, 5] is correct.
4	Output	analysis output [-4, -3, 5] is correct.
1	Information	[Correct Plan Verification]
	(a)	(h)

Figure 2: (a): An illustrated workflow of the solution generation phase in LPW. (b): Example message fragments corresponding to each workflow component for a HumanEval problem (*120th*) with the GPT-3.5 backbone. The detailed messages are available in Section 6.

problem for illustration (part (b)). LPW leverages the self-planning approach introduced by Jiang et al. (2023) to abstract and decompose the problem description Q into a strategic and adaptable plan  $\Pi$ at the start of the solution generation phase. 215

216

217

218

219

220

221

222

223

224

225

227

228

229

230

231

232

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

For a problem in HumanEval described by block (1) in Figure 2, its example solution plan is illustrated at block (3). However, the LLM-generated plan  $\Pi$  may occasionally be incorrect, misguiding subsequent program generation. To avoid this, LPW queries the LLM to verify  $\Pi$  against all visible tests  $T_v$ . The LLM-responded plan verification  $\mathcal{A}(\Pi, T_v)$  delivers a step-by-step analysis, including all intermediate results and final derived outputs for all visible tests  $T_v$  based on  $\Pi$ . For each  $t_v \in T_v$ , its verification  $\mathcal{A}(\Pi, \{t_v\})$  compares the derived output  $t_v^{o'}$  with the ground-truth output  $t_v^{o}$ to assess the correctness of  $\Pi$ , as outlined at block 4. If  $\Pi$  is successfully verified on all visible tests, where in  $\mathcal{A}(\Pi, T_v)$ ,  $t_v^{o'} = t_v^o, \forall t_v \in T_v$ , then the plan verification  $\mathcal{A}(\Pi, T_v)$  is reviewed by the LLM again to ensure the accuracy of all intermediate results, since each intermediate step result is used in locating bugs and providing refinement suggestions when compared with the code runtime information on the failed visible test. If all intermediate outputs in  $\mathcal{A}(\Pi, T_v)$  are validated as correct by the LLM as shown at block 5,  $\mathcal{A}(\Pi, T_v)$  is treated as the intended solution for  $T_v$ . The plan  $\Pi$  and its verification  $\mathcal{A}(\Pi, T_v)$  serve as the output of the solution generation phase, guiding code development and refinements in the code implementation phase.

**Code Implementation**. Figure 3 shows the overall workflow of the code implementation phase in



Figure 3: (a): An illustrated workflow of the code implementation phase in LPW. (b): Example message fragments extending from Figure 2 and corresponding to each workflow component. See Section 6 for detailes.

249

257

259

260

261

262

263

270

271

272

273

276

277

278

279

283

LPW (part (a)), using the same problem from Figure 2 as an illustration (part (b)). LPW develops an initial program f by prompting the LLM with the problem description Q (block (1) in Figure 2), along with plan  $\Pi$  and its verification  $\mathcal{A}(\Pi, T_v)$ from the solution generation phase. Subsequently, LPW queries the LLM to add print statements for each line in f, resulting in  $f_p$ , and then executes  $f_p$  on visible tests  $T_v$ . If  $f_p$  successfully solves  $T_v$ , LPW validates it on the hidden tests  $T_h$  to report Pass@1 accuracy. Otherwise, LPW collects the runtime information on the first failed visible test  $\bar{t_v}$ , indicating that the implementation in f deviates from the specifications in  $\mathcal{A}(\Pi, \{\bar{t_v}\})$ . Blocks 1-3 in part (b) of Figure 3 depict an initial program f(block (1)) that fails on a visible test  $\bar{t_v}$  (block (2)) and its execution trace (block (3)) on  $t_v$  after adding print statements. We omit  $f_p$  from Figure 3 to keep the discussion concise. LPW instructs the LLM to conduct an error analysis by identifying inconsistencies between the intermediate outputs recorded in the execution trace of  $\bar{t_v}$  and the expected intermediate outputs documented in the verification  $\mathcal{A}(\Pi, \{\bar{t_v}\})$ , analyzing causes, and offering refinement suggestions (block (4)). Subsequently, the error analysis and code explanation for f generated by the LLM (block (5)) are concatenated as the prompt to generate the refined program f' (block (6)). The code explanation helps the LLM align the text-based error analysis with the code implementation. LPW replaces f with the refined program f' and revalidates the updated f against the visible tests  $T_v$  to assess the need for further refinements.

**Iterative Updates.** LPW includes two update steps in the solution generation phase to enable

self-correction as indicated by the red arrows in 284 Figure 2: 1) when the plan verification inferred 285 final output differs from the ground-truth output 286 for a visible test, where  $t_v^{o'} \neq t_v^o, \exists t_v \in T_v$  in  $\mathcal{A}(\Pi, T_v)$ , a revised solution plan  $\Pi'$  is included in the LLM response to substitute the original plan; 289 2) when the LLM detects any incorrect interme-290 diate values in  $\mathcal{A}(\Pi, T_v)$  e.g., contextual inconsis-291 tencies, mathematical miscalculations, or logical 292 flaws, LPW prompts the LLM to regenerate the 293 plan verification. These update methods ensure that 294 the solution plan  $\Pi$  and its verification  $\mathcal{A}(\Pi, T_n)$ 295 maintain the necessary precision, as well-formed 296  $\Pi$  and  $\mathcal{A}(\Pi, T_v)$  are essential for accurate code 297 generation (Jiang et al., 2023). In the code imple-298 mentation phase, the code refinement process acts 299 as an update mechanism, replacing the program f300 with the refined program f' when f fails the visible 301 test  $T_v$  as highlighted by the red arrow in Figure 3. 302 Overall, for a problem  $\mathcal{P}$ , LPW iteratively revises 303 the generated plan  $\Pi$  and its verification  $\mathcal{A}(\Pi, T_v)$ , 304 in the solution generation phase, until  $\mathcal{A}(\Pi, T_v)$  in-305 fers correct outputs for all visible tests  $T_v$  and no 306 error intermediate outputs are present in  $\mathcal{A}(\Pi, T_n)$ . 307 Otherwise, LPW reports a failure for  $\mathcal{P}$  when reach-308 ing the maximum iterations. Similarly, in the code 309 implementation phase, LPW iteratively refines the 310 generated program f if bugs exist. This process 311 continues until a refined f successfully solves  $T_v$ , 312 followed by Pass@1 accuracy calculation on hid-313 den tests  $T_h$ , or LPW reports a failure for  $\mathcal{P}$  upon 314 reaching the maximum iteration limit. 315

## 4 **Experiments**

Benchmarks. We first evaluate LPW on the wellestablished text-to-code benchmarks HumanEval, MBPP, HumanEval-ET, and MBPP-ET, where the given context outlines the intended functionality of the program to be synthesized. HumanEval-ET and MBPP-ET introduce approximately 100 additional hidden tests, covering numerous edge cases, for each problem in HumanEval and MBPP, thus being regarded as more reliable benchmarks (Dong et al., 2023a). In HumanEval and HumanEval-ET, we treat the test cases described in the task description as visible tests, typically 2-5 per task. For MBPP, we consider its test set that contains 500 problems with 3 hidden tests per problem. We set the first hidden test as the visible test and treat the other two as hidden, consistent with studies (Chen et al., 2023b; Zhong et al., 2024; Ni et al., 2023; Shi et al.,

316

317

318

319

320

321

322

323

324

325

327

329

330

331

		HumanE	Eval	HumanEva	al-ET	MBPI	P	MBPP-	ET
		Acc ↑	$\Delta \uparrow$	Acc ↑	$\Delta \uparrow$	Acc ↑	$\Delta \uparrow$	Acc ↑	$\Delta \uparrow$
	Baseline	$74.4 \pm 0.8$	_	$66.5 \pm 1.3$	_	$67.4 \pm 0.5$	_	$52.8 \pm 0.3$	_
	SP	$77.4 \pm 0.8$	3.0	$69.5 \pm 0.8$	3.0	$69.2 \pm 0.4$	1.8	$52.4 \pm 0.2$	-0.4
CDT 2.5	MapCoder	$77.4 \pm 1.0$	3.0	$66.5 \pm 1.0$	0.0	$72.0 \pm 0.5$	4.6	$56.6 \pm 0.5$	3.8
GP1-3.5	SD	$81.1 \pm 1.0$	6.7	$72.0 \pm 1.0$	5.5	$71.2 \pm 0.3$	3.8	$56.0 \pm 0.1$	3.2
	LDB	$82.9 \pm 1.0$	8.5	$72.6 \pm 1.0$	6.1	$72.4 \pm 0.3$	5.0	$55.6 \pm 0.2$	2.8
	LPW (ours)	<b>89.0</b> ±0.8	14.6	<b>77.4</b> ±0.8	10.9	<b>76.0</b> ±0.2	8.6	<b>57.6</b> ±0.1	4.8
	Baseline	$73.2 \pm 1.3$	-	$61.0 \pm 1.0$	-	$44.0 \pm 1.2$	-	$35.4 \pm 1.0$	-
	SP	$78.0 \pm 2.0$	4.8	$65.2 \pm 1.0$	4.2	$48.6 \pm 1.4$	4.6	$38.4 \pm 1.4$	3.0
	MapCoder	$83.5 \pm 1.3$	10.3	$73.8 \pm 0.8$	12.8	$71.4 \pm 1.0$	27.4	$55.6 \pm 1.0$	20.2
Llama-3	SD	$81.7 \pm 1.3$	8.5	$68.3 \pm 0.8$	7.3	$63.6 \pm 1.2$	19.6	$50.0{\scriptstyle~\pm1.3}$	14.6
	LDB	$84.1 \pm 1.7$	10.9	$72.0 \pm 0.8$	11.0	$57.2 \pm 1.6$	13.2	$44.8 \pm 1.4$	9.4
	LPW (ours)	$88.4 \pm 1.6$	15.2	<b>76.2</b> ±1.3	15.2	<b>73.6</b> ±1.3	29.6	<b>56.4</b> ±1.2	21.0
	Baseline	$36.0 \pm 1.0$	_	$32.3 \pm 1.0$	_	$39.0 \pm 1.3$	-	$33.2 \pm 1.4$	-
	SP	$40.8 \pm 1.4$	4.8	$34.8 \pm 0.9$	2.5	$46.4 \pm 1.4$	7.4	$37.6 \pm 1.4$	4.4
DI : 2	MapCoder	-	-	-	-	-	-	-	-
Phi-3	SD	$51.2 \pm 1.2$	15.2	$45.7 \pm 1.0$	13.4	$45.8 \pm 1.2$	6.8	$36.6 \pm 1.2$	3.4
	LDB	$65.9 \pm 1.6$	29.9	$54.9 \pm 0.9$	22.6	$52.4 \pm 1.6$	13.4	$42.8 \pm 1.4$	9.6
	LPW (ours)	<b>76.8</b> ±1.3	<b>40.8</b>	<b>62.8</b> ±1.3	30.5	<b>64.0</b> ±1.2	25.0	$48.4 \pm 1.2$	15.2

Table 1: Comparisons of Baseline, SP, MapCoder, SD, LDB, and LPW in terms of Pass@1 accuracy (Acc) and improvement ( $\Delta$ ) with respect to Baseline across benchmarks HumanEval, HumanEval-ET, MBPP, and MBPP-ET with LLMs GPT-3.5, Llama-3, and Phi-3. Acc and  $\Delta$  are measured in percentages. Best results are highlighted in red. The standard deviation ( $\pm$ ) is calculated based on three runs and applies to other experiments when reported.

2022). MBPP-ET uses the same set of problems and visible tests for each problem as MBPP.

Experimental Setup. We compare LPW with the representative code generation approaches Self-Planning (SP) (Jiang et al., 2023), MapCoder (Islam et al., 2024), Self-Debugging (+Expl) (SD) (Chen et al., 2023b), and Large Language Model Debugger (LDB) (Zhong et al., 2024). SP relies solely on the LLM-generated solution plan to produce the program solution in a single effort without refinements. MapCoder, a multi-agent collaborative approach, generates multiple unverified plans and traverses them to produce and refine code based on the current plan. SD uses a rubber duck debugging approach in LLMs, where LLMs are prompted to provide explanations of generated programs as feedback for debugging. LDB, a state-of-the-art LLM debugger, segments generated programs into blocks based on the control flow graph, which facilitates bug detection and the refinement of each program block using runtime information in LLMs. A detailed comparison of different methods is summarized in Appendix Table 15.

We generate a seed program for each problem with the same prompts and parameters introduced by Chen et al. (2023b) for SD and LDB and label the performance of seed programs as Baseline. We note that SD and LDB only perform refinements on the seed program that fails the visible tests. We experiment with various LLMs with different parameter sizes, including GPT-3.5 (turbo-0125, ≥175B), Llama-3 (70B-Instruct), and Phi-3 (14B-Instruct) to evaluate performance and demonstrate that LPW is model-independent. The Pass@1 accuracy is reported. We apply 2-shot prompting in LPW, with a maximum of 12 iterations for both the solution generation and code implementation phases. Similarly, we set the maximum debugging iterations to 12 for SD and LDB. MapCoder generates 3 plans, each with up to 4 refinement iterations. All following experiments adhere to these iteration settings. An empirical discussion on the number of iterations for LPW is available in Appendix B.

365

366

367

369

370

371

373

374

375

376

377

378

379

380

381

382

384

385

386

387

390

391

392

393

394

395

Results on various LLMs. Table 1 presents the Pass@1 accuracy for evaluated approaches, along with their respective improvements over Baseline. LPW outperforms all competing methods across all benchmarks and with various LLM backbones. Compared to LDB, LPW improves Pass@1 accuracy by 6.1%, 4.8% 3.6%, and 2%, on HumanEval, HumanEval-ET, MBPP, and MBPP-ET, respectively, with GPT-3.5 and achieves up to 16.4% improvement on MBPP with Llama-3. These results showcase the effectiveness of the proposed workflow and demonstrate the model-independent benefits of LPW. MapCoder fails on Phi-3 as it requires strict XML-formatted responses, which pose a challenge for Phi-3. The failure analysis for LPW with GPT-3.5 is available in Appendix D.

**Results on Advanced LLM**. To further demonstrate the effectiveness of LPW, we evaluate its performance on the same benchmarks presented in

361

363

364

		HumanEval	HumanEval-ET	MBPP	MBPP-ET	LiveCode	APPS	CodeContests
	Baseline	$91.5 \pm 0.3$	81.7 ±0.3	$78.4 \pm 0.4$	$62.6 \pm 0.2$	$45.7 \pm 0.6$	$41.7 {\pm} 0.9$	$28.0 \pm 0.5$
GPT-40	LDB	$92.1 \pm 0.0$	$81.7 \pm 0.0$	$82.4 \pm 0.3$	$65.4 \pm 0.0$	$54.3 \pm 0.3$	$53.2 \pm 0.3$	$29.3 \pm 0.3$
	LPW (ours)	$\textbf{98.2} \pm 0.3$	<b>84.8</b> ±0.3	$84.8 \pm 0.2$	<b>65.8</b> ±0.1	$\textbf{59.3} \pm 0.6$	<b>62.6</b> ±0.3	$34.7 \pm 0.3$

Table 2: Pass@1 accuracy, in percentages, for Baseline, LDB, and LPW on HumanEval, HumanEval-ET, MBPP, MBPP-ET, LiveCode, APPS and CodeContests benchmarks when using GPT-40 (2024-05-13) as the backbone.

		MBPP-ET↑	MBPP-ET-3↑	$\Delta \uparrow$
	MapCoder	$56.6 \pm 0.5$	$60.6 \pm 0.2$	4.0
	SD	$56.0 \pm 0.1$	$59.2 \pm 0.3$	3.2
GP1-3.5	LDB	$55.6 \pm 0.2$	$57.6 \pm 0.2$	2.0
	LPW (ours)	<b>57.6</b> ±0.1	<b>62.0</b> ±0.2	4.4

Table 3: The impact on Pass@1 accuracy with additional visible tests using the GPT-3.5 backbone. MBPP-ET-3 includes two more visible tests per problem than MBPP-ET.  $\Delta$  represents the accuracy improvement on MBPP-ET-3 over MBPP-ET. Pass@1 accuracy and  $\Delta$ are measured as percentages.

Table 1, along with the contamination-free benchmark, LiveCode, and two competitive benchmarks, APPS and CodeContests, using the advanced LLM GPT-40 as the backbone. LDB is compared due to its second-highest performance in Table 1 using 400 GPT-3.5. We sample 140 problems from LiveCode, 401 published between November 2023 and Septem-402 ber 2024, postdating GPT-4o's cutoff date. For 403 APPS and CodeContests, we use subsets of 139 404 and 150 problems, respectively. The experiment 405 results are shown in Table 2. Similarly, the perfor-406 mance of the seed programs for LDB is referred 407 to as Baseline. LPW outperforms Baseline and 408 409 LDB across all benchmarks, and establishes new state-of-the-art Pass@1 accuracy, notably achiev-410 ing 98.2% on HumanEval. GPT-40 exhibits re-411 duced performance on LiveCode, while LPW reli-412 ably outperforms LDB by 5% accuracy. For APPS 413 and CodeContests, LPW surpasses LDB by around 414 10% and 5% accuracy, highlighting the advantages 415 of LPW in tackling challenging benchmarks. Incor-416 porating the plan verification allows LPW to clarify 417 issues before code generation and efficiently cor-418 rect bugs overlooked by LLMs. In contrast, LDB 419 shows a negligible improvement of only 0.6% and 420 1.3% compared to Baseline on HumanEval and 421 CodeContests, underscoring the limitations of de-422 bugging with coarse feedback. See Appendix G for 423 tasks that LPW fails to address with GPT-40. 424

Learning from Test. We further investigate the impact of the number of visible tests on MapCoder, SD, LDB, and LPW that use visible tests to refine code. We propose a variant of MBPP-ET, denoted as MBPP-ET-3. In MBPP-ET-3, each problem's visible tests are the three hidden tests from MBPP, while the hidden tests are the extended test cases

425

426

427

428

429

430

431



Figure 4: The impact on Pass@1 accuracy with the increased number of code implementation iterations/debugging iterations on the HumanEval benchmark when leveraging GPT-3.5 as the LLM backbone. The shaded area represents the standard deviation.

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

introduced in MBPP-ET. In other words, each problem in MBPP-ET-3 contains two more visible tests than in MBPP-ET. Results in Table 3 show that LPW achieves the highest Pass@1 accuracy of 62.0% on MBPP-ET-3 and the largest improvement of 4.4% over MBPP-ET. LPW produces the step-by-step solution for each visible test to clarify code development logic and inform subsequent refinements, demonstrating superior efficiency in utilizing visible tests to improve performance among evaluated methods.

**Performance Analysis.** Figure 4 evaluates the Pass@1 accuracy of LPW when considering different numbers of code implementation iterations on the HumanEval benchmark with GPT-3.5. For SD and LDB, we allocate the same number of debugging iterations. We omit MapCoder due to its distinct refinement strategy. We note that all evaluated approaches start from iteration 0, representing the Pass@1 accuracy before debugging. In Figure 4, Baseline and SP are plotted as straight lines with 74.4% and 77.4% accuracy, respectively, due to no debugging involved. Baseline and SP serve as the control group to illustrate when debugging methods surpass no-debugging methods. SD and LDB refine incorrect programs in Baseline, surpassing SP after two iterations. LPW starts debugging from an initial 79.9% accuracy, higher than the 77.4% for SP, underscoring the importance of plan verification in initial code generation. LPW surpasses the best performance of SD and LDB after only one

		HumanEval		MBPF	)
		Acc	$\Delta$	Acc	$\Delta$
	LPW	$89.0 \pm 0.8$	_	$76.0 \pm 0.2$	-
CDT 2.5	LPW-V	$86.0 \pm 0.5$	-3.0	$73.2 \pm 0.2$	-2.8
GPT-3.5	LPW-S	$86.0 \pm 1.0$	-3.0	$73.0 \pm 0.3$	-3.0
	LPW-C	$79.9{\scriptstyle~\pm0.8}$	-9.1	$72.2 \pm 0.5$	-3.8

Table 4: Pass@1 accuracy (Acc) for different variants of LPW with GPT-3.5.  $\Delta$  denotes the decrease against LPW. Acc and  $\Delta$  are measured in percentages. See Appendix C for additional ablation study.



Figure 5: Pass@1 accuracy vs. average token cost per program for LDB and LPW on different benchmarks using GPT-40 as the LLM backbone. K is  $10^3$ . The standard deviation is too small to be visible.

iteration, demonstrating its efficient code refinement strategy. LPW gradually refines the code and reaches the highest accuracy by the 10th iteration.

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485 486

487

488

489

490

Ablation Study. Table 4 shows the Pass@1 accuracy of different variants of LPW on the HumanEval and MBPP benchmarks with GPT-3.5. The suffix -V denotes the exclusion of plan verification in both solution generation and code implementation phases; -S stands for the LPW variant that excludes the solution generation phase; while -C represents the removal of the code implementation phase, specifically omitting code refinements. For each problem, LPW-V generates the initial program based on the unverified plan and repairs the program using only code explanation and runtime information. LPW-S refines the seed program from Baseline that fails visible tests, leveraging only code explanation and runtime information. LPW-C generates the program solution based on the plan and its verification without refinements.

The performance decline of LPW-V highlights the significance of plan verification, which serves as the intended solution for visible tests, improving the performance of LLMs in both initial code generation and subsequent refinements. LPW-V considers the unverified plan when drafting initial programs. However, the effect of the unverified plan is limited, as LPW-V shows only slight improvement



Figure 6: Pass@1 accuracy as a function of token consumption for LDB, Repeated Sampling, and LPW on the APPS and CodeContests benchmarks with GPT-40. The standard deviation is omitted for better illustration. The same illustration for the LiveCode benchmark is available in Appendix Figure 8.

on MBPP compared to LPW-S, which excludes both the plan and plan verification. This aligns with the results in Table 1, where Self-Planning shows minimal improvement compared to Baseline. The results of LPW-S and LPW-C show that removing either phase in LPW decreases its performance, indicating that both solution generation and code implementation phases are crucial for optimal performance. See Appendix E and F for a discussion on the accuracy of LLM-generated plans, plan verifications and refined programs in LPW.

491

492

493

494

495

496

497

498

499

500

501

502

509

514

#### 5 **Cost-Performance Analysis**

Figure 5 compares Pass@1 accuracy against the 503 average token cost per program for LDB and LPW 504 across five benchmarks using GPT-40. When an-505 alyzing the cost for LDB, we include the tokens 506 used to generate the seed programs, which account 507 for about 2% of its total token consumption. LDB 508 consumes fewer tokens per problem but achieves lower accuracy. When measured by the accuracy-510 cost ratio, computed as Pass@1 accuracy divided 511 by the total tokens used, LDB shows better ratios 512 on HumanEval and MBPP benchmarks. On Live-513 Code, APPS, and CodeContests benchmarks, LDB and LPW exhibit similar token usage per problem, 515 while LPW displays notably higher accuracy. As a 516 result, LPW realizes higher accuracy-cost ratios of 517 0.60% on LiveCode, 0.43% on APPS, and 0.18% 518 on CodeContests per 1000 tokens, compared to 519 LDB, which reports 0.50% on LiveCode, 0.39% on 520 APPS and 0.14% on CodeContests per 1000 tokens. 521 The lower ratios of LDB arise from insufficient re-522 finements, where multiple ineffective debugging 523 iterations consume significant token resources, yet 524



Figure 7: A case study of LPW on the *120th* problem in HumanEval, extending from Figures 2 and 3, using GPT-3.5. We omit certain components in Figures 2 and 3, e.g., the plan verification check and the initial code, and present incomplete prompts and responses to save space.

the generated program remains flawed.

525

531

532

533

534

541

542

543

545

549

553

554

555

557

561

Figure 6 illustrates the variation in Pass@1 accuracy with token consumption across different approaches on the APPS and CodeContests benchmarks using GPT-40. We introduce Repeated Sampling as an enhanced Baseline. For each problem, it repeatedly samples program solutions from the LLM until either the token consumption exceeds that of LPW, or a solution passes all visible tests and is validated on hidden tests. Repeated Sampling and LDB initially improve accuracy with increased tokens but show negligible improvement after around  $10^{3.8}$  tokens on the APPS benchmark and  $10^{3.3}$  tokens on the CodeContests benchmark. In contrast, LPW starts with high token consumption for initial plan and verification generation, resulting in a sharp accuracy increase that quickly surpasses Repeated Sampling and LDB after around  $10^{4.6}$  tokens on the APPS benchmark and  $10^{4.8}$  tokens on the CodeContests benchmark. Repeated Sampling is allocated the same token budget as LPW, while its final accuracy remains lower than LPW on both benchmarks, highlighting the benefits of plan and plan verification in generating highquality initial code and subsequent refinements. Challenging benchmarks align with LPW usage scenarios, where the precise natural language solution described in the plan and plan verification is essential for logical consistency and understanding non-trivial bugs in the program, particularly when problems involve complex logical reasoning steps.

#### 6 Case Study

Figure 7 illustrates example message fragments from LPW in the *120th* problem of HumanEval using the GPT-3.5 backbone. LPW successfully generates the correct program, while all other approaches fail. This problem requires returning a

sorted array with the maximum k numbers. However, in the problem description (block (a)), the unspecified order in the output array introduces ambiguity, confusing other methods. LPW struggles at the initial solution plan (block (c)), while the issue is addressed in the [Revised Plan], during plan verification (block (d)). The visible test (block (b)) delineates the reverse order in the return array after sorting in descending order. The initial code with print statements (block (f)) fails on the visible test since the array is not reversed. Subsequently, its execution trace is compared with the plan verification (block (e)) to identify this bug, as described in the [Error Analysis] in block (g). The refined code, which first sorts the array in descending order and then reverses the first k elements into ascending order, successfully solves this problem.

562

563

564

565

566

567

568

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

590

591

592

593

594

595

596

598

## 7 Conclusion

We introduce LPW, a large language model programming workflow, designed for text-to-code generation tasks. LPW effectively integrates various advanced code generation techniques within a two-phase development model. A key innovation of LPW is the incorporation of plan verification, which enables LLMs to accurately draft an initial program and effectively correct bugs. We evaluate LPW on well-established text-to-code generation benchmarks across various LLMs. LPW significantly improves code generation accuracy compared to other existing approaches and achieves new state-of-the-art Pass@1 accuracy, with 98.2% on HumanEval, 84.8% on MBPP, 59.3% on Live-Code, 62.6% on APPS, and 34.7% on CodeContests benchmarks using GPT-40 as the backbone. In the future, additional visible tests automatically generated by LLMs (Chen et al., 2023a) can be explored to improve the performance of LPW.

## 8 Limitations

599

617

618

621

622

623

624

625

635

636

637

641

643

644

647

651

Similar to other code generation approaches, LPW is constrained by the imperfect reasoning capacity of LLMs. Strengthening the reasoning capacity of LLMs remains an ongoing challenge. While the plan and plan verification have proven valuable across different benchmarks, they require substantial token consumption. In the future, reducing 606 this consumption remains a critical area for improvement. Besides, although the LLM-generated plan and plan verification demonstrate promising accuracy on current tasks, the accuracy of the gen-610 erated code still lags behind that of the plan and plan verification (Appendix E). Incorporating al-612 ternative solution representations (Zelikman et al., 2023) alongside natural language representations 614 may assist LLMs in translating text-based solutions 615 into program solutions more accurately. 616

#### References

- Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Anthropic. 2024. The claude 3 model family: Opus, sonnet, haiku. *Anthropic AI Hub*. Accessed: 2024-07-18.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021.
  Program synthesis with large language models. arXiv preprint arXiv:2108.07732.
- Kent Beck. 2022. *Test driven development: By example.* Addison-Wesley Professional.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In Proceedings of the 34th Advances in Neural Information Processing Systems, NeurIPS, pages 1877–1901.

Angelica Chen, Jérémy Scheurer, Jon Ander Campos, Tomasz Korbak, Jun Shern Chan, Samuel R. Bowman, Kyunghyun Cho, and Ethan Perez. 2024. Learning from natural language feedback. *Transactions on Machine Learning Research*. 652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

- Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu Chen. 2023a. Codet: Code generation with generated tests. In *Proceedings of the 11th International Conference on Learning Representations*, ICLR, pages 1–19.
- 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*.
- Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. 2023b. Teaching large language models to self-debug. *arXiv preprint arXiv:2304.05128*.
- Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu, Roberta Raileanu, Xian Li, Asli Celikyilmaz, and Jason Weston. 2023. Chain-of-verification reduces hallucination in large language models. *arXiv preprint arXiv:2309.11495*.
- Yihong Dong, Jiazheng Ding, Xue Jiang, Ge Li, Zhuo Li, and Zhi Jin. 2023a. Codescore: Evaluating code generation by learning code execution. *arXiv* preprint arXiv:2301.09043.
- Yihong Dong, Xue Jiang, Zhi Jin, and Ge Li. 2023b. Self-collaboration code generation via chatgpt. *arXiv preprint arXiv:2304.07590*.
- Sumit Gulwani, Oleksandr Polozov, Rishabh Singh, et al. 2017. Program synthesis. *Foundations and Trends*® *in Programming Languages*, 4:1–119.
- Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, and Jacob Steinhardt. 2021. Measuring coding challenge competence with APPS. In *Proceedings of the 35th Advances in Neural Information Processing Systems*, NeurIPS.
- Samuel Holt, Max Ruiz Luyten, and Mihaela van der Schaar. 2024. L2MAC: Large language model automatic computer for extensive code generation. In *Proceedings of the 12th International Conference on Learning Representations*, ICLR, pages 1–61.
- Dong Huang, Qingwen Bu, Jie M Zhang, Michael Luck, and Heming Cui. 2023a. Agentcoder: Multi-agentbased code generation with iterative testing and optimisation. *arXiv preprint arXiv:2312.13010*.
- Kai Huang, Xiangxin Meng, Jian Zhang, Yang Liu, Wenjie Wang, Shuhao Li, and Yuqing Zhang. 2023b. An empirical study on fine-tuning large language

705

- 752 753 754 755
- 757

models of code for automated program repair. In Proceedings of the 38th IEEE/ACM International Conference on Automated Software Engineering, ASE, pages 1162–1174.

- Kai Huang, Zhengzi Xu, Su Yang, Hongyu Sun, Xuejun Li, Zheng Yan, and Yuqing Zhang. 2023c. A survey on automated program repair techniques. arXiv preprint arXiv:2303.18184.
- Md. Ashraful Islam, Mohammed Eunus Ali, and Md Rizwan Parvez. 2024. MapCoder: Multi-agent code generation for competitive problem solving. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics, ACL, pages 4912-4944.
- Maliheh Izadi, Jonathan Katzy, Tim Van Dam, Marc Otten, Razvan Mihai Popescu, and Arie Van Deursen. 2024. Language models for code completion: A practical evaluation. In Proceedings of the 46th IEEE/ACM International Conference on Software Engineering, ICSE, pages 1–13.
- 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.
- Xue Jiang, Yihong Dong, Lecheng Wang, Fang Zheng, Qiwei Shang, Ge Li, Zhi Jin, and Wenpin Jiao. 2023. Self-planning code generation with large language models. ACM Transactions on Software Engineering and Methodology, pages 1-28.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. In Proceedings of the 36th Advances in neural information processing systems, NeurIPS, pages 22199–22213.
- Hung Le, Yue Wang, Akhilesh Deepak Gotmare, Silvio Savarese, and Steven Chu Hong Hoi. 2022. Coderl: Mastering code generation through pretrained models and deep reinforcement learning. In Proceedings of the 36th Advances in Neural Information Processing Systems, NeurIPS, pages 21314–21328.
- Chao Lei, Nir Lipovetzky, and Krista A. Ehinger. 2023. Novelty and lifted helpful actions in generalized planning. In Proceedings of the 16th International Symposium on Combinatorial Search, SoCS, pages 148-152.
- Chao Lei, Nir Lipovetzky, and Krista A Ehinger. 2024. Generalized planning for the abstraction and reasoning corpus. In Proceedings of the 38th AAAI Conference on Artificial Intelligence, AAAI, pages 20168-20175.
- Raymond Li, Loubna Ben allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia LI, Jenny Chim, Qian Liu, Evgenii Zheltonozhskii, Terry Yue Zhuo,

Thomas Wang, Olivier Dehaene, Joel Lamy-Poirier, Joao Monteiro, Nicolas Gontier, Ming-Ho Yee, Logesh Kumar Umapathi, Jian Zhu, Ben Lipkin, Muhtasham Oblokulov, Zhiruo Wang, Rudra Murthy, Jason T Stillerman, Siva Sankalp Patel, Dmitry Abulkhanov, Marco Zocca, Manan Dey, Zhihan Zhang, Urvashi Bhattacharyya, Wenhao Yu, Sasha Luccioni, Paulo Villegas, Fedor Zhdanov, Tony Lee, Nadav Timor, Jennifer Ding, Claire S Schlesinger, Hailey Schoelkopf, Jan Ebert, Tri Dao, Mayank Mishra, Alex Gu, Carolyn Jane Anderson, Brendan Dolan-Gavitt, Danish Contractor, Siva Reddy, Daniel Fried, Dzmitry Bahdanau, Yacine Jernite, Carlos Muñoz Ferrandis, Sean Hughes, Thomas Wolf, Arjun Guha, Leandro Von Werra, and Harm de Vries. 2023. Starcoder: may the source be with you! Transactions on Machine Learning Research.

761

762

763

764

765

768

769

770

771

772

773

774

775

776

778

779

780

781

782

783

784

787

788

789

790

792

793

794

795

796

797

798

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

- Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. 2022. Competition-level code generation with alphacode. Science, 378:1092-1097.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2024. Let's verify step by step. In Proceedings of the 12th International Conference on Learning Representations, ICLR, pages 1–24.
- Feng Lin, Dong Jae Kim, et al. 2024. When Ilm-based code generation meets the software development process. arXiv preprint arXiv:2403.15852.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. 2023. Self-refine: Iterative refinement with self-feedback. In Proceedings of the 37th Advances in Neural Information Processing Systems, NeurIPS, pages 46534-46594.
- AI Meta. 2024. Introducing meta llama 3: The most capable openly available llm to date. Meta AI. Accessed: 2024-07-18.
- Niklas Muennighoff, Qian Liu, Armel Zebaze, Qinkai Zheng, Binyuan Hui, Terry Yue Zhuo, Swayam Singh, Xiangru Tang, Leandro Von Werra, and Shayne Longpre. 2023. Octopack: Instruction tuning code large language models. arXiv preprint arXiv:2308.07124.
- Ansong Ni, Srini Iyer, Dragomir Radev, Veselin Stoyanov, Wen-tau Yih, Sida Wang, and Xi Victoria Lin. 2023. Lever: Learning to verify language-to-code generation with execution. In Proceedings of the 40th International Conference on Machine Learning, ICML, pages 26106–26128.
- OpenAI. 2024. Hello gpt-4o. OpenAI. Accessed: 2024-07-18.
- Rangeet Pan, Ali Reza Ibrahimzada, Rahul Krishna, Divya Sankar, Lambert Pouguem Wassi, Michele

909

910

911

873

874

Merler, Boris Sobolev, Raju Pavuluri, Saurabh Sinha, and Reyhaneh Jabbarvand. 2024. Lost in translation: A study of bugs introduced by large language models while translating code. In *Proceedings of the 46th IEEE/ACM International Conference on Software Engineering*, ICSE, pages 1–13.

818

819

822

824

830

835

836

839

840

841

846

847

849

850

851

852

855

856

857 858

861

862

863

864

868

870

871

872

- Chen Qian, Xin Cong, Cheng Yang, Weize Chen, Yusheng Su, Juyuan Xu, Zhiyuan Liu, and Maosong Sun. 2024. Chatdev: Communicative agents for software development. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*, ACL, pages 15174–15186.
- Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. 2023. Code llama: Open foundation models for code. arXiv preprint arXiv:2308.12950.
- Torsten Scholak, Nathan Schucher, and Dzmitry Bahdanau. 2021. Picard: Parsing incrementally for constrained auto-regressive decoding from language models. *arXiv preprint arXiv:2109.05093*.
- Javier Segovia-Aguas, Sergio Jiménez, and Anders Jonsson. 2024. Generalized planning as heuristic search:
   A new planning search-space that leverages pointers over objects. Artificial Intelligence, 330:104097.
- Freda Shi, Daniel Fried, Marjan Ghazvininejad, Luke Zettlemoyer, and Sida I. Wang. 2022. Natural language to code translation with execution. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP, pages 3533– 3546.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: language agents with verbal reinforcement learning. In Proceedings of the 37th Advances in Neural Information Processing Systems, NeurIPS, pages 8634–8652.
- Zilu Tang, Mayank Agarwal, Alexander Shypula, Bailin Wang, Derry Wijaya, Jie Chen, and Yoon Kim. 2023.
  Explain-then-translate: an analysis on improving program translation with self-generated explanations. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1741–1788.
- Zhao Tian and Junjie Chen. 2023. Test-case-driven programming understanding in large language models for better code generation. *arXiv preprint arXiv:2309.16120*.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971.
- 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. In *Proceedings of the 36th Advances in Neural Information Processing Systems*, NeurIPS, pages 24824–24837.

- Zeqiu Wu, Yushi Hu, Weijia Shi, Nouha Dziri, Alane Suhr, Prithviraj Ammanabrolu, Noah A Smith, Mari Ostendorf, and Hannaneh Hajishirzi. 2023. Finegrained human feedback gives better rewards for language model training. In *Proceedings of the 37th Advances in Neural Information Processing Systems*, NeurIPS, pages 1–26.
- Zhen Yang, Fang Liu, Zhongxing Yu, Jacky Wai Keung, Jia Li, Shuo Liu, Yifan Hong, Xiaoxue Ma, Zhi Jin, and Ge Li. 2024. Exploring and unleashing the power of large language models in automated code translation. In *Proceedings of the 2024 ACM International Conference on the Foundations of Software Engineering*, FSE, pages 1–23.
- Michihiro Yasunaga and Percy Liang. 2021. Break-itfix-it: Unsupervised learning for program repair. In *Proceedings of the 38th International conference on machine learning*, ICML, pages 11941–11952.
- Eric Zelikman, Qian Huang, Gabriel Poesia, Noah Goodman, and Nick Haber. 2023. Parsel: Algorithmic reasoning with language models by composing decompositions. In *Proceedings of the 37th Advances in Neural Information Processing Systems*, NeurIPS, pages 31466–31523.
- Li Zhong, Zilong Wang, and Jingbo Shang. 2024. Ldb: A large language model debugger via verifying runtime execution step-by-step. *arXiv preprint arXiv:2402.16906*.
- Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc V Le, and Ed H. Chi. 2023. Least-to-most prompting enables complex reasoning in large language models. In *Proceedings* of the 11th International Conference on Learning Representations, ICLR, pages 1–61.

#### 912 Appendix

913

914

941

942

943

945

946

951

953

955

957

959

#### A Related Work

#### A.1 Program Synthesis

Program synthesis remains an open challenge of 915 generating a program within a target domain-916 specific language (DSL) from given specifications. 917 One prevalent approach involves searching the 918 large space of possible programs. For example, gen-919 eralized planning whose solution is formalized as 920 a planning program with pointers (Segovia-Aguas et al., 2024; Lei et al., 2023) has demonstrated 922 promising results in synthesizing program solutions for abstract visual reasoning tasks (Lei et al., 924 2024) when the DSL is carefully designed. How-925 ever, hand-crafted DSLs often suffer from limited 926 generalization capacity, and the huge search space diminishes its effectiveness. Recently, large language models trained on vast corpora have excelled in natural language processing (NLP) tasks and have been extended to code generation e.g., GPT-931 series (Achiam et al., 2023; OpenAI, 2024), Llama-932 series (Meta, 2024; Roziere et al., 2023; Touvron et al., 2023), and Claude-series (Anthropic, 2024). 934 LPW leverages the strengths of LLMs in NLP tasks to generate intended solutions in natural lan-936 937 guage. These text-based solutions demonstrate high-quality logical reasoning steps and satisfactory accuracy, thereby effectively aiding subsequent code generation.

## A.2 Prompting Techniques

To imitate the logical chain in human brain when tackling reasoning tasks, prompting methods direct LLMs to decompose problems into solvable sub-problems (Jiang et al., 2023; Zhou et al., 2023; Lightman et al., 2024; Dhuliawala et al., 2023) and progressively infer the correct answer with intermediate outputs, as exemplified by chain-ofthought prompting (Wei et al., 2022; Kojima et al., 2022). Inspired by these studies, LPW decomposes a text-to-code problem into several sub-problems described by the solution plan and follows the chain-of-thought prompting idea to verify the solution plan against visible tests with step-by-step analysis. The generated plan and its verification provide step-by-step natural language instructions for code generation, supporting LLMs in both the initial code development and subsequent refinements.



Figure 8: Pass@1 accuracy as a function of token consumption for LDB, Repeated Sampling, and LPW on the LiveCode benchmark with GPT-40.

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

989

990

991

992

993

994

#### A.3 Code Refinement

Accurate program solutions often require iterative refinements due to model limitations (Zhong et al., 2024; Chen et al., 2023b; Shinn et al., 2023). Various interactive approaches have been proposed to optimize debugging performance in LLMs, such as human feedback (Chen et al., 2024; Le et al., 2022; Wu et al., 2023), trained models (Huang et al., 2023b; Le et al., 2022; Yasunaga and Liang, 2021), LLM-generated explanations (Chen et al., 2023b; Madaan et al., 2023; Shinn et al., 2023; Tang et al., 2023), execution results (Zhong et al., 2024; Holt et al., 2024; Tian and Chen, 2023), and multi-agent collaborations (Islam et al., 2024; Huang et al., 2023a; Qian et al., 2024; Dong et al., 2023b). Current state-of-the-art LLM debuggers, such as Self-Debugging and LDB, repair various seed programs to create program solutions. However, they encounter difficulties when the initial code substantially deviates from the original intent. Besides, without safeguarding, the refined code frequently diverges from the problem specifications. In contrast, LPW develops initial code that adheres to the validated intended solution through plan verification, minimizing deviations from the problem description. The plan verification further guides the code refinement, ensuring alignment with the problem specifications.

## **B** Parameter Study

We conduct an experiment involving 20 iterations for both the solution generation and code implementation phases in LPW. Figure 9 shows the variation in Pass@1 accuracy on the HumanEval benchmark using GPT-3.5. The increased number of iterations in the solution generation phase



Figure 9: Pass@1 accuracy of LPW on the HumanEval benchmark using GPT-3.5 with 20 iterations in both the solution generation and code implementation phases.

		HumanEval		MBPP	
		Acc	$\Delta$	Acc	$\Delta$
CDT 2 5	LPW	89.0	_	76.0	_
GP 1-5.5	LPW-E	87.8	-1.2	75.6	-0.4

Table 5: Pass@1 accuracy (Acc) for the variant of LPW with the GPT-3.5 backbone. The suffix -E stands for the exclusion of code explanation in the code implementation phase. Other metrics remain consistent with those in Table 4.

results in higher initial program accuracy, 81.1%, compared to 79.9% with 12 iterations before refinements. Subsequently, the accuracy steadily increases, reaching the highest value of 89.6% after 18 debugging turns, compared to a maximum of 89.0% with 12 iterations. Overall, a larger number of iterations improves performance in both initial code generation and subsequent refinements. However, the significant token consumption presents challenges for practical applications.

C Additional Ablation Study

995

999

1001

1003

1004

1005

1006

1008

1009

1010

1012

1013

1014

1016

Table 5 shows the performance of the variant of LPW on the HumanEval and MBPP benchmarks using GPT-3.5 as the LLM backbone. The suffix -E denotes removing the code explanation when generating the refined program in the code implementation phase. The code explanation facilitates LLMs in aligning text-based error analysis with code implementation when locating and refining incorrect program lines. LPW-E demonstrates a decrease in Pass@1 accuracy on both the HumanEval and MBPP benchmarks.



Figure 10: Pass@1 accuracy of Baseline, Self-Planning (SP), MapCoder, Self-Debugging (+Expl) (SD), LDB, and LPW across different difficulty levels, *Easy*, *Medium*, and *Hard* on the HumanEval benchmark when leveraging GPT-3.5 as the LLM backbone.

## D Analysis of Unsolved Problems for LPW using GPT-3.5

### D.1 Performance Across Different Difficulty Levels

1017

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

Figure 10 compares the Pass@1 accuracy of competing methods across different difficulty levels, *Easy, Medium*, and *Hard*, on the HumanEval benchmark using GPT-3.5. We apply the method described in Zhong et al. (2024) to generate the difficulty annotations in Figure 10 by querying GPT-3.5 with problem descriptions and canonical solutions. LPW displays convincing performance, exceeding 85% accuracy across all difficulty levels. For the *Hard* level, LPW achieves 85.7% accuracy, while competing approaches show a notable degradation, falling below 70%.

#### **D.2** Failure Reasons

LPW achieves state-of-the-art performance among the evaluated methods and demonstrates superior-1035 ity over other LLM debuggers. We categorize the 1036 failure reasons for LPW on HumanEval with GPT-1037 3.5 into 4 types. Table 6 compares the percentage 1038 of different failure reasons out of the total unsolved problems for LPW based on authors' manual re-1040 view. In LPW, half of the errors result from the 1041 *No Code* type, where the generated solution plan 1042 fails to be verified on the visible tests, or the re-1043 sulting verification includes incorrect intermediate 1044 outputs in the solution generation phase, leading to 1045 failure after reaching the maximum iteration thresh-1046 old. The second most common reason is Missing Conditions, originating from the same issues in the 1048

	Missing Conditions	Differ from Intended Solution	No Code	Others
LPW	33.3	5.6	50.0	11.1

Table 6: The percentage of different failure reasons for LPW on the HumanEval benchmark with GPT-3.5 as the backbone. *Missing Conditions* arises from the same issues in the plan and plan verification. *Differ from Intended Solution* indicates the plan and plan verification are manually classified as correct, while the generated code deviates, resulting in failure. *No Code* represents the absence of valid plan and plan verification in the solution generation phase, leading to failure after reaching the maximum number of iterations. *Others* denotes error program solutions caused by various reasons that differ from the previously listed categories.

		HuamEval	MBPP
	MapCoder	18.9	46.4
GDT 4 5	SD	22.6	36.1
GP1-3.5	LDB	28.6	37.7
	LPW (ours)	44.4	36.7

Table 7: The percentage of problems where MapCoder, Self-Debugging (+Expl) (SD), LDB, and LPW generated programs solve the visible tests but fail the hidden tests, out of total failed problems for each method on HumanEval and MBPP, with GPT-3.5 as the backbone.

		HuamEval	MBPP
	MapCoder	4.3	13
<b>GDT 4 5</b>	SD	4.3	10.4
GP1-3.5	LDB	4.9	10.4
	LPW (ours)	4.9	8.8

Table 8: The percentage of problems where MapCoder, Self-Debugging (+Expl) (SD), LDB, and LPW generated programs pass the visible tests but fail the hidden tests, out of a total of 164 problems in HumanEval and 500 problems in MBPP, with GPT-3.5 as the backbone.

plan and plan verification. For LPW, 5.6% of failures result from the generated program solution differing from the plan and plan verification (*Differ from Intended Solution*). For example, LPW fails in the 91st problem, where the generated program is unable to solve the hidden tests due to deviations from the plan and plan verification. The plan verification clearly specifies splitting the input string into sentences using delimiters ".", "?" or "!", but the generated code only handles the full stop case and ignores "?" and "!".

#### D.3 Failure on Hidden Tests

1050

1052

1053

1054

1055

1056

1058

1059

1060

1061

1062

1063

1064

1066

Tables 7 and 8 show the percentage of problems where MapCoder, Self-Debugging (+Expl) (SD), LDB, and LPW generated program solutions pass the visible tests but fail the hidden tests, out of respectively failed problems and the total number of problems in the HumanEval and MBPP benchmarks using GPT-3.5 as the backbone. In Table 7, 1067 44.4% of failures in LPW result from solving the 1068 visible tests only on the HumanEval benchmark, 1069 since except for the No Code category, other rea-1070 sons discussed in Table 6 could contribute to this 1071 issue. In contrast, less than 30% of failures in Map-1072 Coder, SD and LDB experience this issue on Hu-1073 manEval as the larger number of failed problems in 1074 these methods. In Table 7, MapCoder demonstrates 1075 a higher likelihood of passing only the visible tests, 1076 while all other evaluated approaches show similar percentages on the MBPP benchmark, with the re-1078 maining failures arising from different reasons. We 1079 note that all methods tend to address visible tests 1080 only on the same set of problems in the HumanEval 1081 benchmark, resulting in the similar percentage out 1082 of the total number of problems, as shown in Table 1083 8. In contrast, on the MBPP benchmark, MapCoder 1084 exhibits the highest rate of passing only the visible tests out of the total number of problems, consistent 1086 with the result in Table 7, while LPW demonstrates 1087 the lowest rate, as shown in Table 8. Meanwhile, all methods are prone to addressing visible tests only 1089 on MBPP rather than on HumanEval as indicated 1090 in Table 8. Compared to the detailed problem descriptions in HumanEval, the problem descriptions 1092 in MBPP are concise but lack clarity. For example, 1093 Figure 11 illustrates two problems in MBPP where LPW generated solutions are tailored to the visible 1095 tests but deviate significantly from the canonical 1096 solution. 1097

## E Accuracy of Plans, Plan Verifications, and Programs in LPW using GPT-3.5

1098

1099

1100

#### E.1 Plans and Plan Verifications

We manually investigate the accuracy of solution1101plans and plan verifications generated by GPT-3.51102on the HumanEval benchmark, and the results are1103presented in Table 9. Overall, GPT-3.5 generates1104the correct solution plans and plan verifications in1105natural language for majority of problems. In LPW,1106



Figure 11: Example problems in MBPP.

	Plan and Plan Verification	Correct Plan	Correct Plan Verification
LPW	94.5	92.7	92.7

Table 9: Percentage of problems where the LLM successfully generates the valid plans and plan verifications in the solution generation phase (first column); percentage of problems where the LLM-generated plans are manually classified as correct (middle column), considering no plan cases; and percentage of problems where the LLM-generated plan verifications are manually classified as correct (last column), considering no plan cases. All percentages are reported using GPT-3.5 as the backbone on the HumanEval benchmark, with a total of 164 problems.

GPT-3.5 successfully produces plans and plan veri-1107 fications for 94.5% of the problems. GPT-3.5 gen-1108 erates the correct plans for 92.7% of the problems 1109 and achieves the same accuracy for plan verifica-1110 tions. A common issue in the LLM-generated plan 1111 is the omission of certain conditions. For exam-1112 ple, solution plan frequently overlooks uppercase 1113 1114 situations and negative numbers. We note that the LLM-generated plan verification closely adheres 1115 to the solution plan. When the plan is accurate, 1116 the verification process strictly follows the plan 1117 logic, resulting in a correct analysis. Conversely, 1118 if the plan contains logical errors or omits edge 1119 cases, the verification process replicates these mis-1120 takes. Specifically, for LPW, all correct plans lead 1121 to accurate plan verifications, and vice versa. 1122

#### 1123 E.2 Plan Verifications and Programs

We further manually explore the relationship be-1124 tween plan verification and program solution on 1125 the HumanEval benchmark with GPT-3.5. Table 1126 10 evaluates the conditional probabilities between 1127 wrong code and wrong plan verification, as well as 1128 between correct code and correct plan verification. 1129 Typically, in LPW, accurate plan verifications sig-1130 nificantly contribute to correct program solutions, 1131 1132 whereas incorrect plan verifications inevitably result in errors. LPW generates program solutions 1133 based on plans and plan verifications. Therefore, 1134 any accurate descriptions or mistakes, including 1135 missed conditions, in the plan and plan verification 1136

are propagated to the code. When plan verifications are accurate, 96.1% of program solutions are correct in LPW. The remaining incorrect instances arise from unclear condition statements in plan verification that fail to account for hidden tests, leading to erroneous program solutions.

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

The results from Tables 9 and 10 highlight the impressive capability of LLMs in tackling text-tocode generation tasks when outputs are represented in natural language. Plan and plan verification generation accuracy is typically higher than code generation accuracy, underscoring the rationale behind LPW, which produces the high-quality program solution by leveraging plan and plan verification. It is worth exploring methods to help LLMs overcome the challenges of translating natural language solutions into programs, given the strict lexical, grammatical, and semantic constraints. Integrating alternative solution representations (Zelikman et al., 2023) alongside natural language representations could offer a promising approach.

## F Refinement Consistency in LPW

LPW allows multiple rounds of debugging to re-1159 fine code based on error analysis, generated by 1160 comparing the code execution trace with plan ver-1161 ification on the failed visible test. LPW queries 1162 LLMs to generate refined code accompanied by 1163 a refinement explanation, detailing the modifica-1164 tions implemented to address the errors identified 1165 in the error analysis. For instance, Figures 14 and 1166

	Wrong Code $\leftarrow$ Wrong Plan Verification	$Correct \ Code \leftarrow Correct \ Plan \ Verification$
LPW	100	96.1

Table 10: The relationship between LLM-generated code solutions and plan verifications on the HumanEval benchmark with GPT-3.5. The first column shows the percentage of problems where the LLM generates incorrect code solutions when plan verifications are incorrect; the second column shows the percentage of problems where correct code solutions are generated when plan verifications are correct.



Figure 12: Pass@1 accuracy of Baseline, LDB, and LPW across different difficulty levels, Easy, Medium, and Hard, on the LiveCode benchmark when using GPT-40 as the LLM backbone.

		LiveCode
CDT 4a	LDB	23.4
OF 1-40	LPW (ours)	31.6

Table 11: The percentage of problems where LDB and LPW generated programs solve the visible tests but fail the hidden tests, out of total failed problems for each method in LiveCode, with GPT-40 as the backbone.

7 illustrate two HumanEval problems where LPW 1167 successfully generates the correct program through 1168 1169 refinements informed by the error analysis using the GPT-3.5 backbone. We note that in LPW, if the 1170 refined code is irrelevant to the error analysis, the 1171 entire debugging process degrades to a simple sam-1172 pling approach, contradicting our original intent. 1173 As a result, we manually evaluate the debugging 1174 consistency among the generated error analysis 1175 (part (e)), the refined code (part (f)), and the refine-1176 ment explanation (part (g)), as exampled in Figure 1177 14. LPW demonstrates excellent consistency be-1178 tween error analysis and the refined code, where 1179 only one refined code deviates from the error anal-1180 ysis yet still produces the correct solution, across 1181 1182 all problems solved through debugging. This result validates the effectiveness of the debugging steps 1183 in the code implementation phase for LPW, where 1184 the meaningful error analysis enables LLMs to pro-1185 duce the correct program with precise refinements. 1186

		LiveCode
GPT-40	LDB	10.7
	LPW (ours)	12.9

Table 12: The percentage of problems where LDB and LPW generated programs pass the visible tests but fail the hidden tests, out of a total of 140 problems in Live-Code, with GPT-40 as the backbone.

#### G **Analysis of Unsolved Problems for** LPW using GPT-40

1187

1188

1189

1192

1197

1198

1199

1203

1204

1206

1207

1208

## G.1 HumanEval

LPW achieves 98.2% Pass@1 accuracy on Hu-1190 manEval with the GPT-40 backbone, indicating 1191 only 3 unsolvable problems. We further investigate the reasons behind GPT-4o's failures on the 1193 91st, 132nd, and 145th problems as shown in Fig-1194 ures 15, 16, and 17, and attempt to generate the 1195 correct program solutions. The 91st problem fails 1196 since GPT-40 misinterprets the linguistic distinction between the word and the letter; the 132nd problem's ambiguous description challenges GPT-40; and the incomplete description of the 145th 1200 problem leads to failed plan verifications. LPW 1201 successfully generates correct program solutions 1202 for 2 out of 3 problems, achieving 99.4% Pass@1 accuracy, by involving an additional visible test to validate the intended solution for the 91st problem 1205 and providing a comprehensive problem description for the 145th problem.

#### G.1.1 The 91st Problem

Figure 15 illustrates the 91st problem in Hu-1209 manEval, where the GPT-40 generated code (part 1210 (c)) contains an incorrect condition. The code veri-1211 fies if the sentence starts with the letter "I", which is 1212 inconsistent with the problem description (part (a)) 1213 that requires the sentence to start with the word "I". 1214 The provided visible tests (part (b)) fail to clarify 1215 the correct condition, resulting in the error program 1216 passing the visible tests only. Inspired by the supe-1217 rior learning-from-test capacity discussed earlier, 1218 we convert a failed hidden test into a visible test, 1219

		APPS	CodeContests
CDT 4o	LDB	23.1	27.4
GF 1-40	LPW (ours)	23.1	29.6

Table 13: The percentage of problems where LDB and LPW generated programs solve the visible tests but fail the hidden tests, out of total failed problems for each method in APPS and CodeContests, with GPT-40 as the backbone.

highlighted in red in part (d). Consequently, GPT-40 successfully generates the correct program, as shown in part (e).

### G.1.2 The 145th Problem

1220

1221

1223

1224

1225

1226

1227

1228

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240 1241

1242

1243

1244

1245

1246

1247

1248

1249

1250

1251

1252

1253

1254

1255

1257

1258

Figure 17 displays the *145th* problem, where the incomplete problem description (part (a)) results in incorrect plan verification on visible tests (part (b)), leading to a failure after reaching the iteration threshold. The problem description requires returning a list sorted by the sum of digits but omits the specification regarding the sign of negative numbers. This omission confuses GPT-40, resulting in consistently incorrect sorting when verifying the solution plan on the first visible test. We refine the problem description with a detailed explanation on handling both positive and negative numbers (part (c)), leading to the correct program solution, as shown in part (d).

### G.1.3 The 132nd Problem

Figure 16 illustrates the 132nd problem, where ambiguities in the problem description (part (a)) challenge GPT-40. The problem description lacks clarity on "a valid subsequence of brackets" and fails to specify the meaning of "one bracket in the subsequence is nested". We deduce the intended problem description by prompting GPT-40 with a canonical solution (part (d)). However, the LLM-generated description remains unclear and results in various erroneous programs. Furthermore, adding typically failed hidden tests to the visible tests (part (b)) is also ineffective in clarifying the correct logic. We acknowledge that a clearer description might help generate the correct program. However, some problems are inherently difficult to describe accurately in natural language without careful organization, posing challenges for LLMs.

#### 1256 G.2 LiveCode

LiveCode, a contamination-free dataset, serves as a reliable benchmark for evaluating code genera-

		APPS	CodeContests
CDT 4a	LDB	10.8	19.3
GP1-40	LPW (ours)	8.6	19.3

Table 14: The percentage of problems where LDB and LPW generated programs pass the visible tests but fail the hidden tests, out of a total of 139 problems in APPS and 150 problems in CodeContests, with GPT-40 as the backbone.



Figure 13: Pass@1 accuracy of Baseline, LDB, and LPW across different difficulty levels, *Introductory*, *Interview*, and *Competition*, on the APPS benchmark when using GPT-40 as the LLM backbone.

tion performance. It presents challenges for the advanced LLM GPT-40, with all competing approaches showing performance limitations. We note that LPW consistently surpasses Baseline and LDB by approximately 15% and 5%, respectively, in Pass@1 accuracy and across different difficulty levels, as discussed below, emphasizing the reliability of LPW.

1259

1260

1261

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

## G.2.1 Performance Across Different Difficulty Levels

Figure 12 compares the Pass@1 accuracy of Baseline, LDB, and LPW across varying difficulty levels, *Easy*, *Medium*, and *Hard* on the LiveCode benchmark using GPT-40. LPW achieves the highest Pass@1 accuracy across all levels, surpassing LDB by over 5% accuracy in each level. Compared to Baseline, LPW delivers over 15% higher accuracy at the *Easy* and *Medium* levels and 7.5% higher at the *Hard* level. Conversely, LDB performs similarly to Baseline at the *Hard* level, underscoring its limited refinement capability in addressing more complex tasks.

#### G.2.2 Failure on Hidden Tests

Tables 11 and 12 show the percentage of problems1282where the program solutions generated by LDB and1283LPW pass only the visible tests but fail the hidden1284

tests out of the total number of failed problems and the total number of problems, respectively, on the LiveCode benchmark using GPT-40. For LPW, 31.6% of failures occur when only the visible tests are solved, while for LDB, this percentage is 23.4%, as shown in Table 11. In Table 12, both LDB and LPW generate a similar proportion of solutions that pass only the visible tests. However, LPW tends to solve only the visible tests in 12.9% of problems, compared to 10.7% for LDB.

## G.3 APPS and CodeContests

1285

1286

1287

1288

1290

1291

1292

1293

1294

1295

1297

1298

1299

1300

1301

1302

1303

1304

1305

1306

1307

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

APPS and CodeContests are unstructured benchmarks where visible tests are intermingled with the problem statements and function signatures are excluded. To align input data structure across benchmarks, we instruct GPT-40 to derive the optimal function signature and identify visible tests for each problem in APPS and CodeContests prior to conducting experiments. Example structured problems from APPS and CodeContests are illustrated in Figures 18 and 19. LPW demonstrates significant improvements on APPS and CodeContests, exceeding around 10% and 5% Pass@1 accuracy, respectively, compared to LDB with GPT-40. However, in contrast to the performance on the HumanEval and MBPP benchmarks, where LPW achieves 98.2% and 84.8% Pass@1 accuracy, the 62.6% accuracy on APPS and 34.7% accuracy on CodeContests indicate that even for the advanced LLM GPT-40, code generation remains challenging when addressing complicated programming problems, such as those encountered in collegiate programming competitions like IOI and ACM (Hendrycks et al., 2021).

## G.3.1 Performance Across Different Difficulty Levels

Figure 13 compares the Pass@1 accuracy of Base-1321 line, LDB, and LPW across different difficulty lev-1322 els, Introductory, Interview, and Competition, on 1323 the APPS benchmark using GPT-40. LPW con-1324 sistently dominates in Pass@1 accuracy across all difficulty levels. LPW shows strong performance 1326 on the relatively easier levels, i.e., Introductory and 1327 Interview, surpassing LDB by around 9% and 13% accuracy, respectively, and outperforming Baseline 1329 1330 by over 20% accuracy. For the problems belonging to the most challenging level, Competition, LPW 1331 achieves 34.8% accuracy, compared to 28.3% for 1332 LDB and 17.4% for Baseline. However, all ap-1333 proaches experience a substantial decrease at the 1334

*Competition* level, emphasizing the necessity for further improvements.

1335

1336

1337

1366

1367

1368

1369

## G.3.2 Failure on Hidden Tests

Tables 13 and 14 present the percentage of prob-1338 lems where the generated program solutions from 1339 LDB and LPW solve visible tests but fail hidden 1340 tests out of the total failed problems and the total 1341 number of problems, respectively, on the APPS 1342 and CodeContests benchmarks using GPT-40 as 1343 the backbone. In Table 13, 23.1% of failures 1344 result from passing only the visible tests on the 1345 APPS benchmark, with this percentage increasing 1346 to around 30% on CodeContests for both LDB and 1347 LPW. In Table 14, LDB and LPW display similar 1348 percentages of solving visible tests only on each 1349 benchmark, ranging from around 10% on APPS to 1350 19% on CodeContests. Compared to the results in 1351 Table 8, where LDB and LPW address only visible 1352 tests in around 5% of problems on the HumanEval 1353 benchmark, LDB and LPW exhibit weaker perfor-1354 mance on the more challenging APPS and Code-1355 Contests benchmarks. This is particularly evident 1356 on CodeContests, where the percentage is roughly 1357 twice as high as APPS for both LDB and LPW. In 1358 APPS and CodeContests, each problem averages 1359 approximately 2 visible tests, while CodeContests 1360 includes more comprehensive hidden tests, averag-1361 ing about 23 per problem compared to only around 1362 5 per problem in APPS, increasing the likelihood 1363 of solving only the visible tests. 1364

## G.4 Prompts for LPW

We provide the LLM prompts used in LPW in Prompts 1 to 8. For conciseness, we only include one example in each prompt. Full prompts can be found in our released code.

	Analysis Before Coding		Coding With Debugging			
	Plan	Plan Verification	Code Explanation	Runtime Information	Intended Solution	
SP	1	×	×	×	×	
MapCoder	1	×	×	✓	×	
SD	×	×	1	×	×	
LDB	×	×	1	✓	×	
LPW (ours)	1	<ul> <li>Image: A second s</li></ul>	1	1	1	

Table 15: Features of Self-Planning (SP), MapCoder, Self-Debugging (+Expl) (SD), LDB, and LPW.



Figure 14: The *12th* problem in HumanEval, where LPW with GPT-3.5 generated initial code (part (c)) is unable to solve one of the visible tests (part (d)). The refined code (part (f)) successfully solves both visible and hidden tests based on the error analysis (part (e)). The modification in the refined code aligns with the error analysis, as evidenced by the refinement explanation (part (g)).



Figure 15: The problem description (part (a)) and visible tests (part (b)) of the *91st* problem in HumanEval, where GPT-40 generated code (part (c)) addresses the visible tests but fails the hidden tests. However, after converting a failed hidden test to a visible test (part (d)), GPT-40 successfully generates the correct program (part (e)).



Figure 16: The problem description (part (a)) and visible tests (part (b)) of the *132nd* problem in HumanEval, where the GPT-40 generated error code (part (c)) passes the visible tests yet fails the hidden tests. GPT-40 consistently generates incorrect programs despite providing additional visible tests or refining the problem description.



Figure 17: The problem description (part (a)) and visible tests (part (b)) of the *145th* problem in HumanEval where GPT-40 fails to respond with a valid plan verification, resulting in failure. However, after refining the problem description (part (c)), GPT-40 successfully generates the correct program (part (d)).

3231st Problem (APPS)	/ def case_unification(s: str) -> str:
<pre># Task Given an initial string s, switch case of the minimal possible number of letters to make the whole string written in the upper case or in the lower case. # Input/Output [input] string s String of odd length consisting of English letters. 3 ≤ inputString.length ≤ 99. [output] a string The resulting string. # Example For s = "Aba", the output should be "aba" For s = "ABa", the output should be "ABA"</pre>	<pre># Task Given an initial string s, switch case of the minimal possible number of letters to make the whole string written in the upper case or in the lower case. # Input/Output [input] string s String of odd length consisting of English letters. 3 ≤ inputString.length ≤ 99. [output] a string The resulting string. # Example For s = "Aba", the output should be "aba" For s = "ABa", the output should be "ABA" """</pre>

Figure 18: An example structured APPS problem with a function signature and visible tests, generated by instructing GPT-40 with the unstructured problem description.

	(-) C(- (- 1D 11
137th Problem (CodeContests)	def AkString(k: int, s: str) -> str:
(a) Unstructured Problem A string is called a k-string if it can be represented as k concatenated copies of some string. For example, the string "aabaabaabaab" is at the same time a 1-string, a 2- string and a 4-string, but it is not a 3-string, a 5-string, or a 6-string and so on. Obviously any string is a 1- string. You are given a string s, consisting of lowercase English letters and a positive integer k. Your task is to reorder the letters in the string s in such a way that the resulting string is a k-string.	<ul> <li>A string is called a k-string if it can be represented as k concatenated copies of some string. For example, the string "aabaabaabaab" is at the same time a 1-string, a 2-string and a 4-string, but it is not a 3-string, a 5-string, or a 6-string and so on. Obviously any string is a 1-string. You are given a string s, consisting of lowercase English letters and a positive integer k. Your task is to reorder the letters in the string s in such a way that the resulting string is a k-string.</li> <li>Input: The first input line contains integer k (1 ≤ k ≤ 1000). The second line contains s, all characters in s are lowercase English letters. The string length s satisfies the inequality 1 ≤</li> </ul>
Input: The first input line contains integer k ( $1 \le k \le 1000$ ). The second line contains s, all characters in s are lowercase English letters. The string length s satisfies the linequality $1 \le  s  \le 1000$ , where $ s $ is the length of string s. Output: Rearrange the letters in string s in such a way that the result is a k-string. Print the result on a single	[s] \$\ge 1000\$, where [s] is the length of string s.         Output: Rearrange the letters in string s in such a way that         the result is a k-string. Print the result on a single output         line. If there are multiple solutions, print any of them. If the         solution doesn't exist, print "-1" (without quotes).         Examples
them. If the solution doesn't exist, print "-1" (without	Input 2 aazz Output azaz
quotes). Examples	Input 3 abcabcabz Output -1
Input 2 aazz Output azaz	AkString(2, 'aazz') = 'azaz'
Input 3 abcabcabz Output -1	AkString(3, 'abcabcabz') == '-1'

Figure 19: An example structured CodeContests problem with a function signature and visible tests, generated by instructing GPT-40 with the unstructured problem description.

Listing 1: Prompt for plan generation

 П

===========	
You are a F IN ENGLISH	Python writing assistant that responds with a step-by-step thought process ( ) to solve Python coding problems.
===========	
(ou will be Example] ar Yython cod: function si solution pi to solve th	e provided with a series of examples, where each example begins with [Start of ends with [End Example]. In each example, you will be presented with a ing problem, starting with [Example Problem Description], which includes the ignature and its accompanying docstring. You will then provide a reasonable lan, starting with [Example Start Plan] and ending with [Example End Plan], he given problem.
[Start Exar [Example Pu def encrypt """	nple] roblem Description] c(s):
Create encrypt manner """	a function encrypt that takes a string as an argument and returns a string ed with the alphabet being rotated. The alphabet should be rotated in a such that the letters shift down by two multiplied to two places.
[Example Si Create an a Loop throug Return the [Example En [End Examp]	cart Plan] alphabet, biased by two places multiplied by two. gh the input, find the letter biased by the alphabet. result. nd Plan] Le]
Authors found in ou	s' notes: We omit another example for conciseness. The full prompt can be ur released code
Lastly, you ], which in The phrase to create a	will be given a Python writing problem, beginning with [Problem Descriptioncludes the function signature, its docstring, and any potential constraints "Let's think step by step" will signal the start of the plan. Your task is a solution plan, starting with [Start Plan] and ending with [End Plan].

Listing 2: Prompt for plan verification

=	System
١	You are a logical reasoner tasked with performing a step-by-step analysis to derive the correct solution to the given problem based on the provided plan.
=	User
YSFSLEtVoi rtoe	You will be presented with several plan verification examples, each starting with [ Start Example] and ending with [End Example]. In each example, you will be given a Python writing problem, starting with [Example Problem Description], followed by the solution plan starting with [Example Solution Plan], and its verification process beginning with [Example Plan Verification for X] for a test case X, starting with [ Example Test Cases]. During the verification process, intermediate variables that need to be recorded are clearly identified at the outset, starting with [Record Analysis]. Whenever the value of a recorded intermediate variable is updated, the new result is clearly displayed, beginning with [Record]. After the verification, the derived result is compared to the correct test result, starting with [Results Compare]. If the derived result matches the test result, the output will be [Correct Plan]. If the derived result differs from the test result, the output will be [Incorrect Plan], followed by the reasons for the discrepancy, starting with [Incorrect Reasons], and the revised correct solution plan, beginning with [Start Revised Solution Plan] and ending with [ End Revised Solution Plan].
[ [ c	[Start Example] [Example Problem Description] def prime_number(n: int): """

In range 0 to 100, returns n-th number that is a prime.

...... 1438 1439 [Example Solution Plan] 1440 1. Iterate number through 0 to 100. 1441 2. Check each number, if it's prime. 1442 Keep track of the count of prime numbers found.
 Stop when we find the n-th prime number. 1443 1444 5. Return the nth prime number. 1445 1446 [Example Test Cases] assert prime\_number(3)==5 1447 [Example Plan Verification for assert prime\_number(2)==3] 1448 [Record Analysis] 1449 The return value is the nth prime number, so all nth prime numbers need to be clearly 1450 recorded! 1451 1. Call the function prime\_number(2). 1452 2. According to line 1 in solution plan, Iterate number through 0 to 100. 1453 3. According to line 2 in solution plan, Check if 0 is prime. It's not. 1454 4. Move to next number 1. 1455 5. According to line 2 in solution plan, Check if 1 is prime. It's not. 1456 6. Move to next number 2. 1457 7. According to line 2 in solution plan, Check if 2 is prime. It is a prime. 1458 8. According to line 3 in solution plan, the count of prime numbers is 1. 1459 [Record]: 1st prime number is 2 1460 9. Move to next number 3. 1461 10. According to line 2 in solution plan, Check if 3 is prime. It is a prime. 1462 11. According to line 3 in solution plan, the count of prime numbers is 2. 1463 [Record]: 2nd prime number is 3 1464 12. According to line 4 in solution plan, Stop when we find the 2nd prime number. 1465 13. According to line 5 in solution plan, Return the 2nd prime number, which is 3. 1466 [Results Compare] 1467 The test correct output is 3. The logic analysis output is 3. 3=3. Thus, the plan is 1468 verified to correctly handle all test cases. 1469 [Correct Plan] 1470 [End Example] 1471 1472 ... Authors' notes: We omit another example for conciseness. The full prompt can be 1473 found in our released code. ... 1474 1475 Finally, you will be given a problem description, beginning with [Problem Description], 1476 along with your generated solution plan, starting with [Solution Plan], to solve the [ 1477 Problem Description], and multiple test cases starting with [Test Cases]. The phrase Let's verify the plan" will indicate the beginning of the verification process, 1478 1479 followed by your verification steps to confirm whether your generated plan can pass all 1480 test cases. 1481 1482 For each test case, the verification must include [Record Analysis] to track the 1483 intermediate variables at the beginning. If any intermediate variable value is updated 1484 during the reasoning process, the updated value should be clearly displayed, starting 1485 with [Record]. Please include [Results Compare] to assess the derived outcome against 1486 the correct test output. If the derived result matches the test result, output [Correct 1487 Plan] and proceed to the next test case. If the derived result does not match the test 1488 result, output [Incorrect Plan], followed by the reasons for the discrepancy, starting 1489 with [Incorrect Reasons]. Finally, provide the revised solution plan, starting with [ 1490 Start Revised Solution Plan] and ending with [End Revised Solution Plan], to complete 1491 the process. 1493 Listing 3: Prompt for plan verification check 1404

======================================	1495
	1496
You are a logical reasoner. Your goal is to identify any incorrect logic within the	1497
logic verification process.	1498
	1499
User	1500
	1501
You will be given several examples demonstrating how to evaluate a logic verification	1502
process. Each example will begin with [Start Example] and end with [End Example]. In	1503
each example, you will find the following:	1504
	1505
[Example Problem Description] outlining the Python writing problem:	1506

[Example Solution Plan] describing the approach to solve the problem; [Example Plan Verification for X], applying the solution plan to a specific test case X. In this process, the intermediate variables to be tracked are analyzed at the start, marked by [Record Analysis]. Whenever the value of a recorded intermediate variable is updated, its new value is displayed starting with [Record]. The [Results Compare] section compares the verification derived result with the correct test output; [Example Verification Check for X], this section evaluates, step by step, whether the logic verification process for test case X is correct. If the verification is correct, the output will be [Correct Plan Verification], and please proceed to the next example. If the verification is incorrect, explanation should be provided and [Incorrect Plan Verification] will be the output to conclude the evaluation [Start Example] [Example Problem Description] def addOne(message: str): You are given a large integer represented as an integer array digits, where each digits[i] is the ith digit of the integer. The digits are ordered from most significant to least significant in left-to-right order. The large integer does not contain any leading 0's. Increment the large integer by one and return the resulting array of digits. [Example Solution Plan] 1. Convert the list of digits into a number. 2. Increment the number by one. 3. Convert the incremented number back into a list of digits and return it. [Example Plan Verification for assert addOne([1,2,3])==[1,2,4]] [Record analysis] The return value is the incremental resulting array of digits, so the incremental resulting array of digits needs to be clearly recorded! According to line 1 in solution plan, convert [1,2,3] to the number 123. According to line 2 in solution plan, Increment 123 by one to get 124. According to line 3 in solution plan, convert 124 back into the list [1,2,4] [Record]: incremental resulting array is [1,2,4] According to line 3 in solution plan return incremental resulting array [1,2,4]. [Results Compare] The test correct output is [1,2,4]. The logic analysis output is [1,2,4]. [1,2,4]=[1,2,4]. So the plan is verified to correctly handle all test cases. [Correct Plan] [Example Verification Check for assert ddOne([1,2,3])==[1,2,4]]: "Convert [1,2,3] to the number 123" is correct! "Increment 123 by one to get 124" is correct! since 123+1=124 "Convert 124 back into the list [1,2,4]" is correct! "return incremental resulting array [1,2,4]" is correct! In [Results Compare] "The test correct output = [1,2,4]" is correct! "The logic analysis output = [1,2,4]" is correct! The results comparison "[1,2,4]=[1,2,4]" is correct! All analysis steps are correct! [Correct Plan Verification] [Example Plan Verification for assert addOne([-1,2])==[-1,1]] [Record analysis] The return value is the incremental resulting array of digits, so the incremental resulting array of digits needs to be clearly recorded! According to line 1 in solution plan, convert [-1,2] to the number 12. According to line 2 in solution plan, Increment 12 by one to get 13.

1507 1508

1509 1510

1511

1512

1513

1514 1515 1516

1517

1518 1519

1520 1521

1522

1523

1524 1525

1526 1527 1528

1529

1532 1533 1534

1536

1537

1538 1539 1540

1541 1542

1543

1545

1546

1547

1548

1549 1550 1551

1553 1554

1556

1557

1558

1559

1560

1561 1562

1563

1565

1567 1568 1569

1570

1571 1572

1573

1574

1575

According to line 3 in solution plan, convert 13 back into the list [1,3] [Record]: incremental resulting array is [1,3] According to line 3 in solution plan return incremental resulting array [1,3]. [Results Compare] The test correct output is [-1,1]. The logic analysis output is [-1,1]. [-1,1]=[-1,1]. So the plan is verified to correctly handle all test cases. [Correct Plan] [Example Verification Check for assert addOne([-1,2])==[-1,1]]: "Convert [-1,2] to the number 12" is incorrect. The analysis doesn't correctly interpret the -1 and assumes all values are positive, the sequence -1, 2 should form -12. "Increment 12 by one to get 13" is correct, but as established, the initial conversion should not yield 12. "Convert 13 back into the list [1,3]" is correct! "Return incremental resulting array [1,3]" is correct! In [Results Compare] "The test correct output = [-1,1]" is correct! "The logic analysis output = [-1,1]" is incorrect! The logic analysis result is [1,3] mentioned in the verification "return incremental resulting array [1,3]". The results comparsion "[-1,1]=[-1,1]" is incorrect! The logic analysis result is [1,3] and [-1,1] is not equal [1,3]. The logic verification process for addOne([-1,2])==[-1,1] is incorrect. The analysis doesn't correctly interpret the -1 and assumes all values are positive, the sequence -1, 2 should form -12. The logic analysis output = [-1,1] is incorrect! It is [1,3]. The results comparison is incorrect since [-1,1] is not equal [1,3]. [Incorrect Plan Verification] [End Example] ... Authors' notes: We omit another example for conciseness. The full prompt can be found in our released code. ... Finally, you will be given a problem description, beginning with [Problem Description], followed by your generated solution plan, starting with [Solution Plan], to address the [Problem Description]. You will then work through multiple Plan Verification, each starting with [Plan Verification for X], where X represents a test case. At the start of the verification process, [Record Analysis] examines the intermediate variables that should be tracked. During the logic verification, the tag [Record] indicates any updates to the values of the recorded intermediate variables. The [Results Compare] section documents the comparison between the verification derived result and the expected test output. The phrase "Let's evaluate the verification" will indicate the start of the evaluation for each verification process. This will be followed by your step-by-step verification check to assess whether each intermediate output in the verification process is correct, starting with [Verification Check for X], as shown in the examples. If all intermediate results in the verification process are correct, the output will be [ Correct Plan Verification], and you will proceed to the next verification. If the verification process is incorrect, an explanation should be provided, and [Incorrect Plan Verification] will be output to conclude the evaluation. Listing 4: Prompt for initial code generation 

You are a Python writing assistant that only responds with Python programs to solve a Python writing problem.

You will receive several examples, each structured as follows, beginning with [Start Example] and ending with [End Example]. Within each example, you will find a Python programming problem starting with [Example Problem Description] and a solution plan starting with [Example Solution Plan]. Additionally, you will receive plan

verifications for specific test cases. For each test case X, the plan verification is labeled as [Example Plan Verification for X], providing a detailed logical breakdown and variable value updates, which are recorded starting with [Record]. Following the verification, you will encounter the example-generated program starting with [Example Generated Program]. The program, marked from [Start Program] to [End Program], is generated based on the solution plan and plan verification, ensuring that the program's execution aligns with the plan verification when test case X is used as input. [Start Example] [Example Problem Description] from typing import List def get\_closest\_transition\_character(word): You are given a word. Your task is to find the closest transition character from the right side of the word(case sensitive). The transition character is lowercase and the character after it is uppercase. If no such lowercase character is found, return an empty string. >>> get\_closest\_transition\_character("eAsy") == "s" ..... [Example Solution Plan] 1. Reverse iterate through the characters of the word starting from the last character from the right. 2. For each character, check if the current character is lowercase and if the character after it is uppercase. If step 2 is satisfied, 4. return the current vowel character. 5. If no such vowel is found, return an empty string. [Example Plan Verification for assert get\_closest\_transition\_character("eAsy")=="s"] [Record analysis] The return value is the closest transition character, so the closest transition character should be recorded! Call the function get\_closest\_vowel("eAsy"). 2. According to line 1 in the solution plan, reverse iterate the word, from the last character to the first character, so the last character is "y" 3. According to line 2 in the solution plan, "y" is a lowercase but the character after "y" is "s" and "s" is a lowercase. 4. Move to the next character based on the reverse iterate, so the character is "s". 5. According to line 2 in the solution plan, "s" is a lowercase and the character after 's' is 'A', and 'A' is uppercase. 6. According to line 3 in the solution plan, step 2 is satisfied, [Record]: the closest transition character 's 8. According to line 4 in the solution plan, return the current lowercase character 's' [Example Generated Program] [Start Program] from typing import List def get\_closest\_transition\_character(word): "" You are given a word. Your task is to find the closest transition character from the right side of the word(case sensitive). The transition character is lowercase and the character after it is uppercase. >>> get\_closest\_transition\_character("eAsy") == "s" ..... # reverse iterate the word for i in range (len(word)-1,-1,-1): current\_character=word[i] if current\_character.islower(): if i!=0: after\_character=word[i-1] if after\_character.isupper(): return current\_character return "" [End Program]

1646

1647

1648 1649

1650

1651

1653 1654

1655 1656

1657

1658 1659 1660

1661

1662

1664

1665 1666 1667

1668 1669

1671

1672

1673

1674 1675 1676

1677 1678

1679

1681 1682

1684

1685 1686

1687

1688 1689

1690

1693

1695

1696

1697 1698

1699

1700

1701

1702

1704

1705 1706

1707

1708

1709

1710

1711

1712

[End Example]

... Authors' notes: We omit another example for conciseness. The full prompt can be found in our released code. ...

1716 1717

1718

1719 1720

1721

1722

1723

1724

1725 1726

1727

1728

1729

1730

1731

1732

1733

1735

Finally, you will be provided with a Python writing problem, starting with [Problem Description]. A solution plan will follow, beginning with [Solution Plan]. Next, you will receive several plan verifications. For each test case X, the plan verification, starting with [Plan Verification for X] provides detailed logical reasoning steps to solve it.

Once the plan verification is provided, the "Let's generate the program" flag indicates the start of Python program generation. You will then need to generate the Python program solution for the problem. The plan verification serves as a constraint during program generation. It is essential to ensure that the execution of the generated program remains consistent with [Plan Verification for X] when using test case X as input. Additionally, the generated program should incorporate all conditions noted in [ Plan Verification for X] to solve test case X. Please ONLY output the generated Python program, starting with [Start Program] and ending with [End Program].

Listing 5: Prompt for print state
-----------------------------------

1736 1737 1738 You are a Python writing assistant that only responds with Python programs with PRINT 1739 statements. 1740 1741 1742 1743 You'll be provided with several examples structured as follows, beginning with [Start 1744 Example] and ending with [End Example]. In each example, you will be given a sample Python program, starting with [Example Python Program]. You will also receive several 1745 1746 plan verifications for specific test cases. For a test case X, its plan verification, 1747 starting with [Example Plan Verification for X], includes a worded description of the 1748 logic used to solve test case X. During the verification, the intermediate variable 1749 that needs to be tracked is clearly identified, starting with [Record Analysis] at the 1750 beginning, and any updates to its value are recorded, starting with [Record]. 1751 1752 Following this, you will be shown a Python program that includes detailed print 1753 statements, starting with [Example Python Program with Print Statements]. These print 1754 statements illustrate how the values of the intermediate variables (described in the 1755 plan verification) are modified during program execution, as well as how other 1756 variables in the program change. These examples will guide you on where and how to add 1757 print statements in your Python program. 1758 1759 [Start Example] 1760 1761 [Example Python Program] 1762 from typing import List 1763 def get\_closest\_transition\_character(word): 1764 "" You are given a word. Your task is to find the closest transition character from 1765 the right side of the word(case sensitive). The transition character is lowercase 1766 and the character after it is uppercase. 1767 >>> get\_closest\_transition\_character("eAsy") == "s"
""" 1768 1769 for i in range (len(word)-1,-1,-1): 1770 current\_character=word[i] 1771 if current\_character.islower(): 1772 if i!=0: 1773 after\_character=word[i-1] 1774 if after\_character.isupper(): 1775 return current\_character 1776 return "" 1778 [Example Plan Verification for assert get\_closest\_transition\_character("eAsy")=="s"] 1779 [Record analysis] 1780 The return value is the closest transition character, so the closest transition 1781 character should be recorded! 1782 1783 1. Call the function get\_closest\_vowel("eAsy"). 1784

1785 2. According to line 1 in the solution plan, reverse iterate the word, from the last character to the first character, so the last character is "y" 1786 1787 3. According to line 2 in the solution plan, "y" is a lowercase but the character after "y" is "s" and "s" is a lowercase. 1788 4. Move to the next character based on the reverse iterate, so the character is "s". 5. According to line 2 in the solution plan, "s" is a lowercase and the character after 1789 1790 's' is 'A', and 'A' is uppercase. 1791 6. According to line 3 in the solution plan, step 2 is satisfied, 1792 7. [Record]: the closest transition character 's 1793 1794 8. According to line 4 in the solution plan, return the current lowercase character 's' 1795 1796 [Example Python Program with Print Statements] 1797 from typing import List def get\_closest\_transition\_character(word):
 """ You are given a word. Your task is to find the closest transition character from 1798 1799 the right side of the word(case sensitive). The transition character is lowercase 1800 and the character after it is uppercase. 1801 >>> get\_closest\_transition\_character("eAsy") == "s" 1803 1804 1805 print(f"Reverse iterate the word {word}") 1806 for i in range (len(word)-1,-1,-1): current\_character=word[i] 1807 print(f"current character at index {i} is {word[i]}") if current\_character.islower(): print(f"current character {word[i]} is lowercase") 1810 if i!=0: 1811 1812 print(f"There is a character after {word[i]}") after\_character=word[i-1] 1813 1814 print(f"character after {word[i]} is {word[i-1]}") 1815 if after\_character.isupper(): 1816 print(f"character is {word[i-1]} is uppercase") print(f"[Record]: the closest transition character {word[i]}") 1817 print(f"Return the closest transition character {word[i]}") 1818 1819 return current\_character 1821 print(f"no such lowercase character is found, return an empty string") return "" [End Example] 1823 1824 1825 ... Authors' notes: We omit another example for conciseness. The full prompt can be 1826 found in our released code. ... 1827 1828 Finally, you will be provided with a Python program, starting with [Python Program], 1829 along with several plan verifications for specific test cases. For each test case X, the plan verification, starting with [Plan Verification for X], includes a detailed description of the logic used to solve test case X. In the plan verification, the intermediate variables to be tracked are clearly analyzed at the beginning, starting 1832 with [Record Analysis], and any updates to these variable values are recorded, starting 1834 with [Record]. 1835 The phrase "Let's add print statements" signals the start of the process to incorporate 1836 1837 print statements into the provided Python program. Your task is to add print 1838 statements that track how the variables in the program change. Ensure that the intermediate variable values (as outlined in the plan verification) are printed using 1839 1840 these statements. Output your program with print statements, starting with [Start Program] and ending with [End Program]. 1842 Listing 6: Prompt for code explanation 1843 1845 You are a Python interpreter tasked with providing an explanation for each line of the 1847 Python program. 1848 1849 1850 You will be given a Python program, and your task is to generate a word-by-word 1851 1852 explanation describing the effect of each line in the program. You will be provided 1853 with several examples, each beginning with [Start Example] and ending with [End Example

]. In each example, you will receive a Python programming problem, starting with [ 1854 Example Problem Description], along with an example Python program, marked as [Example 1855 Python Program], which is generated to solve the given problem. Additionally, you will 1856 be provided with an explanation of each line in the example Python program, starting 1857 with [Example Explanation For Each Line]. 1858 1859 [Start Example] 1860 [Example Problem Description] 1861 def encode(message): 1862 1863 Write a function that takes a message, and encodes in such a way that replaces all 1864 letters in the message with the letter that appears 1 place after of that letter in 1865 the english alphabet and then delete letter is a vowel. 1866 The last letter z is execluded in the message. Assume the input only containing 1867 letters. 1868 ..... 1869 1870 [Example Python Program] 1871 def encode(message): 1872 encoded\_message = "" 1873 for char in message: 1874 next\_char = chr(ord(char) + 1) 1875 if next\_char not in "aeiouAEIOU": 1876 encoded\_message += next\_char 1877 return encoded\_message 1878 1879 [Example Explanation For Each Line] 1880 The Python function encode is designed to manipulate a given string (referred to as a message) by replacing the current letter with the letter that appears 1 place after it 1882 in the English alphabet and then skipping the letter if it is a vowel: 1883 1884 Function Definition (def encode(message):): Defines a function named encode that 1885 accepts one parameter, message. This parameter is intended to be a string that will be 1886 processed to create an encoded message. 1887 1888 Initialize Encoded Message (encoded\_message = ""): Initializes a variable 1889 encoded\_message as an empty string. This variable will store the encoded version of the 1890 message as characters are processed and approved for inclusion. 1891 1892 For Loop (for char in message:): Iterates over each character in the message. Each 1893 character is processed individually. 1894 1895 Calculate Next Character (replacechar = chr(ord(char) + 1)): For each character in the 1896 message, this line calculates its next character that appears 1 place after it in the 1897 English alphabet. It converts the character to its ASCII value with ord(char), 1898 increments this value by 1, and then converts it back to a character with chr(). 1899 1900 Check if the resulting character is a Vowel (if replacechar in "aeiouAEIOU":): Check if 1901 the resulting character (replacechar) after incrementation is a vowel (either 1902 uppercase or lowercase is checked here). If it is a vowel, the continue statement is 1903 executed. 1904 1905 Add Character to Encoded Message (else: encoded\_message += replacechar): If replacechar 1906 is not a vowel, it is appended to encoded\_message. This builds up the final encoded 1907 string with the modified characters. 1908 1909 Return Encoded Message (return encoded\_message): After processing all characters in the 1910 original message, the function returns the fully encoded string which consists of all 1911 non-vowel characters that are the successors of the original characters in the ASCII 1912 sequence. 1913 [End Example] 1914 1915 ... Authors' notes: We omit another example for conciseness. The full prompt can be 1916 found in our released code. ... 1917 1918 Finally, you will be presented with a problem description, starting with [Problem 1919 Description], and your generated Python program, starting with [Python Program], which 1920 is meant to solve the [Problem Description]. After this, the "Let's generate the 1921 explanation" flag will signal the beginning of the explanation phase. Your task is to generate a word-by-word explanation for each line in the Python program, following the 1923

format shown in the previous examples. Please skip the explanation for any line that is a print statement. Output your explanation starting with [Start Explanation] and ending with [End Explanation].

1924

1925

1929

1930 1931

1932

1933 1934 1935

1936 1937

1938

1939

1940

1941

1942 1943

1944

1945 1946

1947

1948

1949

1950

1952

1953

1954 1955

1956

1957

1958 1959

1960

1964

1966 1967

1968

1969 1970

1971

1973

1974

1975

1976

1977

1978 1979 1980

1981 1982

1984

1986

1987 1988

1989

1990 1991

1992

Listing 7: Prompt for error analysis You are a logical reasoner. You will be provided with two logical reasoning processes: [Plan Verification] and [Error Execution Trace]. Your task is to identify any errors in the [Error Execution Trace] by comparing it with the [Plan Verification]. You will be provided with several examples, each starting with [Start Example] and ending with [End Example]. In each example, you will receive a Python programming problem, starting with [Example Problem Description], along with an example of an incorrect Python program, marked as [Example Error Program], generated for that problem. You will also be provided with a detailed execution trace of the error program on the failed test case X, labeled as [Example Error Execution Trace for X], including the intermediate variable values. Additionally, you will be provided with an example of the correct logical reasoning process, labeled as [Example Plan Verification for X]. This process outlines the necessary steps to solve test case X accurately, including condition checks and recording intermediate variable updates, starting with [Record]. Next, [Example Discrepancy Analysis] provides a comparison between the Example Plan Verification and the Example Error Execution Trace, highlighting output differences and identifying where the Error Execution Trace deviates from correctness. Finally, [Example Error Analysis] summarizes the errors identified in the [Example Discrepancy Analysis] and proposes solutions to correct them. [Start Example] [Example Problem Description] def is\_palindrome(num):
 """ check if a given integer is a palindrome. [Example Error Program] def is\_palindrome(num): num\_str = str(abs(num)) return num\_str == num\_str[::-1] [Example Error Execution Trace for assert is\_palindrome(-121)==False] 1. Convert the integer -121 to the string "121" 2. The integer string "121" is equal to the reversed string "121", the result is True 3. Return True [Example Plan Verification for assert is\_palindrome(-121)==False] [Record analysis] The return value is the checking result about a given integer is a palindrome, so the checking result should be clearly recorded! 1. Call the function is\_palindrome(-121). change integer to string, it is "-121"
 check whether the string "-121" is equal to its reversed string "121-", the checking result is False 4. [Record]: checking result = False 5. Return checking result False [Example Discrepancy Analysis] In the plan verification, the recorded value is the checking result: Let's trace the "checking result" value in the plan verification when it is first-time recorded (SKIP INITIALIZATION). In the plan verification, the value of checking result is first-time recorded in Line 4 after executing lines:

1. Call the function is\_palindrome(-121). 1993 change to integer to the string, it is "-121"
 check whether the string "-121" is equal to its reversed string "121-", the checking 1994 1995 result is False 1996 4. [Record]: checking result = False 1997 1998 In the plan verification, the first-time update changes the checking result value to 1999 False. 2000 2001 Let's trace the "checking result" value in the Error Execution Trace. 2002 In Error Execution Trace, the value of checking result is first-time recorded in Line 2 2003 after executing lines 2004 1. Convert the integer -121 to the string "121" 2. The integer string "121" is equal to the reversed string "121", the result is True 2006 2007 In Error Execution Trace, the first-time update changes the checking result value to 2008 True. 2010 The checking result value in the plan verification and Error Execution Trace are NOT 2011 the same, due to False NOT equaling True when the checking result value is first 2012 updated. 2013 2014 Let's carefully analyse the reason with step-by-step thinking: 2015 In lines 1-4 in the plan verification, the integer -121 is first converted to the string "-121". Then "-121" is compared with its reversed string "121-". "-121" is NOT 2016 2017 equaling "121-" so the result is False 2018 2019 In lines 1-2 in Error Execution Trace, the integer -121 is first converted to the 2020 string "121". This is different from the plan verification where converting -121 to string is "-121" rather than "121". Then "121" is compared with its reversed string 2021 "121". "121" is equaling "121" so the result is True. 2023 2024 [Example Error Analysis] 2025 The error execution trace incorrectly converts the negative integer to its negative integer string. The negative signal is missed. For example, negative integer -121 2027 should be converted to string "-121" but not "121. To fix this error, the negative 2028 number must be considered and its negative sign should be contained when converted to 2029 string. Such as negative integer -121 should be converted to string "-121". 2031 [End Example] 2033 ... Authors' notes: We omit another example for conciseness. The full prompt can be 2034 found in our released code. ... 2036 Finally, you will be presented with a problem description, starting with [Problem 2037 Description], along with your generated error program, starting with [Error Program], which attempts to solve the [Problem Description]. You will also receive a detailed execution trace, including intermediate variable values, for the failed test case X, 2040 starting with [Error Execution Trace for X]. This trace is generated by the error 2041 program. Additionally, you will be provided with a correct logical reasoning process, labeled as [Plan Verification for X], which outlines the necessary steps to solve test 2042 case X accurately, including condition checks and recording intermediate variable 2044 updates, starting with [Record]. 2045 2046 Following this, the "Let's do analysis" flag will indicate the start of the analysis 2047 phase. Your task is to analyze where the [Error Execution Trace for X] deviates from the [Plan Verification for X], as demonstrated in the examples. This analysis should be 2049 output starting with [Discrepancy Analysis]. Finally, you should provide a summary of the errors identified in the [Discrepancy Analysis], including the reasons for these mistakes (IN ENGLISH) and suggestions on how to correct them, starting with [Error 2051 2052 Analysis]. 2053

Listing 8: Prompt for code refinement

======================================
You are a Python program fixer. You need to correct an error Python program based on the provided information.
User

2057

You will receive several examples, each structured as follows, starting with [Start Example] and ending with [End Example]. Within each example, you will find a Python programming problem, beginning with [Example Problem Description], followed by an error program provided under [Example Error Program] for the given problem. You will then receive an explanation for the error program, including a line-by-line explanation starting with [Example Error Program Explanation]. Additionally, an error analysis will be provided, starting with [Example Error Analysis ], describing the issues in the error program and offering suggestions for refinement. You will then be provided with the refined Python program under [Example Refined Program], based on the error analysis. Following that, a refinement explanation, starting with [Example Refinement Explanation], will be given to show which lines of the program were changed and explain the reasons for those changes. [Start Example] [Example Problem Description] def is\_palindrome(num): check if a given integer is a palindrome. [Example Error Program] def is\_palindrome(num): num\_str = str(abs(num)) return num\_str == num\_str[::-1] [Example Error Program Explanation] Function Definition (def is\_palindrome(num):): This line defines a function named is\_palindrome that takes one parameter, num. This parameter is expected to be an integer. Convert Number to Absolute String (num\_str = str(abs(num))): A variable num\_str is initialized with the absolute value of num converted to a string. The abs() function removes the sign from num if it's negative, ensuring the palindrome check is based solely on the digits. Check Palindrome and Return (return num\_str == num\_str[::-1]): This line checks if the string representation of num\_str is the same forwards and backwards. It uses the slicing technique [::-1] to reverse the string. If num\_str is equal to its reversed version, the function returns True, indicating the number is a palindrome. Otherwise, it returns False. [Example Error Analysis] The error execution trace incorrectly converts the negative integer to its negative integer string. The negative signal is missed. For example, negative integer -121 should be converted to string "-121" but not "121. To fix this error, the negative number must be considered and its negative sign should be contained when converted to string. [Example Refined Program] def is\_palindrome(num): num\_str = str(num) return num\_str == num\_str[::-1] [Example Refinement Explanation] Program line (num\_str = str(abs(num))) is changed to (str(num)) to convert the negative integer to its negative integer string by deleting the abs function to keep the negative representation as mentioned in the the error analysis. (str(num)) can correctly convert negative integer -121 to string "-121". [End Example] ... Authors' notes: We omit another example for conciseness. The full prompt can be found in our released code. ... You will be presented with a Python writing problem, starting with [Problem Description ]. The error program will be provided under [Error Program], followed by an explanation

2064

2066

2067

2069

2070 2071

2072

2073

2074

2075

2076 2077

2078 2079

2081 2082

2084

2086

2088

2090

2091

2092

2093

2094

2096

2099

2100 2101

2102

2103

2104 2105 2106

2107

2108 2109

2110

2115

2116 2117

2118

2120

2122

2123

2124 2125

2126 2127

2128

2129

of each line, starting with [Error Program Explanation]. You will then receive an	2132
error analysis, starting with [Error Analysis], which describes the issues in the error	2133
program and provides refinement suggestions.	2134
	2135
The repair process will begin with the phrase "Let's correct the program." Based on the	2136
error analysis, generate the refined program. Output your refined program, starting	2137
with [Start Refined Program] and ending with [End Refined Program], ensuring that ONLY	2138
the Python code is included between these markers. Finally, provide a refinement	2139
explanation, starting with [Refinement Explanation], detailing how the program was	2140
modified to align with the error analysis.	2142