MultiPoT: Multilingual Program of Thoughts Harnesses Multiple Programming Languages

Anonymous ACL submission

Abstract

Program of Thoughts (PoT) is an approach characterized by its executable intermediate 002 steps, which ensure the accuracy of the numerical calculations in the reasoning process. Currently, PoT primarily uses Python. However, relying solely on a single language may result in suboptimal solutions and overlook the potential benefits of other programming languages. In this paper, we conduct comprehensive experiments on the programming languages used in PoT and find that no single language consistently delivers optimal performance across all tasks and models. The effectiveness of each language varies depending on the specific scenarios. Inspired by this, we propose a task and model agnostic approach called MultiPoT, which harnesses strength and diversity from 017 various languages. Experimental results reveal that it significantly outperforms Python Self-Consistency. Furthermore, it achieves comparable or superior performance compared to the 021 best monolingual PoT in almost all tasks across 022 all models. In particular, MultiPoT achieves more than 4.6% improvement on average on both Starcoder and ChatGPT (gpt-3.5-turbo).

1 Introduction

037

041

Program of Thoughts (PoT) aims to prompt Code Large Language Models (Code LLMs) to decompose complex problems into successive executable codes (Gao et al., 2023; Chen et al., 2022). The computational process of the final result is decoupled from Code LLMs and accurately executed by an external interpreter. PoT significantly reduces mathematical computation errors and improves reasoning performance (Wang et al., 2023a). Subsequently, benefiting from its flexibility and scalability, it is gradually applied to a broader spectrum of fields like image understanding (Surís et al., 2023) and robotic control (Li et al., 2023a). Nowadays, PoT has become a key method for enabling intelligence in agents (Yang et al., 2024). Today is the last day of the first quarter of 2008. What is the date one year ago from today?



Figure 1: Comparison of PoT with different PLs. Python's 'timedelta' lacks support for year computation, leading to a leap year (2008 has 366 days) error by subtracting 365 days. R and JavaScript directly compute the year and get the correct answer.

Despite significant progress, PoT has a notable limitation: to the best of our knowledge, **all research on PoT focuses on Python**. However, since Code LLMs are capable of multilingual generation,¹ and most of the reasoning tasks are languageindependent, many other programming languages (PLs) can also be applied to PoT, especially when considering their unique strength. From the perspective of **tasks**, different languages represent PoT in different forms. Figure 1 shows that the representation and calculation of dates in R is more concise than in Python. This can reduce the complexity when Code LLMs generate PoTs. From the perspective of **models**, their multilingual ability is inconsistent. For instance, Deepseek Coder's

043

044

045

046

047

051

052

¹In this paper, our "multilingual" represents multiple programming languages, not natural languages.

C++ outperforms Python on the code generation task (Guo et al., 2024). It is natural to wonder whether this phenomenon also occurs on reasoning tasks. Therefore, a crucial question is raised with these perspectives: *Is Python truly the optimal language for all tasks and models for PoT?* Relying on Python may lead to a local optimum. In Figure 1, Python's 'timedelta' does not support 'year', resulting in a miscalculation for the leap year. In contrast, R and JavaScript yield the correct answer.

057

058

059

061

062

063

087

094

100

101

102

103

104

105

106

108

Motivated by this, we conduct comprehensive experiments for multilingual PoTs. Beyond Python, we select four PLs: three widely used general languages (JavaScript, Java, and C++) and a niche but comprehensive language (R). For a comprehensive comparison, we identify five distinct sub-tasks within reasoning tasks: math applications (Cobbe et al., 2021; Patel et al., 2021; Miao et al., 2020), math (Hendrycks et al., 2021), tabular, date, and spatial (Suzgun et al., 2022). We select four backbone LLMs: three strongest Code LLMs (Starcoder (Li et al., 2023b), Code Llama (Roziere et al., 2023), and Deepseek Coder (Guo et al., 2024)) and code-capable ChatGPT (gpt-3.5-turbo). Under both greedy decoding and Self-Consistency (Wang et al., 2022) settings, we answer that "Python is not always the optimal choice, as the best language depends on the specific task and model being used."

In addition to the analysis contribution, to leverage the strength of multiple PLs, we further introduce a simple yet effective approach, called Multi-PoT (Multilingual Program of Thoughts). Multi-PoT is a task and model agnostic approach, which uses LLMs to synchronously generate PoTs with various PLs and subsequently integrates their results via a voting mechanism. The use of multiple PLs also provides greater diversity and reduces the probability of repeating the same errors compared to single-language sampling. Experimental results demonstrate that MultiPoT outperforms Python Self-Consistency significantly. Furthermore, MultiPoT achieves great performance across nearly all tasks and models. It effectively matches or even surpasses the top-performing languages in each specific scenario, and outperforms on task and model averages. Especially on both ChatGPT and Starcoder, MultiPoT performs the best on four out of five tasks, with only a slight underperformance on the remaining task, and shows an improvement of over 4.6% compared to the best monolingual PoT on average.

Our contributions are summarized below:

• We conduct comprehensive experiments of PoTs with different PLs across various reasoning tasks and models, revealing that the choice of PL is dependent on tasks and models.

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

140

141

142

143

144

145

146

147

148

149

150

152

153

154

155

156

- We introduce a task and model agnostic approach called MultiPoT, which integrates multilingual PoTs and leverages strength and diversity across various PLs.
- Experimental results show that MultiPoT outperforms Python Self-Consistency and matches or surpasses the best language of each scenario. On both the model and task averages, MultiPoT enhances performance.

2 Related Work

2.1 Program of Thoughts

In-context learning (Brown et al., 2020; Chowdhery et al., 2023) is a natural ability that emerges when models scale up (Kaplan et al., 2020; Wei et al., 2022b). It appends demonstrations before the question, guiding the LLMs to generate outputs in specific format (Liu et al., 2023; Wei et al., 2022a). CoT (Wei et al., 2022c) is a specific form of incontext learning, whose demonstrations consist of intermediate steps imitating the human thought process. It significantly enhances the model's reasoning capabilities (Yang et al., 2023) but suffers from errors associated with numerical calculations (Madaan and Yazdanbakhsh, 2022). CoT always uses Self-Consistency (Wang et al., 2023c) to increase the probability of getting the correct answer by sampling and voting.

PoT (Chen et al., 2022; Gao et al., 2023) is an extension of CoT to avoid incorrect calculation. It represents intermediate steps as comments and code and executes the entire program with an interpreter to obtain answers. PoT not only excels in reasoning tasks but has rapidly extended to practical applications, including image understanding and robotic control (Surís et al., 2023; Li et al., 2023a). It has become a key method for agents to perform complex reasoning and tool invocation (Yang et al., 2024). It is important to note that all previous PoT work only use Python. For the first time, we are exploring PoTs that use multiple PLs.

2.2 Usage of Multiple Program Languages

The training datasets naturally include a variety of PLs, endowing Code LLMs with the ability to handle multilingual programming (Kocetkov et al.,



Figure 2: Comparative overview of MultiPoT and Self-Consistency. MultiPoT constructs prompts for each PL, ensuring a consistent reasoning process while aslo considering the distinct coding styles. It then integrates these PLs: generating multilingual PoTs based on prompts, executing them to gather results, and finally voting for the answer. In contrast to Self-Consistency's single-language focus, MultiPoT leverages multiple PLs.

2022; Nguyen et al., 2023; Gao et al., 2020; Nijkamp et al., 2023; Chen et al., 2021). This capabil-158 ity extends code tasks like generation, optimization, 159 translation, and repair to other languages beyond 160 Python (Gimeno et al., 2023; Shypula et al., 2023; Zhang et al., 2023; Wu et al., 2023). Despite the progress, current multilingual research (Jin et al., 2023; Joshi et al., 2023; Khare et al., 2023) mainly 164 focuses on code-related tasks, neglecting the potential of PLs as tools to assist in other tasks. Additionally, these studies often treat each language separately without interaction. Our study pioneers the use of multiple PLs in reasoning tasks and introduces a novel integrated approach, leveraging 170 the collective strength and diversity of various PLs to enhance overall performance.

Methodology 3

157

161

163

165

167

171

172

Figure 2 provides an overview of MultiPoT and 174 compares it with Self-Consistency. MultiPoT in-175 volves two key steps. Firstly, it constructs a unique 176 prompt for each PL while ensuring semantic con-177 sistency and formal diversity among them. (Sec-178 tion 3.1). Secondly, it integrates multiple PLs. Specifically, generate PoT in the corresponding 181 language from the prompt of each PL, followed by executing to obtain results and voting for the final answer. (Section 3.2). Unlike Self-Consistency, which relies on a single PL, MultiPoT integrates a range of PLs, utilizing their strength and diversity. 185

3.1 **Multilingual Prompts Construction**

To prompt Code LLMs to generate PoTs with multiple PLs, we first construct demonstrations with PoTs across different PLs. To ensure fairness, all PLs share the same questions in demonstrations. This requires that PoTs in each PL are semantically identical, which means the same reasoning process and execution results. However, there is a great diversity among PLs. The representation of the reasoning process should closely align with the typical styles of each PL in the pre-training data. This is crucial to fully utilize the capabilities of Code LLMs (Wang et al., 2023b). Therefore, when transforming the reasoning process into code, we consider the formal diversity among PLs.

187

188

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

209

210

211

212

213

214

215

We take Built-in Content, Special Syntax, Type Definition, and Varibale Naming as examples. In Figure 2, (1) while Python can directly employ the 'sort' function, C++ has to load it from the 'algorithm' library. Regarding variables, Python's 'list' resembles C++'s 'vector' more than the array. (2) List comprehension like 'sum(1 for penguin in penguins if penguin["age"] < 8)' is a standard syntax in Python. However, a straightforward for-loop is the common practice in other PLs. (3) Static PLs such as C++ require to define the variable type. We carefully define 'int' and 'double' variables to ensure computational accuracy and enhance flexibility by defining 'struct'. (4) We keep the naming styles of each

PL. For instance, Python uses Snake Case, whereas Java favors Camel Case ('secondPenguin'). Appendix A.3 shows the demonstrations.

Based on semantic consistency, we successfully craft multilingual demonstrations exhibiting some formal variations. By adding the question after the demonstration, we get the prompt for each PL.

3.2 Intergration

216

217

218

219

224

226

229

230

232

233

236

237

240

241

242

245

246

247

248

257

259

260

261

263

While Self-Consistency enhances performance by sampling to explore more reasoning paths, it can lead to repeated errors across different samples. In contrast, MultiPoT constructs multilingual prompts and generates PoTs in multiple PLs, significantly increasing the diversity of results.

Specifically, after constructing prompts for each PL, models generate corresponding PoTs, while tracking cumulative probabilities. These probabilities indicate the model's confidence in each answer, with higher probabilities denoting greater confidence. PoTs are then executed and results are collected. The final answer is determined by voting on these results. In cases of tied votes, answers with higher cumulative probabilities are favored. The integration of multiple PLs introduces more potential correct answers and reduces the probability of the same errors in candidate results.

4 Experiment Setup

4.1 Programming Languages

When selecting PLs to compare with Python, we focus on diversity. JavaScript is the most popular language on GitHub (GitHub, 2023) and has less overlap in application with Python, particularly in the ML/AI domains. R has an extensive range of packages, which is similar to Python, but has much less data in pre-training data since it is a niche language. Python, R, and JavaScript are dynamic languages that do not require explicit variable type definitions. To incorporate the diversity of language types, we choose the two most common static languages, Java and C++. The latter is closer to low-level programming and has fewer extension packages. We do not include C due to its high similarity with C++. These five languages offer a diverse range of application scenarios, data volumes, and language types compared to Python.

4.2 Tasks

We select representative and discriminating tasks. We initially select four tasks from Gao et al. (2023): Math Application (Appl.), Date, Tabular and Spatial, and add the task Math. Appl. contains elementary-level mathematical application problems (GSM8K (Cobbe et al., 2021), SVAMP (Patel et al., 2021), Asdiv (Miao et al., 2020)). Date, Tabular, and Spatial are extracted from BBH-Hard (Suzgun et al., 2022) (Date Understanding, Penguins in a Table, Reasoning about Coloured Objects). These tasks assess understanding and reasoning about temporal sequences, structured text, and spatial positioning respectively. Math, consisting of the transformed MATH (Hendrycks et al., 2021) dataset. The difference between Math and Appl. lies in the level of difficulty. Math is more challenging and directly describes the math question without scenarios. These five tasks do not encompass all reasoning tasks and datasets but are discriminating and representative enough to test the capabilities of the model in different PLs. The additional details of the tasks are in the Appendix A.1.

264

265

266

267

269

270

271

272

273

274

275

276

277

278

279

281

282

286

287

289

290

291

292

293

294

295

297

298

299

301

302

303

304

305

306

307

308

309

310

311

312

4.3 Backbone LLMs

As the previously used code-davinci family is no longer accessible, we select four backbone LLMs, including the three strongest Code LLMs: **Starcoder** (15B), **Code Llama** (34B), and **Deepseek Coder** (33B). We select the base versions. The experiments of the Python version are discussed in Section 6.2.5, and the results indicate that they do not affect our conclusions and methodology. **Chat-GPT** is also utilized as a representative of codecapable NL LLMs, invoking through the API of gpt-3.5-turbo. By choosing these backbone LLMs with different sizes and characteristics, we can obtain more realistic and credible results.

4.4 Inference Details

We combine Chen et al. (2022) and Gao et al. (2023)'s prompt templates for few-shot inference. We fix the questions from the previous work and write code in the respective PLs. The number of questions in each task is shown in Appendix A.1. When sampling for Self-Consistency, we follow Chen et al. (2022) and set t = 0.4, $top_p = 1$. For a fair comparison with MultiPoT which integrates five languages, we set k = 5.

5 Results

In this section, we first discover from the results of greedy decoding that Python is not the best language for all tasks and all models. There is no such perfect language. The performance of each



Figure 3: The performance of three models across five tasks in five different programming languages. AVG denotes the average performance of a programming language across all tasks. Each language performance is expressed as a ratio to the highest-performing language for that specific task. The center of the circle represents 50%.

Language		Code LLMs			ChatGPT							
88-	Appl.	Math	Date	Tabular	Spatial	AVG	Appl.	Math	Date	Tabular	Spatial	AVG
Python	58.51	23.62	42.37	83.00	73.87	56.27	80.75	39.74	46.61	94.63	91.70	70.69
R	57.04	22.61	47.70	85.46	71.20	56.80	79.37	34.86	55.01	89.93	92.85	70.40
C++	60.80	22.61	32.79	86.35	75.87	55.68	79.46	39.90	47.70	91.95	86.65	69.13
Java	60.11	23.75	43.81	87.92	75.82	58.28	80.63	42.65	51.22	87.92	86.70	69.82
JavaScript	60.14	24.35	42.82	83.89	71.58	56.56	81.25	36.07	55.01	92.62	90.15	71.02

Table 1: The performance of Code LLMs and ChatGPT for greedy decoding for five languages on five tasks. Code LLMs are the average results for Starcoder, Code Llama, and Deepseek Coder. **AVG** means the average performance of the language on five tasks. **Bold** denotes the highest performance on the task.

PL varies greatly depending on the task and model (Section 5.1). After Self-Consistency, the performance discrepancy still exists. Finally, by integrating multiple languages, MultiPoT significantly outperforms Python. Furthermore, its performance matches or exceeds the best monolingual PoTs in almost all scenarios and achieves improvement on task and model averages (Section 5.2).

313

314

315

316

317

318

319

321

323

326

331

335

To reduce evaluation bias, our analysis of the correlation between PLs and tasks will be based on the averaged results from the three Code LLMs. Due to the significant differences between Chat-GPT and these Code LLMs, we will not include its results in the average calculation.

5.1 Comparison among programming languages

Python is not the optimal language choice. Table 1 shows that the best PL for Code LLMs and ChatGPT on average is Java and JavaScript, respectively. Figure 3 illustrates that Python does not achieve the best performance on any of the tasks for any of the Code LLMs. On Deepseek Coder, Python is even the worst-performing PL on average.

On ChatGPT, although Python performs best on Tabular, it falls short by 2.9% and 8.4% compared to the best PL on Math and Date respectively. The preference for Python among humans may be due to its simple syntax and high readability, but it is a subjective bias that PoT only needs it. Relying on Python leads to a suboptimal outcome.

However, it is important to note that **there is no one-size-fits-all language for all tasks and models**. The gap between PLs is significant when considering each task and model.

The performance of each programming language is task-dependent. *Different tasks are suitable for different languages*. Table 1 indicates that, except for Python on Code LLMs and C++ on Chat-GPT, all PLs excel in at least one task. What's more, on ChatGPT, except for JavaScript, each language also ranks as the least effective in at least one task. *A language that performs exceptionally well in one task might underperform in another.* For instance, R demonstrates superior performance on Date for both Code LLMs and ChatGPT, yet it is the least effective on Appl. and Math.

The performance of each programming lan-

336

337

	ChatGPT						Star	coder				
	Appl.	Math	Date	Table	Spatial	AVG	Appl.	Math	Date	Table	Spatial	AVG
Python	82.31	45.76	47.70	94.63	93.60	72.80	47.04	19.69	34.96	79.19	70.00	50.18
R	80.95	40.61	58.81	93.29	94.60	73.65	44.21	17.74	37.13	77.85	65.90	48.57
C++	81.40	43.77	49.05	93.29	88.45	71.19	47.34	16.74	18.70	82.55	70.95	47.26
Java	81.79	45.33	53.39	92.62	88.80	72.39	47.97	16.76	35.23	78.52	69.50	49.60
JavaScript	82.58	40.64	56.10	96.64	93.30	73.85	48.40	19.15	36.31	80.54	72.95	51.47
MultiPoT	84.33	49.92	58.54	98.66	95.30	77.35	49.67	20.41	40.38	87.25	71.55	53.85
			Code	Llama					Deepse	ek Coder		
Python	68.63	27.95	50.68	92.62	77.55	63.48	70.65	37.64	44.72	93.96	89.80	67.35
R	66.80	26.65	58.27	93.29	79.05	64.81	69.22	33.59	53.12	93.29	92.60	68.36
C++	71.33	24.99	43.36	93.29	80.45	62.68	72.32	33.94	39.57	95.30	93.40	66.91
Java	70.10	27.93	56.91	93.96	81.80	66.14	72.10	35.35	55.56	93.96	88.75	69.14
JavaScript	68.97	26.16	50.41	87.25	80.35	62.63	71.89	35.60	52.57	93.29	86.10	67.89
MultiPoT	71.17	27.97	58.54	93.96	79.60	66.24	72.32	37.55	54.47	95.30	91.70	70.27

Table 2: Self-Consistency and MultiPoT results of four LLMs on five tasks and AVG. Δ represents the percentage improvement of MultiPoT compared to the best monolingual PoT performance.

guage is model-dependent. *Code LLMs and ChatGPT differ significantly.* Table 1 shows that C++, JavaScript, and Java, excel on Appl., Math, and Spatial respectively on Code LLMs, but rank second-to-last on ChatGPT. *Even within Code LLMs, disparities between models are evident.* Figure 3 shows that Code Llama has a clear preference for Java, which keeps the top two ranks across all tasks, yet is not observed on the remaining models. On Deepseek Coder, C++ leads on average, whereas notably ranks last on the other models. R ranks second on Spatial on Deepseek Coder, but the worst on the other two Code LLMs.

These variations demonstrate that different PLs exhibit distinct strength due to complex factors such as task suitability and model preference.

5.2 MultiPoT

361

362

363

364

371

373

374

375

384

387

392

Self-Consistency does not eliminate performance disparities between languages. Table 2 presents the Self-Consistency results. *The inherent strength of different languages persist*. The optimal PL on each scenario is generally consistent with greedy decoding results, except Python emerges as the superior language on Math on all models. Despite Self-Consistency significantly improving the performance, it does not smooth out the performance gap among PLs. *A single language offers limited diversity*. When faced with tasks outside its strength, monolingual samples often make the same mistakes repeatedly, resulting in incorrect answers being chosen through voting.

Different from Self-Consistency relying on a single PL, **MutliPoT** integrates multiple PLs. It

not only **leverages the distinct strength** of each PL, covering more questions (Section 6.2.1) but also **utilizes their greater diversity** to reduce the probability of repeating the same errors.

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

MultiPoT significantly outperforms Python on almost all scenarios. It enhances performance in tasks or models where Python is weak. Across the four models, MultiPoT improves upon Python's performance on Date by at least 15%, and in average (AVG) performance by 4.33% to 7.32%. Furthermore, *MultiPoT also capitalizes on Python's strength.* On Math, where Python excels, MultiPoT also achieves the best results, except in Deepseek Coder, where it slightly trails Python but remains significantly ahead of other languages.

MultiPoT achieves comparable or superior performance to the best monolingual results across all tasks and models. *It is task-agnostic*. It surpasses Self-Consistency on four tasks, ranking second on the remaining task, regardless of whether on Code LLMs average (Table 9) or Chat-GPT. *MultiPoT is also model-agnostic*. It is the top performer across all LLMs on Tabular. On Appl. and Math, MultiPoT surpasses all single-language PoTs across three LLMs. Additionally, on AVG, MultiPoT outperforms the best monolingual result on all four models. Particularly on ChatGPT and Starcoder, it exhibits an improvement of over 4.6%.

The performance of PLs depends on the task and model. Analyzing the interplay of PL, task, and model in practical applications is challenging. Therefore, MultiPoT is a task and model agnostic method of choice due to its consistently high performance across scenarios.



Figure 4: The reasoning ability, code generation ability, and percentage in pre-training data for different languages. Generation lacks data for R. The horizontal coordinates of each model are ranked according to the rise in reasoning performance (excluding R).

6 Discussion

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

6.1 Reasoning Ability of Different Languages

In Section 5.1, we note that the ranking of the average performance of PL varies on each model. The language distribution in the pre-training data of Starcoder and Deepseek Coder offers insights into whether data amount impacts reasoning capabilities. Moreover, we are interested in examining whether code generation and reasoning of multilingual ability are aligned. To assess code generation ability, we utilize the results of each model on the Multilingual HumanEval benchmark, focusing on the pass@1 metric for the four available languages, excluding R due to a lack of evaluation data.

Data distribution influences but does not completely determine reasoning ability. Figure 4 shows the relative relationships among reasoning performance of C++, Python, and Java are consistent with data distribution on Starcoder. However, R demonstrates unexpectedly strong performance, which has an extremely low percentage in both models. C++ has less data amount than Java on Deepseek Coder, but better reasoning performance. This suggests that there are other factors affecting performance besides data distribution.

Code generation abilities do not always align with reasoning abilities. We compare the four languages excluding R in Figure 4. On ChatGPT, the reasoning and code generation abilities of C++, Java, and Python align perfectly. For Deepseek Coder, C++ leads in both abilities. However, an

	StarC.	C. Llama	Deep.C.	GPT
Python	61.03	73.23	75.80	77.62
R	58.86	75.11	76.02	79.00
C++	59.75	72.82	75.80	77.82
Java	61.32	75.62	78.06	78.08
JavaScript	62.60	74.15	76.62	77.65
MultiPoT	64.52	75.71	78.41	83.94

Table 3: The average coverage rate on five tasks of Self-Consistency and MultiPoT on each model.

Stability Metric	Starcoder	Deepseek Coder
Default	53.85	70.27
Length Short	53.36	69.99
Length Long	53.16	69.76
Random	53.71	69.99
Data Amount Little	53.18	70.20
Data Amount Large	53.55	69.43
Δ	0.69	0.84

Table 4: The performance of MultiPoT with different sorting methods. Length Short/Long represents the ascending/descending order according to the length of PoTs, respectively. Δ denotes the range of change.

opposite trend is observed in Deepseek Coder's Python, JavaScript, and Java, where the two abilities diverge significantly. This suggests that while there is consistency, the differences in tasks can lead to substantial variations in results. 458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

6.2 MultiPoT Analysis

6.2.1 Highest Coverage Rate

Unlike the voting mechanism which requires a majority for the correct answer, the coverage rate only needs the answer to appear in PoT results. Table 3 demonstrates that MultiPoT achieves the highest coverage rates on all four models. The monolingual sampling covers less than the multilingual attempts, highlighting that the strength of different PLs exists. MultiPoT effectively utilizes the strength of different PLs and has the highest upper bound.

6.2.2 Stable Performance

When results are tied, the top-ranked result is selected. Different sorting methods reflect the stability. Table 4 shows the performance fluctuation is less than 1% across various sorting criteria, including PoT length, randomness, or data amount from pre-training, compared to the default cumulative probability sorting. This indicates that Multi-PoT consistently selects the correct answer directly, with few instances of ties with incorrect answers. This also suggests a lower probability of different

Model	Method	Appl.	Math	Date	Table	Spatial	AVG
Base	Python MultiPoT	68.63 71.17	27.95 27.95	50.68 58.54	92.62 93.96	77.55 79.60	63.48 66.24
Python	Python MultiPoT	69.54 70.67	28.46 27.46	48.24 55.83	91.28 92.62	74.65 76.70	62.43 64.65

Table 5: The performance of Python Self-Consistency and MultiPoT on Code Llama Base and Code Llama Python.



Figure 5: The impact of the number of integrating PLs. We test the different order of adding languages.

Туре	Starcoder	ChatGPT
All Dynamic	50.41	74.92
Dynamic + Static	51.87	75.77

Table 6: The impact of different language type combinations on MultiPoT. All Dynamic indicates that the three languages are all dynamic, and Dynamic+Static indicates a combination of dynamic and static languages.

PoTs making the same errors.

485

486

487

488

489

490

491

492

493

494

495

496

497

6.2.3 More Languages Better

We investigate the impact of the number of PLs on the performance of MultiPoT. On both Starcoder and Deepseek Coder, we incrementally add languages in both ascending and descending order of data amount. The results demonstrate that the performance of MultiPoT improves with an increasing number of PLs, regardless of the order. This suggests that our method is highly scalable and performance can be further enhanced by incorporating more PLs.

6.2.4 More Language Types Better

498 Python, R, and JavaScript are dynamic languages,
499 while C++ and Java are static. To investigate
500 whether a diverse set of language types enhances
501 MultiPoT's performance, we focus on three PLs.
502 On Starcoder and ChatGPT, JavaScript emerges
503 as the highest-performing dynamic language, sur-

passing Java, which leads between the static languages. Consequently, we integrate JavaScript, Python, and R as All Dynamic and combine Java, Python, and R to represent Dynamic + Static. The results in Table 6 indicate that replacing the higherperforming JavaScript with the lower-performing Java improves performance. This suggests that more language types can provide more diversity to MultiPoT, thereby further enhancing performance. 504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

6.2.5 Python Model Also Works

Our prior experiments with Code LLMs utilize the Base version. However, Code LLMs also have a Python-specific version trained with additional Python corpora. Evaluating MultiPoT on this Python version, as shown in Table 5, we find that Python Self-Consistency improves on Appl. and Math but declines on the other tasks compared to the Base model. Moreover, MultiPoT still outperforms Python Self-Consistency on all tasks except Math, highlighting the adaptability of MultiPoT. Notably, MultiPoT's performance on the Python model is lower across all tasks than on the Base model. This suggests that extensive training on monolingual corpora might diminish the Base model's multilingual abilities on reasoning tasks.

7 Conclusion

Regarding the reliance on Python in PoT, we conduct extensive experiments across various models and tasks using multiple programming languages (Python, R, C++, Java, JavaScript). Our findings demonstrate that Python is not always the best choice, and the optimal language depends on the specific task and the model used. Building on this insight, we introduce MultiPoT, a simple yet effective multilingual integrated method that capitalizes on the strength and diversity of different languages. MultiPoT significantly outperforms Python and achieves comparable or superior performance to the best monolingual outcomes in nearly all scenarios. With its high stability and potential for further expansion, MultiPoT offers a promising avenue for future research.

604 605 606 607 608 609 610 611 612 613 614 615 616 617 618 619 620 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 637 638 639 640 641 642 643 644 645

646

647

648

649

650

651

652

653

654

599

600

601

602

603

Limitations

546

561

562

566

568

569

570

571

573

574

575

576

577 578

579

580

581

583

584

587

588

589 590

591

594

595

598

Our study is comprehensive, but has certain limita-547 tions that we plan to address in future research. Due 548 to computational resource constraints, we confine 549 our experiments to a select number of commonly used programming languages (PLs). While these PLs are representative, they do not encompass the 552 entire spectrum of languages used in programming. 553 Future research could investigate the advantages of incorporating a broader range of programming languages. This may reveal further insights and improve the relevance of our findings.

Ethical Considerations

Our research utilizes publicly available models and datasets with proper citations and adheres to the usage guidelines of ChatGPT, minimizing the risk of generating toxic content due to the widely-used, non-toxic nature of our datasets and prompts.

References

- 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 Advances in Neural Information Processing Systems, volume 33, pages 1877–1901. Curran Associates, Inc.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. Evaluating large language models trained on code.

- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2022. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *CoRR*, abs/2211.12588.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayana Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2023. Palm: Scaling language modeling with pathways. Journal of Machine Learning Research, 24(240):1-113.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. 2020. The pile: An 800gb dataset of diverse text for language modeling.
- Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. Pal: Program-aided language models. In *International Conference on Machine Learning*, pages 10764–10799. PMLR.
- Felix Gimeno, Florent Altché, and Rémi Leblond. 2023. Alphacode 2 technical report. Technical report, AlphaCode Team, Google DeepMind.
- GitHub. 2023. The state of open source and ai. https://github.blog/ 2023-11-08-the-state-of-open-source-and-ai/.
- Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Y. Wu, Y. K. Li, Fuli Luo, Yingfei Xiong, and Wenfeng Liang. 2024. Deepseek-coder: When the large language model meets programming – the rise of code intelligence.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and 656

710

711

657

Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. In *Thirtyfifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2).*

- Matthew Jin, Syed Shahriar, Michele Tufano, Xin Shi, Shuai Lu, Neel Sundaresan, and Alexey Svyatkovskiy. 2023. Inferfix: End-to-end program repair with llms.
- Harshit Joshi, José Cambronero Sanchez, Sumit Gulwani, Vu Le, Gust Verbruggen, and Ivan Radiček. 2023. Repair is nearly generation: Multilingual program repair with llms. 37:5131–5140.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models.
 - Avishree Khare, Saikat Dutta, Ziyang Li, Alaia Solko-Breslin, Rajeev Alur, and Mayur Naik. 2023. Understanding the effectiveness of large language models in detecting security vulnerabilities.
- Denis Kocetkov, Raymond Li, Loubna Ben Allal, Jia Li, Chenghao Mou, Carlos Muñoz Ferrandis, Yacine Jernite, Margaret Mitchell, Sean Hughes, Thomas Wolf, et al. 2022. The stack: 3 tb of permissively licensed source code. *arXiv preprint arXiv:2211.15533*.
- Chengshu Li, Jacky Liang, Fei Xia, Andy Zeng, Sergey Levine, Dorsa Sadigh, Karol Hausman, Xinyun Chen, Li Fei-Fei, and brian ichter. 2023a. Chain of code: Reasoning with a language model-augmented code interpreter. In *NeurIPS 2023 Foundation Models for Decision Making Workshop*.
- Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. 2023b. Starcoder: may the source be with you! *arXiv preprint arXiv:2305.06161*.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pretrain, prompt, and predict: A systematic survey of prompting methods in natural language processing. ACM Comput. Surv., 55(9).
- Aman Madaan and Amir Yazdanbakhsh. 2022. Text and patterns: For effective chain of thought, it takes two to tango. *arXiv preprint arXiv:2209.07686*.
- Shen-Yun Miao, Chao-Chun Liang, and Keh-Yih Su. 2020. A diverse corpus for evaluating and developing english math word problem solvers. In *Proceedings* of the 58th Annual Meeting of the Association for Computational Linguistics, pages 975–984.
- Dung Nguyen, Le Nam, Anh Dau, Anh Nguyen, Khanh Nghiem, Jin Guo, and Nghi Bui. 2023. The vault: A comprehensive multilingual dataset for advancing code understanding and generation. In *Findings* of the Association for Computational Linguistics: EMNLP 2023, pages 4763–4788, Singapore. Association for Computational Linguistics.

Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, and Caiming Xiong. 2023. Codegen: An open large language model for code with multi-turn program synthesis. 713

714

715

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

739

740

741

742

743

744

745

746

747

748

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

- Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. Are nlp models really able to solve simple math word problems? In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2080–2094.
- 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*.
- Alexander Shypula, Aman Madaan, Yimeng Zeng, Uri Alon, Jacob Gardner, Milad Hashemi, Graham Neubig, Parthasarathy Ranganathan, Osbert Bastani, and Amir Yazdanbakhsh. 2023. Learning performanceimproving code edits.
- Dídac Surís, Sachit Menon, and Carl Vondrick. 2023. Vipergpt: Visual inference via python execution for reasoning.
- Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, et al. 2022. Challenging big-bench tasks and whether chain-of-thought can solve them. *arXiv preprint arXiv:2210.09261*.
- Dingzirui Wang, Longxu Dou, Wenbin Zhang, Junyu Zeng, and Wanxiang Che. 2023a. Exploring equation as a better intermediate meaning representation for numerical reasoning.
- Xingyao Wang, Sha Li, and Heng Ji. 2023b. Code4Struct: Code generation for few-shot event structure prediction. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3640– 3663, Toronto, Canada. Association for Computational Linguistics.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023c. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*.
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. 2022a. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations*.

Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022b. Emergent abilities of large language models.

770

772 773

774

776

777

781

783

784

785

787

790

791 792

793

794

796

797

799

800

801

802

- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022c. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.
- Yi Wu, Nan Jiang, Hung Viet Pham, Thibaud Lutellier, Jordan Davis, Lin Tan, Petr Babkin, and Sameena Shah. 2023. How effective are neural networks for fixing security vulnerabilities. In *Proceedings of the* 32nd ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSTA '23. ACM.
 - Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V. Le, Denny Zhou, and Xinyun Chen. 2023. Large language models as optimizers.
- Ke Yang, Jiateng Liu, John Wu, Chaoqi Yang, Yi R. Fung, Sha Li, Zixuan Huang, Xu Cao, Xingyao Wang, Yiquan Wang, Heng Ji, and Chengxiang Zhai. 2024. If Ilm is the wizard, then code is the wand: A survey on how code empowers large language models to serve as intelligent agents.
- Jiyang Zhang, Pengyu Nie, Junyi Jessy Li, and Milos Gligoric. 2023. Multilingual code co-evolution using large language models. In *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, ESEC/FSE 2023, page 695–707, New York, NY, USA. Association for Computing Machinery.

806

807

810

811

812

813

814

815

816

817

818

819

820

821

822

823

A Appendix

A.1 Tasks

Subset	#Original	#Filtered
Algebra	1,187	1,068
Counting & Probability	474	474
Geometry	479	466
Intermediate Algebra	903	721
Number Theory	540	528
Prealgebra	871	842
Precalculus	546	370
SUM	5,000	4,469

Table 7: After filtering, the statistics of MATH dataset.

In MATH, the answers to the problems are expressed using LaTeX. It's too hard to construct prompts in other languages that meet all the requirements, we select those that can be calculated to a single number, excluding problems with interval or formula-based answers.

Task	#Data	#Shots
Appl.	4,415	3
Math	4,469	3
Date	369	6
Tabular	149	3
Spatial	2,000	3

Table 8: Summarization of selected reasoning tasks.

Here are the details of our selected tasks, including the number of questions in each task (#Data) and the number of shots in demonstrations. Appl. comprises the GSM8K (Cobbe et al., 2021), SVAMP (Patel et al., 2021), and Asdiv (Miao et al., 2020) datasets. These datasets contain elementarylevel math problems set in specific application scenarios, focusing on mathematical abstraction and modeling skills, with relatively low difficulty. Since they are the same type of questions, we merge them in one task.

A.2 Results

	Appl.	Math	Date	Table	Spatial
Python	62.11	28.43	43.45	88.59	79.12
R	60.08	25.99	49.50	88.14	79.18
C++	63.66	25.23	33.88	90.38	81.60
Java	63.39	26.68	49.23	88.81	80.02
JavaScript	63.09	26.97	46.43	87.02	79.80
MultiPoT	64.39	28.64	51.13	92.17	80.95

Table 9: The average performance of three Code LLMs for Self-Consistency and MultiPoT in each task.

Table 9 shows that on the average performance824of three Code LLMs, MultiPoT surpasses all Self-825Consistency on four tasks, and is only lower826slightly than C++ on Spatial.827

828

829

830

831

A.3 Prompts

Here are our multilingual prompts. We show prompts of Tabular(3-shots) as an example and prompts for other tasks are in the released code. Question: Here is a table where the first line is a header and each subsequent line is a penguin: name, age, height (cm), weight (kg) Louis, 7, 50, 11 Bernard, 5, 80, 13 Vincent, 9, 60, 11 Gwen, 8, 70, 15 For example: the age of Louis is 7, the weight of Gwen is 15 kg, the height of Bernard is 80 cm. We now add a penguin to the table: James, 12, 90, 12 How many penguins are less than 8 years old? Use the Python program to solve the problem. The reasoning progress is clearly shown in the program. # Define the table of penguins with their attributes penguins = [{"name": "Louis", "age": 7, "height": 50, "weight": 11}, {"name": "Bernard", "age": 5, "height": 60, "weight": 11}, {"name": "Gwen", "age": 8, "height": 70, "weight": 15}, {"name": "James", "age": 12, "height": 90, "weight": 12}] # Count the number of penguins that are less than 8 years old penguins_less_than_8_years = sum(1 for penguin in penguins if penguin["age"] < 8) ans = penguins_less_than_8_years

Figure 6: Python Prompt of the first question.

```
Question: Here is a table where the first line is a header and each subsequent line is a penguin:
name, age, height (cm), weight (kg)
Louis, 7, 50, 11
Bernard, 5, 80, 13
Vincent, 9, 60, 11
Gwen, 8, 70, 15
For example: the age of Louis is 7, the weight of Gwen is 15 kg, the height of Bernard is 80 cm.
Which is the youngest penguin?
Use the Python program to solve the problem. The reasoning progress is clearly shown in the program.
# Define the data for the penguins
penguins = [
        {"name": "Louis", "age": 7, "height": 50, "weight": 11},
        {"name": "Vincent", "age": 9, "height": 80, "weight": 13},
        {"name": "Gwen", "age": 8, "height": 70, "weight": 15}
]
# Find the youngest penguin
youngest_penguin = min(penguins, key=lambda p: p["age"])
ans = youngest_penguin["name"]
```



Question: Here is a table where the first line is a header and each subsequent line is a penguin: name, age, height (cm), weight (kg) Louis, 7, 50, 11 Bernard, 5, 80, 13 Vincent, 9, 60, 11 Gwen, 8, 70, 15 For example: the age of Louis is 7, the weight of Gwen is 15 kg, the height of Bernard is 80 cm. What is the name of the second penguin sorted by alphabetic order? Use the Python program to solve the problem. The reasoning progress is clearly shown in the program. # Creating a list of dictionaries to represent each penguin and their attributes penguins = [{"name": "Louis", "age": 7, "height": 50, "weight": 11}, {"name": "Uncent", "age": 5, "height": 80, "weight": 13}, {"name": "Gwen", "age": 8, "height": 70, "weight": 15}] # Sorting the list of penguins by their names in alphabetical order sorted_penguins = sorted(penguins, key=lambda p: p["name"]) # Extracting the name of the second penguin in the sorted list ans = sorted_penguins[1]["name"]

Figure 8: Python Prompt of the third question.

```
Question: Here is a table where the first line is a header and each subsequent line is a penguin:
Additional and the second seco
 For example: the age of Louis is 7, the weight of Gwen is 15 kg, the height of Bernard is 80 cm.
We now add a penguin to the table:
James, 12, 90, 12
How many penguins are less than 8 years old?
Use the C++ program to solve the problem. The reasoning progress is clearly shown in the program.
#include <iostream>
#include <vector>
using namespace std;
// Define a structure for storing penguin data
 struct Penguin {
              string name;
              int age;
              int height;
              int weight;
};
 int main() {
              // Create a vector of Penguin structures
              vector<Penguin> penguins = {
                           {"Louis", 7, 50, 11},
{"Bernard", 5, 80, 13},
{"Vincent", 9, 60, 11},
{"Gwen", 8, 70, 15},
{"James", 12, 90, 12}
              }:
              // Count penguins less than 8 years old
              int count = 0:
              for (const auto& penguin : penguins) {
                            if (penguin.age < 8) {</pre>
                                         count++;
                            }
              int ans = count;
              // Print the result
              cout << ans << endl;</pre>
              return 0:
}
```

Figure 9: C++ Prompt of the first question.

```
Question: Here is a table where the first line is a header and each subsequent line is a penguin:
name, age, height (cm), weight (kg)
Louis, 7, 50, 11
Bernard, 5, 80, 13
Vincent, 9, 60, 11
Gwen, 8, 70, 15
For example: the age of Louis is 7, the weight of Gwen is 15 kg, the height of Bernard is 80 cm.
Which is the youngest penguin?
Use the C++ program to solve the problem. The reasoning progress is clearly shown in the program.
#include <iostream>
#include <string>
#include <vector>
using namespace std;
struct Penguin {
     string name;
     int age;
     int height;
     int weight;
};
int main() {
     // Create a list of penguins
     vector<Penguin> penguins = {
         {"Louis", 7, 50, 11},
{"Bernard", 5, 80, 13},
{"Vincent", 9, 60, 11},
{"Gwen", 8, 70, 15}
     };
     // Initialize variables to track the youngest penguin
     Penguin youngest = penguins[0];
     // Iterate through the list to find the youngest penguin
     for (const Penguin& penguin : penguins) {
          if (penguin.age < youngest.age) {</pre>
               youngest = penguin;
          }
     }
     string ans = youngest.name;
     // Output the name of the youngest penguin
     cout << ans << endl;</pre>
```

```
return 0;
}
```

Figure 10: C++ Prompt of the second question.

```
Question: Here is a table where the first line is a header and each subsequent line is a penguin:
name, age, height (cm), weight (kg)
Louis, 7, 50, 11
Bernard, 5, 80, 13
Vincent, 9, 60, 11
Gwen, 8, 70, 15
For example: the age of Louis is 7, the weight of Gwen is 15 kg, the height of Bernard is 80 cm.
What is the name of the second penguin sorted by alphabetic order?
Use the C++ program to solve the problem. The reasoning progress is clearly shown in the program.
#include <iostream>
#include <vector>
#include <algorithm>
using namespace std;
struct Penguin {
     string name;
     int age;
     int height;
     int weight;
};
int main() {
     // Creating a list of penguins
     vector<Penguin> penguins = {
           ("Louis", 7, 50, 11},
{"Bernard", 5, 80, 13},
{"Vincent", 9, 60, 11},
{"Gwen", 8, 70, 15}
     };
     // Sorting the penguins by name using a simpler lambda function
     sort(penguins.begin(), penguins.end(), [](const Penguin &a, const Penguin &b) {
          return a name < b name;</pre>
     });
     // Displaying the name of the second penguin after sorting
     string ans = penguins[1].name;
     cout << ans << endl;</pre>
     return 0;
}
```

Figure 11: C++ Prompt of the third question.

```
Question: Here is a table where the first line is a header and each subsequent line is a penguin:
name, age, height (cm), weight (kg)
Louis, 7, 50, 11
Bernard, 5, 80, 13
Vincent, 9, 60, 11
Gwen, 8, 70, 15
For example: the age of Louis is 7, the weight of Gwen is 15 kg, the height of Bernard is 80 cm.
We now add a penguin to the table:
James, 12, 90, 12
How many penguins are less than 8 years old?
Use the Java program to solve the problem. The reasoning progress is clearly shown in the program.
import java.util.ArrayList;
import java.util.List;
public class PenguinAgeCounter {
      public static void main(String[] args) {
           List<Penguin> penguins = new ArrayList<>();
           // Add penguins to the list
           // Add penguins to the tist
penguins.add(new Penguin("Louis", 7, 50, 11));
penguins.add(new Penguin("Bernard", 5, 80, 13));
penguins.add(new Penguin("Vincent", 9, 60, 11));
penguins.add(new Penguin("Gwen", 8, 70, 15));
penguins.add(new Penguin("James", 12, 90, 12));
           // Count the number of penguins less than 8 years old
           int count = 0;
           for (Penguin penguin : penguins) {
    if (penguin.age < 8) {</pre>
                      count++;
                 }
           }
            int ans = count;
           System.out.println(ans);
      }
}
class Penguin {
      String name;
      int age;
      int height; // in cm
int weight; // in kg
      public Penguin(String name, int age, int height, int weight) {
           this.name = name;
           this.age = age;
           this.height = height;
           this.weight = weight;
      }
}
```

Figure 12: Java Prompt of the first question.

```
Question: Here is a table where the first line is a header and each subsequent line is a penguin:
uuestion: mere is a table where the first line is a header and each subsequent line is a penguin
name, age, height (cm), weight (kg)
Louis, 7, 50, 11
Bernard, 5, 80, 13
Vincent, 9, 60, 11
Gwen, 8, 70, 15
For example: the age of Louis is 7, the weight of Gwen is 15 kg, the height of Bernard is 80 cm.
Which is the vuergest penguin?
Which is the youngest penguin?
Use the Java program to solve the problem. The reasoning progress is clearly shown in the program.
import java.util.ArrayList;
import java.util.List;
public class PenguinAgeFinder {
     public static void main(String[] args) {
          // Create a list of Penguin objects
          List<Penguin> penguins = new ArrayList<>();
          penguins.add(new Penguin("Louis", 7, 50, 11));
penguins.add(new Penguin("Bernard", 5, 80, 13));
penguins.add(new Penguin("Vincent", 9, 60, 11));
          penguins.add(new Penguin("Gwen", 8, 70, 15));
           // Call the method to find the youngest penguin
          Penguin youngest = findYoungestPenguin(penguins);
          // Print the name of the youngest penguin
          String ans = youngest.name;
          System.out.println(ans);
     }
     // Method to find the youngest penguin
     public static Penguin findYoungestPenguin(List<Penguin> penguins) {
          Penguin youngest = penguins.get(0);
          for (Penguin penguin : penguins) {
                if (penguin.age < youngest.age) {</pre>
                    youngest = penguin;
                }
          }
          return youngest;
     }
}
// Define the Penguin class with relevant attributes
class Penguin {
     String name;
     int age;
     int height; // in cm
int weight; // in kg
      // Constructor to initialize a Penguin object
     Penguin(String name, int age, int height, int weight) {
           this.name = name;
          this.age = age;
          this.height = height;
          this.weight = weight;
     }
}
```

Figure 13: Java Prompt of the second question.

```
Question: Here is a table where the first line is a header and each subsequent line is a penguin:
Question: mere is a table where the
name, age, height (cm), weight (kg)
Louis, 7, 50, 11
Bernard, 5, 80, 13
Vincent, 9, 60, 11
Gwen, 8, 70, 15
For example: the age of Louis is 7
For example: the age of Louis is 7, the weight of Gwen is 15 kg, the height of Bernard is 80 cm.
What is the name of the second penguin sorted by alphabetic order?
Use the Java program to solve the problem. The reasoning progress is clearly shown in the program.
import java.util.*;
public class Main {
     public static void main(String[] args) {
          // Step 2: Create Penguin objects
          // Step 2: Create Penguin objects
Penguin louis = new Penguin("Louis", 7, 50, 11);
Penguin bernard = new Penguin("Bernard", 5, 80, 13);
Penguin vincent = new Penguin("Vincent", 9, 60, 11);
          Penguin gwen = new Penguin("Gwen", 8, 70, 15);
          // Step 3: Add Penguins to a list
          List<Penguin> penguins = new ArrayList<>();
penguins.add(louis);
          penguins.add(bernard);
          penguins.add(vincent);
          penguins.add(gwen);
          // Step 4: Sort the list by name
          Collections.sort(penguins, Comparator.comparing(Penguin::getName));
          // Step 5: Find and print the name of the second penguin
          Penguin secondPenguin = penguins.get(1); // Lists are zero-indexed
          String ans = secondPenguin.getName();
          System.out.println(ans);
     }
}
// Step 1: Define the Penguin class
class Penguin {
     private String name;
     private int age;
     private int height; // in cm
     private int weight; // in kg
     // Constructor
     public Penguin(String name, int age, int height, int weight) {
          this.name = name;
this.age = age;
          this.height = height;
          this.weight = weight;
     }
     // Getters
     public String getName() {
         return name;
     }
}
```

Figure 14: Java Prompt of the third question.

```
Question: Here is a table where the first line is a header and each subsequent line is a penguin:
Louis, 7, 50, 11
Bernard, 5, 80, 13
Vincent, 9, 60, 11
Gwen, 8, 70, 15
For example: the age of Louis is 7, the weight of Gwen is 15 kg, the height of Bernard is 80 cm.
We now add a penguin to the table:
James, 12, 90, 12
How many penguins are less than 8 years old?
Use the Javascript program to solve the problem. The reasoning progress is clearly shown in the program.
// Step 1: Create an array of penguin objects
let penguins = [
     { name: "Louis", age: 7, height: 50, weight: 11 },
{ name: "Bernard", age: 5, height: 80, weight: 13 },
{ name: "Vincent", age: 9, height: 60, weight: 11 },
     { name: "Gwen", age: 8, height: 70, weight: 15 }
];
// Step 2: Add James to the array
penguins.push({ name: "James", age: 12, height: 90, weight: 12 });
// Step 3: Initialize a counter
let count = 0;
// Step 4: Iterate over the array
for (let penguin of penguins) {
    // Step 5: Check if the penguin's age is less than 8
     if (penguin.age < 8) {</pre>
         count++;
     }
}
let ans = count;
```

Figure 15: JavaScript Prompt of the first question.

```
Question: Here is a table where the first line is a header and each subsequent line is a penguin:
name, age, height (cm), weight (kg)
Louis, 7, 50, 11
Bernard, 5, 80, 13
Vincent, 9, 60, 11
Gwen, 8, 70, 15
For example: the age of Louis is 7, the weight of Gwen is 15 kg, the height of Bernard is 80 cm.
Which is the youngest penguin?
Use the Javascript program to solve the problem. The reasoning progress is clearly shown in the program.
// Step 1: Create an array of penguins
const penguins = [
     { name: "Louis"
     { name: "Louis", age: 7, height: 50, weight: 11 },
{ name: "Bernard", age: 5, height: 80, weight: 13 },
{ name: "Vincent", age: 9, height: 60, weight: 11 },
     { name: "Gwen", age: 8, height: 70, weight: 15 }
];
// Step 2: Find the youngest penguin
let youngest = penguins[0];
for (let i = 1; i < penguins.length; i++) {</pre>
    if (penguins[i].age < youngest.age) {
    youngest = penguins[i];</pre>
     }
}
ans = youngest.name;
```



```
Question: Here is a table where the first line is a header and each subsequent line is a penguin:
Louis, 7, 50, 11
Bernard, 5, 80, 13
Vincent, 9, 60, 11
Gwen, 8, 70, 15
For example: the age of Louis is 7, the weight of Gwen is 15 kg, the height of Bernard is 80 cm.
What is the name of the second penguin sorted by alphabetic order?
Use the Javascript program to solve the problem. The reasoning progress is clearly shown in the program.
// Define the array of penguins
let penguins = [
     { name: "Louis", age: 7, height: 50, weight: 11 },
{ name: "Bernard", age: 5, height: 80, weight: 13 },
{ name: "Vincent", age: 9, height: 60, weight: 11 },
     { name: "Gwen", age: 8, height: 70, weight: 15 }
];
// Sort the penguins array by the name property
penguins.sort((a, b) => {
     if (a.name < b.name) {</pre>
         return -1;
     }
     if (a.name > b.name) {
          return 1;
     }
     return 0;
});
// The name of the second penguin after sorting
let ans = penguins[1].name;
```

Figure 17: JavaScript Prompt of the third question.

```
Question: Here is a table where the first line is a header and each subsequent line is a penguin:
name, age, height (cm), weight (kg)
Louis, 7, 50, 11
Bernard, 5, 80, 13
Vincent, 9, 60, 11
Gwen, 8, 70, 15
For example: the age of Louis is 7, the weight of Gwen is 15 kg, the height of Bernard is 80 cm.
We now add a penguin to the table:
James, 12, 90, 12
How many penguins are less than 8 years old?
Use the R program to solve the problem. The reasoning progress is clearly shown in the program.
# Create the data frame
penguins <- data.frame(
name = c("Louis", "Bernard", "Vincent", "Gwen", "James"),
age = c(7, 5, 9, 8, 12),
height = c(50, 80, 60, 70, 90),
weight = c(11, 13, 11, 15, 12)
)
# Filter and count the penguins younger than 8 years
num_penguins_younger_than_8 <- nrow(subset(penguins, age < 8))
ans = num_penguins_younger_than_8
```

Figure 18: R Prompt of the first question.

```
Question: Here is a table where the first line is a header and each subsequent line is a penguin:
name, age, height (cm), weight (kg)
Louis, 7, 50, 11
Bernard, 5, 80, 13
Vincent, 9, 60, 11
Gwen, 8, 70, 15
For example: the age of Louis is 7, the weight of Gwen is 15 kg, the height of Bernard is 80 cm.
Which is the youngest penguin?
Use the R program to solve the problem. The reasoning progress is clearly shown in the program.
# Create a data frame representing the penguins
penguins <- data.frame(
    name = c("Louis", "Bernard", "Vincent", "Gwen"),
    age = c(7, 5, 9, 8),
    height_cm = c(50, 80, 60, 70),
    weight_kg = c(11, 13, 11, 15)
)
# Find the youngest penguin by locating the minimum age
youngest_index <- which.min(penguins$age)
ans = penguins$name[youngest_index]
```

Figure 19: R Prompt of the second question.

```
Question: Here is a table where the first line is a header and each subsequent line is a penguin:
name, age, height (cm), weight (kg)
Louis, 7, 50, 11
Bernard, 5, 80, 13
Vincent, 9, 60, 11
Gwen, 8, 70, 15
For example: the age of Louis is 7, the weight of Gwen is 15 kg, the height of Bernard is 80 cm.
What is the name of the second penguin sorted by alphabetic order?
Use the R program to solve the problem. The reasoning progress is clearly shown in the program.
# Create a data frame with the penguins' information
penguins <- data.frame(
    name = c("Louis", "Bernard", "Vincent", "Gwen"),
    age = c(7, 5, 9, 8),
    height = c(15, 80, 60, 70),
    weight = c(11, 13, 11, 15)
)
# Sort the data frame by the 'name' column
sorted_penguins <- penguins[order(penguins$name),]
# Extract the name of the second penguin in the sorted list
ans <- sorted_penguins$name[2]</pre>
```

