# **SCIAGENT: Tool-augmented Language Models for Scientific Reasoning**

**Anonymous ACL submission** 

# Abstract

Scientific reasoning poses an excessive challenge for even the most advanced Large Language Models (LLMs). To make this task more practical and solvable for LLMs, we introduce 005 a new task setting named tool-augmented sci-006 entific reasoning. This setting supplements LLMs with scalable toolsets, and shifts the focus from pursuing an omniscient problem 009 solver to a proficient tool-user. To facilitate the research of such setting, we construct a tool-augmented training corpus named MATH-012 FUNC which encompasses over 30,000 samples and roughly 6,000 tools. Building on MATH-FUNC, we develop SCIAGENT to retrieve, understand and, if necessary, use tools for scientific problem solving. Additionally, we craft a benchmark, SCITOOLBENCH, spanning five 017 scientific domains to evaluate LLMs' abilities with tool assistance. Extensive experiments on SCITOOLBENCH confirm the effectiveness of 021 SCIAGENT. Notably, SCIAGENT-MISTRAL-7B surpasses other LLMs with the same size by more than 13% in absolute accuracy. Fur-024 thermore, SCIAGENT-DEEPMATH-7B shows much superior performance than ChatGPT.

# 1 Introduction

026

027

Scientific reasoning (Ouyang et al., 2023; Zhao et al., 2023) aims to comprehend and make decisions regarding problems among STEM (Science, Technology, Engineering and Mathematics) domains. It is a fundamental aspect of intelligence, a demanding capability of Large Language Models (LLMs), and a notoriously challenging task. For instance, even GPT-4 (OpenAI, 2023) achieves only 50% and 35% accuracy on TheoremQA (Chen et al., 2023b) and SciBench (Wang et al., 2023b), respectively. Regarding open-source LLMs such as LLaMA-2 (Touvron et al., 2023) and CodeLlama (Rozière et al., 2023), their performances are only about 10% accuracy or even less.



Figure 1: Two paradigms for scientific reasoning. Different colors represent different scientific domains. **Left:** Collecting annotations and fine-tuning LLMs domain by domain. **Right:** Our proposed *tool-augmented* setting. LLMs are fine-tuned on math-related, tool-augmented samples (color in red). When adapting LLMs to a specific domain, a pluggable and domain-specific toolset is attached. No additional fine-tuning is further required.

The challenge in scientific reasoning arises from the need for both mathematical (math) and domainspecific reasoning abilities. To address the physical problem in Figure 3, for example, it is necessary to both understand *Malus' law* (domain knowledge) for analyzing the intensity of polarized light, and possess quantitative ability for calculating the light intensity ratios. A natural approach involves collecting annotations and fine-tuning LLMs to enhance their math and domain-specific reasoning abilities, as depicted in Figure 1 (left). However, annotating scientific reasoning problems is extremely expensive. What is worse, adapting LLMs to a new domain demands a fresh round of annotation and fine-tuning, rendering this approach impractical.

043

044

045

046

047

050

051

054

056

057

059

060

061

062

063

In this paper, we draw inspirations from *tool learning* (Qin et al., 2023a) to enhance LLMs' scientific reasoning capabilities. Instead of solving scientific problem from scratch, humans have summarized and wrapped various points as generalized and well-documented functions in scientific computing softwares, such as Matlab, WolframAlpha, SymPy, *etc.* These functions<sup>1</sup>, which could be

<sup>&</sup>lt;sup>1</sup>In this work, tools refer to Python functions. We use tools and functions interchangeably unless otherwise specified.

equivalently viewed as external tools, greatly facil-064 itate math-adept users to solve difficult scientific 065 problems. In analogy with humans, we do not pur-066 sue an omniscient solver across various scientific domains. Instead, we assume the access to domainspecific toolsets and purse a unified, generalized LLM-based tool-user as shown in the Figure 1 (right). This approach tackles domain-specific reasoning challenges by enabling LLMs learn to use a reusable and scalable toolkit. It alleviates the reasoning challenges of LLMs by concentrating solely on enhancing their tool-use abilities. These abilities are not only easier to acquire but also applicable across a variety of scientific fields. By 077 attaching domain-specific toolsets, our tool-users can be readily adapted to different fields without the need for additional in-domain fine-tuning.

> This work focuses on developing and benchmarking the ability of LLMs in scientific reasoning with the help of tools. We envision a scenario where LLMs have access to a domain-specific toolset, comprising various specialized functions. Upon this scenario, we propose a complete framework of dataset construction, model training and evaluation. Given a scientific question, LLMs are supposed to retrieve functions from the toolset and optionally incorporate functions into the formulated solution. We employ an automatic pipeline featuring GPT-4 to compile a large-scale, mathrelated, tool-augmented training corpus named as MATHFUNC. This corpus is designed to enable LLMs to learn both essential math skills and how to retrieve, understand and use functions properly. As a result, MATHFUNC contains 31,375 samples and equipped with a toolset encompassing 5,981 generalized and well-documented functions. We detail this training corpus in Section 3.

094

100

102

104

105

107

108

109

110

111

We fine-tune open-source LLMs on MATHFUNC to develop tool-augmented agents named SCIA-GENT detailed in Section 4. As shown in Figure 3, SCIAGENT firstly generate a *high-level planning* in response to a given question. The agents then use this plan, along with the question, to retrieve functions from the given toolset. Leveraging these retrieved functions, the agents further complete the *low-level action* integrating natural language and Python code. Finally the agents execute the code to complete the problem at hand.

112To benchmark the tool-use abilities in scientific113reasoning, we develop a new benchmark named114SCITOOLBENCH as described in Section 5. Build-115ing upon TheoremQA (Chen et al., 2023b) and

SciBench (Wang et al., 2023b), it has 856 questions covering five domains: Mathematics, Physical, Chemistry, EECS, and Finance. It also contains five domain-specific toolsets comprising a total of 2,446 functions. We evaluate SCIAGENT on SCITOOLBENCH and another benchmark derived from CREATOR-challenge (Qian et al., 2023). Experimental results demonstrate that our agents present remarkable scientific reasoning capabilities. Notably, SCIAGENT-MISTRAL-7B surpasses the best comparable open-source LLMs by an absolute 13.4% accuracy, and SCIAGENT-DEEPMATH-7B outperforms ChatGPT by a large margin. We also conduct an extensive analysis of the benefits and limitations of SCIAGENT series, providing valuable insights for future research.

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

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

158

159

161

162

164

# 2 Preliminary

Related Work. Current methods (Chen et al., 2023b; Xu et al., 2023b; Ouyang et al., 2023), especially those based on open-source LLMs, perform far from satisfactory on scientific reasoning benchmarks (Chen et al., 2023b; Wang et al., 2023b). We attribute it to the scarcity of annotated samples across diverse scientific domains. As a comparison, LLMs present much more remarkable performance on math problems (Yue et al., 2023b; Gou et al., 2023b; Azerbayev et al., 2023) due to the abundant training corpora and/or annotations. Different from concurrent work (Zhang et al., 2024) which collects physics and chemistry annotations, we do not pursue a problem-solver on some specific scientific domains. Instead, we consider to develop a generalized tool-user being proficient on solving diverse scientific problems with the aid of tools. Following previous work on math domain (Qian et al., 2023; Cai et al., 2023; Yuan et al., 2023a), the tools here refer to Python functions. Please see more detailed literature review in Appendix A.

**Task Formulation.** Given a scientific domain D (e.g., physics), tool-augmented scientific reasoning task assumes access to (1) a question  $q \in D$  and (2) a toolset  $F_D$ .  $F_D$  encompasses large amounts of well-documented, domain-specific functions  $\{f_1, ..., f_m\}$ . Our objective is to develop an agent  $\mathcal{M}$  which selectively use functions in  $F_D$  to enhance the answering for the question q.

# **3** Training Corpus: MATHFUNC

To our best knowledge, there are no readily available tool-augmented datasets in scientific reason-



Figure 2: Automatic pipeline for MATHFUNC construction. Please view it starting from the bottom left corner and proceed clockwise. We disentangle the constructions of toolset (dashed lines) and function-augmented samples (solid lines) for more generalized annotations. We do not visualize the function-free samples for simplicity.

ing domains. Therefore, we construct a corpus named MATHFUNC teaching LLMs to better understand and use functions. MATHFUNC is composed of (1) a toolset  $F^2$  including 5,981 generalized, well-documented, math-related functions and (2) a dataset D encompassing 31,375 samples in which solutions call the function from the toolset if necessary (*e.g.*, ④ in Figure 2). We build this corpus based on MATH (Hendrycks et al., 2021b) training set because we expect to teach LLMs both math skills and tool-use abilities.

165

166

168

169

171

173

174

175

176

181

182

183

185

188

189

190

193

194

196

197

198

**Sample Format.** Each sample is a quintuple  $(q, G_q, F_q, S_q, a_q)$ . Here q is a question,  $G_q$  is the planning,  $F_q$  is the function set filtered from the toolset  $(F_q \subset F, |F_q| \ll |F|)$ ,  $S_q$  is the solution and  $a_q$  is the answer.  $S_q$  interleaves rationales  $E_q^3$  and programs  $P_q$  which optionally call functions in  $F_q$  to facilitate the problem solving.

We employ an automatic pipeline to construct MATHFUNC. We illustrate the pipeline in Figure 2 and detail the process in the following subsections.

### 3.1 Planning and Toolset Construction

This module is depicted in the top-left side of Figure 2. Given a question q and its ground-truth solution (written in pure natural language) in MATH training set, we ask GPT-4 to generate (1) a highlevel planning  $G_q$  to analyze this question, (2) one or more well-documented functions  $\tilde{F}_q$  and (3) a solution  $\tilde{S}_q$  calling the functions above. The prompt used is shown in Appendix F.1. In the prompt, we emphasize that the functions should be as **composable and generalized** as possible. Specifically, we do not hope that each question generates only one ad-hoc function (which could only be used by this question). Instead, we expect GPT-4 to generate functions that follow the points in the planning  $G_q$ and can be reused by other questions. Following previous work (Qian et al., 2023; Pan et al., 2023), we provide the error feedback to GPT-4 if the solutions fail to execute, and ask GPT-4 to rectify the errors in  $\tilde{F}_q$  or  $\tilde{S}_q$ . We repeat this procedure until successful execution or reaching maximum loop limitation. The prompt used for rectification is shown in Appendix F.2.

We collect  $G_q$  (① in Figure 2, the same below) and add  $\tilde{F}_q$  to the toolset (②) for question q if the rectified solution  $\tilde{S}_q$  leads to the correct answer  $\tilde{a}_q$ . Regarding the toolset, it is iterated on all questions and finally accumulated as below:

$$F = \bigcup_{q \in D} \tilde{F}_q \cdot I(\tilde{a}_q \text{ is correct})$$
 214

201

202

203

204

206

209

210

211

212

213

215

216

217

218

219

220

221

224

227

228

229

231

232

233

### 3.2 Function-augmented Solutions

To collect function-augmented solution  $S_q$  and  $F_q$ , a natural idea is to directly use the  $\tilde{S}_q$  and  $\tilde{F}_q$  generated above. However, we find that  $\tilde{S}_q$  tends to be contrived and specifically tailored to fit the requirements of function-calling. Moreover, some functions in  $F_q$  tend to be ad-hoc<sup>4</sup>. For examples, the function f(x, y) in Figure 2 merely parameterizes the hyperbola for a specific question. Therefore we disentangle the construction of toolset and function-augmented solutions. Given the developed toolset, we design a cross-retrieval strategy to retrieve more generalized functions  $F_q$  and generate more qualified solutions  $S_q$ . Specifically, we remove  $F_q$  from F temporarily and then retrieve new functions  $F_q \subseteq (F \setminus \tilde{F}_q)$  for question q. This strategy eliminates the likelihood of calling ad-hoc functions from  $\tilde{F}_q$  in  $S_q$ . See examples of retrieved

<sup>&</sup>lt;sup>2</sup>We remove the domain-specific subscript D for expression simplicity. The same below.

<sup>&</sup>lt;sup>3</sup>Here  $E_q$  is written in natural language but formatted as the annotation lines in the program.

<sup>&</sup>lt;sup>4</sup>Despite we instruct GPT-4 to avoid generating ad-hoc functions, there are still some ad-hoc functions in  $\tilde{F}_q$ 



Figure 3: The model architecture of SCIAGENT. Given a domain-specific toolset  $\checkmark$ , our agent answers the question through four consecutive modules. (1) **Planning** : provides a high-level plan for this problem. (2) **Retrieval**  $\bigcirc$ : retrieves related functions from attached toolset. (3) **Action** : generates a low-level solution interleaving rationale and program. The program uses the retrieved functions if necessary. (4) **Execution** : calls Python executor to run the program and outputs the final answer. Not included in this figure for simplicity.

functions, all of which are derived from other questions, in the right side of Figure 2.

**Retriever.** The cross-retrieval strategy necessities a retriever because it is impractical to enumerate thousands of functions in  $F \setminus \tilde{F}_q$ . We train a dense retriever R (③ in Figure 2). We concatenate the question q and the generated planning  $G_g$  as the query, and view the generated functions  $\tilde{F}_q$  as the keys. See details about R in Appendix B.1.

**Solution Generation.** Upon the toolset F and the retriever R, we retrieve three functions as  $F_q$ :

$$F_q = R([q, G_q]; F \setminus F_q)$$

Then we employ GPT-4 to write solutions which optionally call functions in  $F_q$  to generate the solution  $S_q$  (④). The prompt used is illustrated in Appendix F.3. We explicitly point out in the prompt that  $f \in F_q$  should be called if and only if when they do lower the difficulty of problem solving. It mitigates the over-exploitation of function calling in  $S_q$  and increases the robustness of models finetuned on these samples. Specifically, we firstly use GPT-4 with greedy decoding to generate solutions. For those failing to yield correct answers, we further apply nucleus sampling (Holtzman et al., 2020) with 5 repeat times and 0.6 temperature. We filter wrong solutions and collect remaining 6,229 samples as our function-augmented solutions.

In parallel, we use GPT-4 to generate functionfree solutions. Though not indispensable, we expect them to further enhance the math reasoning, and accordingly the scientific reasoning, abilities of LLMs. We collect a total of 24,946 functionfree solutions nucleus sampling with 5 repeat times and 0.6 temperature. These samples share similar format as ToRA-corpus (Gou et al., 2023b), and do not retrieve/use any functions, *i.e.*,  $F_q = \emptyset$ .

### 4 Model: SCIAGENT

We develop SCIAGENT for tool-augmented scientific reasoning task. It could make plan, retrieve functions, and leverage retrieved functions to facilitate the reasoning. We describe its inference procedure and training approach as below. 270

271

272

273

274

277

279

283

285

287

288

290

291

292

293

296

297

299

300

301

302

303

304

## 4.1 Overview

As shown in Figure 3, SCIAGENT comprises four successive modules.

**Planning.** This module provides a high-level profile for each question:  $G_q = \mathcal{M}_{\text{planning}}(q)$ . Such planning instructs a more targeted retrieval process. **Retrieval.** Given the question and generated planning  $G_q$ , the retriever  $\mathcal{M}_{\text{retrieval}}$  is introduced to retrieve related functions from the domain-specific toolset:  $F_q = \mathcal{M}_{\text{retrieval}}([q, G_q]; F_D) \subseteq F_D$ .

Action. This module aims to generate low-level solutions. Specifically, the agent produces  $S_q = \mathcal{M}_{action}(q; F_q)$ . The solution  $S_q$  is interleaved with natural language rationale  $E_q$  and program snippet  $P_q$ . The program  $P_q$  call retrieved functions with proper arguments if necessary.

**Execution.** This module is simply a Python Executor to run the program  $P_q$  for the final answer:  $a_q = Python-Executor(P_q)$ .

### 4.2 Training

Language models are used in three out of four modules in SCIAGENT: planning, retrieval and action. Rearding retrieval, we directly use the retriever Rfine-tuned in Section 3.2 as  $\mathcal{M}_{retrieval}$ . For planning and action modules, they share the same LLMs:  $\mathcal{M} = \mathcal{M}_{planning} = \mathcal{M}_{action}$ . We fine-tune  $\mathcal{M}$  with different instructions to make it act as planning and action modules, respectively. We construct instructions from  $d = (q, G_q, F_q, S_q, a_q)$  in MATHFUNC.

4

261

265

269

305 
$$D_{\text{planning}} = \{ (I_{\text{plan}}(q), G_q) | d \in D \}$$
  
306 
$$D_{\text{action}} = \{ (I_{\text{action}}(q, F_q), S_q) | d \in D \}$$

Here  $I_{\text{plan}}$  and  $I_{\text{action}}$  are instruction templates for planning and action modules. We show these instructions in Appendix B.2, and mix up them as the training set  $D = (D_{\text{planning}} \bigcup D_{\text{action}})$ . Then we apply imitation learning on D to fine-tune  $\mathcal{M}$ .

$$L_{\mathcal{M}} = \sum_{(X,Y)\in D} -\log \mathcal{P}(Y|X)$$

**Implementation** We detail the training process of (1) the retriever  $\mathcal{M}_{\text{retrieval}}$  and (2) the planner and actor  $\mathcal{M}$  in Appendix B.1 and B.2, respectively.

# **5** Benchmark: SCITOOLBENCH

There currently exists no benchmark assessing the scientific reasoning capabilities of LLMs **when aided by tools**. To address this gap, we develop a benchmark called SCITOOLBENCH. Our benchmark covers five domains: *Mathematics (math)*<sup>5</sup>, *Physics, Chemistry, Finance, Electrical Engineering and Computer Science (EECS)*. Each domain is composed of a set of questions and a domainspecific toolset. The toolset consists of abundant generalized, high-quality and well-documented functions. We expect LLMs to retrieve, understand and, if necessary, use functions in it for reasoning.

Table 1: The statistics of our benchmark. **#Func**: Number of functions. **#Pos./ #Neg.**: The number of positive/negative functions in the toolset. **FPQ** (function per question): The number of derived positive functions from each question.

	# Question	# Func	# Pos. / # Neg.	Avg. FPQ
Math	434	1072	511 / 561	1.47
Physics	156	534	243 / 291	1.63
Chemistry	118	366	155 / 211	1.34
Finance	66	253	97 / 156	1.62
EECS	82	221	97 / 124	1.68
All	856	2446	1103 / 1343	1.51

### 5.1 Dataset Overview.

The statistics of SCITOOLBENCH are presented in Table 1. It comprises a total of 856 questions and 2,446 functions spanning across 5 scientific domains. Notably, SCITOOLBENCH differs from previous tool-based benchmarks, such as Creation



Figure 4: Left: Histogram of FPQ (function per question). Higher values indicate greater composability. **Right**: Histogram of function occurrence. Higher values indicate more generalization and wider application.

Challenge (Qian et al., 2023), in several aspects: (1) Our benchmark encompasses a diverse range of scientific domains. (2) The tools provided are both composable and generalized across different questions. As indicated in Table 1, each question requires an average of 1.51 functions for resolution. And as shown in Figure 4, over 500 functions are designed to be applicable to two or more questions, such as integrate\_function in math domain, coulombs\_law in physical domain, and calculate\_pressure\_van\_der\_waals in chemistry domain. It signifies that the functions in our toolset are not ad-hoc solutions tailored for specific questions. Instead, the effective utilization of the toolset demands significant reasoning abilities of tool-augmented LLMs. Thus we claim this benchmark challenging and practical.



Figure 5: Semi-automatic annotation pipeline for SCI-TOOLBENCH. (S): GPT-4. (a): Human annotator.

### 5.2 Dataset Annotation

We design a pipeline shown in Figure 5 to annotate the benchmark. It employs both GPT-4 and human annotators to combine their merits. We introduce it briefly as below and leave details in Appendix D. **Question Filtering**: We curate questions from TheoremQA (Chen et al., 2023b) and SciBench (Wang et al., 2023b) to collect 856 questions (① in Figure 5, the same below) in our benchmark.

**Toolset Construction**: We construct domainspecific toolsets via two cascade modules: positive and negative function construction. We define positive functions (②) as functions directly deriving from questions. The candidate positive functions 353

354

355

356

357

358

360

361

362

363

364

366

336

337

312

312

315

316

317

321

324

<sup>&</sup>lt;sup>5</sup>Our benchmark contains college-level questions on calculus, differential equations, group theory, *etc*, which are different from the questions in our training corpus MATHFUNC.

Model	Size	Toolset	CDEATION	SCITOOLBENCH					
			CREATION	Math	Physics	Chemistry	Finance	EECS	All
ChatCDT		×	54.6	33.4	19.2	18.6	53.0	25.6	29.6
ChatGP1	-	1	59.8	32.0	31.4	33.9	53.0	48.8	35.4
CDT 4		×	60.0	52.8	42.9	47.5	65.2	35.4	49.5
OF 1-4	-	1	69.8	63.1	63.5	63.6	80.3	80.5	66.2
LLaMA2	7B	1	12.6	4.3	10.9	8.4	13.6	11.0	8.3
CodeLlama	7B	×	17.7	6.5	0.6	5.1	4.9	7.6	5.1
CodeLlama	7B	$\checkmark$	26.1	9.2	8.3	10.2	24.2	25.6	11.9
Llemma	7B	×	26.4	10.4	4.5	8.5	10.6	7.3	8.8
Llemma	7B	$\checkmark$	34.3	16.4	21.2	14.4	36.4	22.0	19.1
Mistral	7B	×	30.1	11.3	4.5	7.6	16.7	6.1	9.5
Mistral	7B	$\checkmark$	27.6	13.1	13.5	14.4	34.8	19.5	15.6
Deepseek-Coder	7B	×	36.8	20.3	8.3	5.9	22.7	12.2	15.5
Deepseek-Coder	7B	$\checkmark$	31.3	21.0	15.4	10.2	30.3	36.6	20.7
Deepseek-Math	7B	×	44.7	26.5	19.2	17.8	27.3	20.7	23.5
Deepseek-Math	7B	$\checkmark$	41.3	24.2	24.4	25.4	43.9	42.7	27.7
ToRA-Coder	7B	×	29.7	26.3	4.5	6.8	9.1	24.4	18.1
ToRA-Coder	7B	$\checkmark$	21.4	21.7	4.5	5.1	13.6	15.9	15.1
MAmmoTH-Coder	7B	$\checkmark$	21.6	14.8	18.5	11.0	25.8	40.0	19.7
SCIAGENT-CODER	7B	1	53.0	30.0	28.3	24.6	39.3	<u>57.3</u>	32.2
SCIAGENT-MISTRAL	7B	1	54.0	31.3	28.8	22.9	<u>51.5</u>	61.0	34.1
SCIAGENT-DEEPMATH	7B	1	60.4	41.2	54.5	44.9	57.5	51.2	46.3
LLaMA2	13B	1	23.3	12.2	11.5	6.8	22.7	14.6	12.4
CodeLlama	13B	×	23.0	9.9	3.2	1.7	9.1	6.1	7.1
CodeLlama	13B	$\checkmark$	38.9	12.7	14.7	7.6	33.3	34.1	16.0
ToRA-Coder	13B	×	30.9	28.6	3.8	4.2	16.7	30.5	20.0
ToRA-Coder	13B	1	28.0	32.0	2.6	11.9	24.2	35.4	23.6
MAmmoTH-Coder	13B	$\checkmark$	34.7	21.4	18.6	11.0	25.8	39.0	21.5
SCIAGENT-CODER	13B	1	<u>54.4</u>	<u>35.0</u>	<u>32.1</u>	<u>28.8</u>	42.4	51.2	<u>35.7</u>

Table 2: Main results on two benchmarks. We highlight our SCIAGENT series in **blue**. The best results (among all open-source LLMs, the same below) are in bold face and the second best are underlined.

(2) are firstly generated from GPT-4. Then human annotators carefully check them and rewrite and/or remove the unqualified ones. We further automatically construct negative functions (3) based on positive functions to reduce the shortcuts in our benchmark. We finally combine both positive and negative functions as the toolset in our benchmark.

## 6 Experiments

# 6.1 Setup

367

368

372

374

375

388

We conduct experiments on SCITOOLBENCH to evaluate the tool-augmented scientific reasoning abilities of LLMs. We also employ CREATION Challenge (Qian et al., 2023) as the second benchmark. It comprises a total of 2,047 samples, with each sample consisting of a question and a groundtruth function. We aggregate all functions to assemble the toolset (thus including 2,047 functions). We report accuracy as the metric in all experiments.

## 6.2 Baselines

We compare SCIAGENT series with eight opensource LLMs: (1) LLaMA-2 (Touvron et al., 2023), (2) CodeLlama (Rozière et al., 2023), (3) Mistral (Jiang et al., 2023), (4) Llemma (Azerbayev et al., 2023), (5) Deepseek-Coder (Guo et al., 2024), (6) Deepseek-Math (Shao et al., 2024), (7) MAmmoTH-Coder (Yue et al., 2023b) and (8) ToRA-Coder (Gou et al., 2023b). We also list the performance of ChatGPT and GPT-4 for reference. We provide all LLMs the same retriever in Section 3.2 to retrieve functions from toolset (if attached). Please see more details in Appendix C.

## 6.3 Main Results

We fine-tune CodeLlama, Mistral and Deepseek-Math for yielding SCIAGENT-CODER, SCIAGENT-MISTRAL and SCIAGENT-DEEPMATH. We show their results in Table 2 and observe: (1) Almost all LLMs present improved performance, *i.e.*, 5.3% absolute and 61.6% relative accuracy increase on average, when supplemented with toolsets. It validates the promise of the tool-augmented setting for scientific reasoning. (2) The models finetuned on math-related datasets from CodeLlama, *i.e.*, ToRA- and MAmmoTH-Coder, perform better than CodeLlama itself by 5.5% abosolute accuracy. It presents the importance of essential math skills among diverse scientific domains. (3)

411

412

390

391

	Planning	Function-augmented solutions	Function-free solutions	Retriever	Accu All	racy (7B) wo. math	Accu All	racy (13B) wo. math
SciAgent-Coder	<ul> <li>✓</li> </ul>	✓(cross-retrieval)	1	$\checkmark$	32.2	34.6	35.7	36.5
Intermediate variants 1-3	X X X	✓(cross-retrieval) ✓(direct-use) ✗	55	5 5 5	$\begin{array}{c c} 30.3 \\ \hline 17.8 \\ 26.3 \end{array}$	<u>33.9</u> 17.3 26.1	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\frac{34.4}{31.0}$ 31.7
CodeLlama wo. retriever	× ×	× ×	× ×	×	11.9 5.1	14.7 3.8	16.0   7.1	19.4 4.3

Table 3: Ablation study on SCITOOLBENCH. We report the accuracy of samples across (1) all domains, (2) four domains excluding the math domain (wo. math).

413 Our agents consistently outperform other opensource LLMs by a large margin. Notably, SCIA-414 GENT-CODER surpasses the most competitive base-415 line, MAmmoTH-Coder, by absolute accuracy of 416 12.5% and 14.2% on the 7B and 13B versions. (4) 417 Our strongest agent, SCIAGENT-DEEPMATH-7B, 418 substantially outperforms ChatGPT with toolset 419 (46.3% v.s. 35.4%) and shows comparable results 420 to GPT-4 without toolset (46.3% v.s. 49.5%). How-421 ever, it still falls significantly behind GPT-4 when 422 both are provided with the same tools. Such gap 423 highlights the challenges of tool-augmented scien-424 tific reasoning (and our benchmark). 425

### 6.4 Ablation Study

426

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

We investigate the effectiveness of components in our training data and agent modules. The specific variants we considered are as follows. (1) We remove the planning module in the agent. (2) We additionally drop the cross-retrieval strategy introduced in Section 3.2. In its place, we construct function-augmented solutions directly from  $\tilde{F}_q$  and  $\tilde{S}_q$ . (3) We further remove all function-augmented solutions from our training data, and only keep the solutions without function callings (functionfree solutions). (4) We do not fine-tune agents but merely use CodeLlama as  $\mathcal{M}_{action}$  for inference. (5) We drop the retriever to disable the LLMs' tool-use abilities. Equivalently, it degrades to the baseline of CodeLlama + PoT (Chen et al., 2023a) prompting.

We illustrate the performance of our agents and their ablated variants in Table 3. We observe that (1) Planning module significantly improves scientific reasoning abilities. As detailed and targeted queries for the retriever, the generated plannings increase the relatedness of retrieved functions. For instance, the function's Recall@3 increases from 48.3% to 53.2% in physics domain, and from 37.3% to 39.8% in chemistry domain. (2) The use of the cross-retrieval strategy is essential. Otherwise, the function-augmented solutions directly from  $\tilde{F}_q$  and  $\tilde{S}_a$  degrade the performance because they are too artificial and ad-hoc to teach LLMs using functions properly. (3) The absence of function-augmented solutions results in a performance drop (row 1 v.s. row 4 in Table 3) of 5.9% and 5.3% in absolute accuracy for 7B and 13B LLMs, respectively. It underscores the critical role of function-augmented solutions to enhance LLMs' tool-use abilities, and the necessity of our MATHFUNC corpus. (4) The removal of function-free solutions (row 4 v.s. row 5) leads to an absolutely 14.4% accuracy decrease. Specifically focusing on non-math samples, there is a notable performance drop of about 12% as well. This clearly demonstrates the fundamental importance of math skills in diverse scientific reasoning tasks, and highlights how our math-related samples enhance LLMs' capabilities in this area. (5) Performance significantly declines when the retriever is removed. It illustrates that the retrieval module is crucial for accessing the appropriate functions from large-scale toolsets.

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

### 6.5 Analysis

Robustness of Toolsets. We acknowledge the construction and maintenance of toolsets is sometime challenging. Therefore, we stress the importance of our agents' robustness. If a sub-par toolset were provided, an robust agent should at the very least perform comparably, if not better, than other competitive LLMs without tool-use. To evaluate the robustness of SCIAGENT-CODER, we simulate two sub-par settings. (1) weak-related: for each question, we restrict the agents from retrieving functions that are directly derived from it. This setting greatly decreases the likelihood of retrieving a proper function from the toolset. (2) unrelated: we completely remove the domain-specific toolset in SCITOOLBENCH. As a substitution, we provide the unrelated toolset constructed in MATHFUNC.

We compare our agents with two competitive LLMs, *i.e.*, ToRA-Coder and MAmmoTH-Coder,

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

Table 4: Accuracy on SCIAGENT with sub-par toolsets. **WR**: weak-related toolsets. **UR**: unrelated toolsets. **NA**: No toolset. The subscripts indicate the difference from the best LLMs (wo. toolsets) each column.

Model	Toolset	Accura All	<b>cy (7B)</b> wo.math	Accura   All	<b>cy</b> (13B) wo. math
SciAgent -Coder	WR UR	$18.8_{\pm 0.7} \\ 14.7_{\pm 3.7}$	$18.0_{+8.3} \\ 10.7_{+1.0}$	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$19.9_{+7.6} \\ 14.7_{+2.4}$
MAmmo-C ToRA-C	NA NA	12.7 18.1	9.0 9.7	16.4 20.0	12.3 11.1

in above two settings. As shown in Table 4, (1) SCIAGENT series with unrelated toolsets present comparable performance with the two LLMs. In other words, our tool-augmented agents are unlikely to degrade the performance even under the extreme scenarios. (2) Our agents with weakrelated toolsets significantly outperform the two LLMs, which further validates the robustness.

493

494

495

496

497

498

499

502

504

508

509

510

511

512

514

515

516

517

518

The Effect of Retriever Quality. We explore the effect of retriever quality on the ending performance. We substitute our fine-tuned retriever in SCIAGENT series by two competitive variants: SimCSE (Gao et al., 2021) and Contriever (Izacard et al., 2021). As shown in Figure 6 (top), our retriever surpasses the other two. It shows that finetuning on the math domain benefits the retrieval of tools in the generalized scientific domains.



Figure 6: **Top**: Performance of SCIAGENT-CODER on SCITOOLBENCH with different retriever variants. **Bottom**: Relationship between the performance and the hit@3 of retrieved functions (artificially controlled).

We further dive deep into the relationship between the hit ratio of tools and the agents' performance. To this end, we manually control the hit@3 ratio by artificially adding/removing the positive functions to/from the retrieved list. Results in Figure 6 (bottom) show a clearly positive correlation between the hit ratio and the task accuracy. It illustrates that the retrieved functions facilitate the reasoning of scientific problems. However, we still observe a limit (40% accuracy) when the hit ratios reaching 100%, showing the challenge of scientific reasoning even when aided by tools. We hope the future work to bridge this performance gap.



Figure 7: The performance of SCIAGENT-CODER (w. toolset) and MAmmoTH-Coder (wo. toolset) on samples which (1) use and (2) not use retrieved functions.

How the Retrieved Functions Benefit. To assess how the retrieved functions aid in the reasoning process of LLMs, we divided the samples into two subsets based on whether our agents use the retrieved functions to solve the problems. We evaluate the performance of these two subsets respectively, comparing with MAmmoTH-Coder series (without tool-use). The results in Figure 7 reveal a two-fold benefit: (1) For samples where functions are explicitly called to solve the questions, our agents demonstrate a substantial 25% improvement in absolute accuracy over LLMs that do not have access to functions. (2) Even for samples that do not directly use functions in their written program, we still observe a slight improvement. It suggests that our agents are capable of learning from retrieved functions as a reference, and then imitate these functions to write their own programs. For instance, example in Figure 12 shows the agents learn how to use scipy.integrate by observing the retrieved function average\_value\_of\_function(...).

# 7 Conclusion

This work proposes *tool-augmented* scientific reasoning, a task aiming to solve challenging scientific problems aided by generalized and scalable tools. To facilitate and evaluate the scientific tooluse abilities of LLMs, we construct a math-related, tool-augmented training corpus MATHFUNC and a benchmark SCITOOLBENCH covering 5 scientific domains. Additionally, we develop open-source agents, SCIAGENT series, as competitive baselines. Extensive experiments reveal that our agents exhibit tool-use abilities exceeding ChatGPT in scientific reasoning tasks.

577 578

580

581

582

584

585

586

590

591

592

594

595

596

597

598

599

601

# Limitations

The primary limitation of our work comes from the 559 way we compile the toolsets in SciToolBench. These tools are constructed directly based on the benchmark's questions, raising concerns about potential information leakage. To address this, we 563 invest significant human effort in our annotation 564 process as detailed in Appendix D.2. We manually 565 review and, if necessary, revise all derived functions to ensure their generalizability and quality. 567 As shown in Figure 6 (bottom), our agents achieve only about 40% accuracy when we provide each question the exact function from which it derives 570 (*i.e.*, 100% hit ratio). It not only highlights the inherent challenge of scientific reasoning tasks, but 572 also suggests that our benchmark suffers minimal impact from the potential information leakage.

We partly attribute this limitation to the absence of a training corpus among scientific (excluding math) domains. The scarcity of annotated solutions for scientific reasoning problems makes it unfeasible to set aside a portion of questions in our benchmark for tool creation. In future work, we plan to collect diverse and high-quality scientific annotations which enable us to develop a more practical and robust tool-augmented benchmark.

# **Ethics Statement**

We ensure that SCITOOLBENCH was constructed in compliance with the terms of use of all source materials and with full respect for the intellectual property and privacy rights of the original authors of the texts. We also provide details on the characteristics and annotation steps of SCITOOLBENCH in Section 5 and Appendix D. We believe our created datasets do not cause any potential risks.

# References

- Zhangir Azerbayev, Hailey Schoelkopf, Keiran Paster, Marco Dos Santos, Stephen McAleer, Albert Q. Jiang, Jia Deng, Stella Biderman, and Sean Welleck. 2023. Llemma: An open language model for mathematics.
- Andres M Bran, Sam Cox, Oliver Schilter, Carlo Baldassari, Andrew D White, and Philippe Schwaller. 2023.
   Chemcrow: Augmenting large-language models with chemistry tools.
- Tianle Cai, Xuezhi Wang, Tengyu Ma, Xinyun Chen, and Denny Zhou. 2023. Large language models as tool makers.

Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2023a. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *Transactions on Machine Learning Research*.

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

655

656

657

658

659

660

661

- Wenhu Chen, Ming Yin, Max Ku, Pan Lu, Yixin Wan, Xueguang Ma, Jianyu Xu, Xinyi Wang, and Tony Xia. 2023b. TheoremQA: A theorem-driven question answering dataset. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7889–7901, Singapore. Association for Computational Linguistics.
- Zhipeng Chen, Kun Zhou, Beichen Zhang, Zheng Gong, Xin Zhao, and Ji-Rong Wen. 2023c. ChatCoT: Tool-augmented chain-of-thought reasoning on chatbased large language models. In *Findings of the Association for Computational Linguistics: EMNLP* 2023, pages 14777–14790, Singapore. Association for Computational Linguistics.
- Ethan Chern, Haoyang Zou, Xuefeng Li, Jiewen Hu, Kehua Feng, Junlong Li, and Pengfei Liu. 2023. Generative ai for math: Abel. https://github.com/ GAIR-NLP/abel.
- Yin Fang, Xiaozhuan Liang, Ningyu Zhang, Kangwei Liu, Rui Huang, Zhuo Chen, Xiaohui Fan, and Huajun Chen. 2023. Mol-instructions: A large-scale biomolecular instruction dataset for large language models.
- Shen Gao, Zhengliang Shi, Minghang Zhu, Bowen Fang, Xin Xin, Pengjie Ren, Zhumin Chen, Jun Ma, and Zhaochun Ren. 2023. Confucius: Iterative tool learning from introspection feedback by easy-to-difficult curriculum.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. SimCSE: Simple contrastive learning of sentence embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6894–6910, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujiu Yang, Nan Duan, and Weizhu Chen. 2023a. Critic: Large language models can self-correct with tool-interactive critiquing.
- Zhibin Gou, Zhihong Shao, Yeyun Gong, yelong shen, Yujiu Yang, Minlie Huang, Nan Duan, and Weizhu Chen. 2023b. Tora: A tool-integrated reasoning agent for mathematical problem solving.
- 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.
- Shibo Hao, Tianyang Liu, Zhen Wang, and Zhiting Hu. 2023. Toolkengpt: Augmenting frozen language models with massive tools via tool embeddings.

Dan Hendrycks, Collin Burns, Steven Basart, Andy

Zou, Mantas Mazeika, Dawn Song, and Jacob Stein-

hardt. 2021a. Measuring massive multitask language

understanding. In 9th International Conference on

Learning Representations, ICLR 2021, Virtual Event,

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul

problem solving with the math dataset. NeurIPS.

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and

Yejin Choi. 2020. The curious case of neural text

degeneration. In 8th International Conference on

Learning Representations, ICLR 2020, Addis Ababa,

Ethiopia, April 26-30, 2020. OpenReview.net.

Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei

Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu,

Chuancheng Lv, Yikai Zhang, Jiayi Lei, Yao Fu,

Maosong Sun, and Junxian He. 2023. C-eval: A

multi-level multi-discipline chinese evaluation suite

Gautier Izacard, Mathilde Caron, Lucas Hosseini, Se-

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Men-

sch, Chris Bamford, Devendra Singh Chaplot, Diego

de las Casas, Florian Bressand, Gianna Lengyel, Guil-

laume Lample, Lucile Saulnier, Lélio Renard Lavaud,

Marie-Anne Lachaux, Pierre Stock, Teven Le Scao,

Thibaut Lavril, Thomas Wang, Timothée Lacroix,

Qiao Jin, Yifan Yang, Qingyu Chen, and Zhiyong Lu.

Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick

Lewis, Ledell Wu, Sergey Edunov, Dangi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-

domain question answering. In Proceedings of the

2020 Conference on Empirical Methods in Natural

Language Processing (EMNLP), pages 6769–6781,

Online. Association for Computational Linguistics.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-

dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,

Luke Zettlemoyer, and Veselin Stoyanov. 2019.

Roberta: A robustly optimized bert pretraining ap-

Yuliang Liu, Xiangru Tang, Zefan Cai, Junjie Lu,

Yichi Zhang, Yanjun Shao, Zexuan Deng, Helan Hu,

Zengxian Yang, Kaikai An, Ruijun Huang, Shuzheng

Si, Sheng Chen, Haozhe Zhao, Zhengliang Li, Liang

Chen, Yiming Zong, Yan Wang, Tianyu Liu, Zhi-

wei Jiang, Baobao Chang, Yujia Qin, Wangchunshu

Zhou, Yilun Zhao, Arman Cohan, and Mark Gerstein.

2023. Ml-bench: Large language models leverage

open-source libraries for machine learning tasks.

2023. Genegpt: Augmenting large language models

with domain tools for improved access to biomedical

mation retrieval with contrastive learning.

and William El Sayed. 2023. Mistral 7b.

bastian Riedel, Piotr Bojanowski, Armand Joulin,

and Edouard Grave. 2021. Unsupervised dense infor-

for foundation models.

information.

proach.

Arora, Steven Basart, Eric Tang, Dawn Song, and

Jacob Steinhardt. 2021b. Measuring mathematical

Austria, May 3-7, 2021. OpenReview.net.

- 673 674 675 676 677
- 678 679
- 686
- 687
- 691

- 697 698
- 701

703

- 710
- 711 712
- 713 714

715 716 717

> 718 719

Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. 2023. Chameleon: Plug-and-play compositional reasoning with large language models.

720

721

722

723

724

726

727

728

729

730

732

734

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, Qingwei Lin, Shifeng Chen, and Dongmei Zhang. 2023. Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct.

OpenAI. 2023. Gpt-4 technical report.

- Siru Ouyang, Zhuosheng Zhang, Bing Yan, Xuan Liu, Jiawei Han, and Lianhui Qin. 2023. Structured chemistry reasoning with large language models.
- Liangming Pan, Alon Albalak, Xinyi Wang, and William Wang. 2023. Logic-LM: Empowering large language models with symbolic solvers for faithful logical reasoning. In Findings of the Association for Computational Linguistics: EMNLP 2023, pages 3806-3824, Singapore. Association for Computational Linguistics.
- Shishir G. Patil, Tianjun Zhang, Xin Wang, and Joseph E. Gonzalez. 2023. Gorilla: Large language model connected with massive apis.
- Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng, Yujia Xie, Yu Hu, Qiuyuan Huang, Lars Liden, Zhou Yu, Weizhu Chen, and Jianfeng Gao. 2023. Check your facts and try again: Improving large language models with external knowledge and automated feedback.
- Cheng Qian, Chi Han, Yi Fung, Yujia Qin, Zhiyuan Liu, and Heng Ji. 2023. CREATOR: Tool creation for disentangling abstract and concrete reasoning of large language models. In Findings of the Association for Computational Linguistics: EMNLP 2023, pages 6922-6939, Singapore. Association for Computational Linguistics.
- Yujia Qin, Shengding Hu, Yankai Lin, Weize Chen, Ning Ding, Ganqu Cui, Zheni Zeng, Yufei Huang, Chaojun Xiao, Chi Han, Yi Ren Fung, Yusheng Su, Huadong Wang, Cheng Qian, Runchu Tian, Kunlun Zhu, Shihao Liang, Xingyu Shen, Bokai Xu, Zhen Zhang, Yining Ye, Bowen Li, Ziwei Tang, Jing Yi, Yuzhang Zhu, Zhenning Dai, Lan Yan, Xin Cong, Yaxi Lu, Weilin Zhao, Yuxiang Huang, Junxi Yan, Xu Han, Xian Sun, Dahai Li, Jason Phang, Cheng Yang, Tongshuang Wu, Heng Ji, Zhiyuan Liu, and Maosong Sun. 2023a. Tool learning with foundation models.
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2023b. Toolllm: Facilitating large language models to master 16000+ real-world apis.

882

883

884

885

886

- 775

- 793
- 795
- 797

- 805
- 815 816
- 810

803

807

817

818

819

821

822

823

824

825

827 828

- Samyam Rajbhandari, Olatunji Ruwase, Jeff Rasley, Shaden Smith, and Yuxiong He. 2021. Zero-infinity: Breaking the gpu memory wall for extreme scale deep learning. In Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, SC '21, New York, NY, USA. Association for Computing Machinery.
- Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2023. Code llama: Open foundation models for code.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models.
- Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. 2023. Hugginggpt: Solving ai tasks with chatgpt and its friends in huggingface. In Advances in Neural Information Processing Systems.
- Yifan Song, Weimin Xiong, Dawei Zhu, Wenhao Wu, Han Qian, Mingbo Song, Hailiang Huang, Cheng Li, Ke Wang, Rong Yao, Ye Tian, and Sujian Li. 2023. Restgpt: Connecting large language models with real-world restful apis.
- Liangtai Sun, Yang Han, Zihan Zhao, Da Ma, Zhennan Shen, Baocai Chen, Lu Chen, and Kai Yu. 2023. Scieval: A multi-level large language model evaluation benchmark for scientific research.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and finetuned chat models.

- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2019. Representation learning with contrastive predictive coding.
- Ke Wang, Houxing Ren, Aojun Zhou, Zimu Lu, Sichun Luo, Weikang Shi, Renrui Zhang, Linqi Song, Mingjie Zhan, and Hongsheng Li. 2023a. Mathcoder: Seamless code integration in llms for enhanced mathematical reasoning.
- Xiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu Zhang, Satyen Subramaniam, Arjun R. Loomba, Shichang Zhang, Yizhou Sun, and Wei Wang. 2023b. Scibench: Evaluating college-level scientific problem-solving abilities of large language models.
- Xingyao Wang, Zihan Wang, Jiateng Liu, Yangyi Chen, Lifan Yuan, Hao Peng, and Heng Ji. 2023c. Mint: Evaluating llms in multi-turn interaction with tools and language feedback.
- Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. 2023. Visual chatgpt: Talking, drawing and editing with visual foundation models.
- Qiantong Xu, Fenglu Hong, Bo Li, Changran Hu, Zhengyu Chen, and Jian Zhang. 2023a. On the tool manipulation capability of open-source large language models.
- Yiheng Xu, Hongjin Su, Chen Xing, Boyu Mi, Qian Liu, Weijia Shi, Binyuan Hui, Fan Zhou, Yitao Liu, Tianbao Xie, Zhoujun Cheng, Siheng Zhao, Lingpeng Kong, Bailin Wang, Caiming Xiong, and Tao Yu. 2023b. Lemur: Harmonizing natural language and code for language agents.
- Rui Yang, Lin Song, Yanwei Li, Sijie Zhao, Yixiao Ge, Xiu Li, and Ying Shan. 2023. Gpt4tools: Teaching large language model to use tools via self-instruction.
- Da Yin, Faeze Brahman, Abhilasha Ravichander, Khyathi Chandu, Kai-Wei Chang, Yejin Choi, and Bill Yuchen Lin. 2023. Lumos: Learning agents with unified data, modular design, and open-source llms.
- Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. 2023. Metamath: Bootstrap your own mathematical questions for large language models. ArXiv preprint, abs/2309.12284.
- Lifan Yuan, Yangyi Chen, Xingyao Wang, Yi R. Fung, Hao Peng, and Heng Ji. 2023a. Craft: Customizing llms by creating and retrieving from specialized toolsets.
- Zheng Yuan, Hongyi Yuan, Chengpeng Li, Guanting Dong, Keming Lu, Chuanqi Tan, Chang Zhou, and Jingren Zhou. 2023b. Scaling relationship on learning mathematical reasoning with large language models.

Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. 2023a. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi.

887

890

891

896

899

900

901 902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

- Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. 2023b. Mammoth: Building math generalist models through hybrid instruction tuning.
- Dan Zhang, Ziniu Hu, Sining Zhoubian, Zhengxiao Du, Kaiyu Yang, Zihan Wang, Yisong Yue, Yuxiao Dong, and Jie Tang. 2024. Sciglm: Training scientific language models with self-reflective instruction annotation and tuning.
  - Wenxuan Zhang, Sharifah Mahani Aljunied, Chang Gao,
     Yew Ken Chia, and Lidong Bing. 2023a. M3exam:
     A multilingual, multimodal, multilevel benchmark
     for examining large language models.
- Yifan Zhang, Jingqin Yang, Yang Yuan, and Andrew Chi-Chih Yao. 2023b. Cumulative reasoning with large language models.
- Yilun Zhao, Hongjun Liu, Yitao Long, Rui Zhang, Chen Zhao, and Arman Cohan. 2023. Knowledgemath: Knowledge-intensive math word problem solving in finance domains.
- Aojun Zhou, Ke Wang, Zimu Lu, Weikang Shi, Sichun Luo, Zipeng Qin, Shaoqing Lu, Anya Jia, Linqi Song, Mingjie Zhan, and Hongsheng Li. 2023. Solving challenging math word problems using gpt-4 code interpreter with code-based self-verification.

# A Detailed Related Work

### A.1 Scientific Reasoning

920

921

922

923

925

926

930

931

932

934

935

936

937

941

942

943

944

947

948

949

951 952

953

955

957

959

961

962

963

964 965

966

967

969

Scientific reasoning can be roughly categorized into two branches: (1) mathematical reasoning and (2) reasoning across other scientific domains.

Mathematical Reasoning. Mathematical (math) reasoning has attracted much more attentions recently. Thanks to abundant training datasets and corpus, there are intensive studies for more powerful math-oriented LLMs by prompt engineering (Qian et al., 2023; Zhang et al., 2023b; Zhou et al., 2023), instruction-tuning (Yuan et al., 2023b; Yue et al., 2023b; Gou et al., 2023b; Yu et al., 2023; Wang et al., 2023a) and even pre-training (Luo et al., 2023; Azerbayev et al., 2023; Chern et al., 2023). Regarding instruction-tuning, we notice that recent studies have automatically constructed high-quality instructions from GPT-4, *i.e.*, finetuning open-source LLMs by Program-of-thought (PoT; Chen et al. 2023a) prompting. It enables open-source LLMs to present remarkable performance, even comparable with GPT-4.

Reasoning across Other Domains. There have been intensive works on scientific LLMs (Bran et al., 2023; Jin et al., 2023; Fang et al., 2023) and benchmarks (Hendrycks et al., 2021a; Huang et al., 2023; Zhang et al., 2023a; Yue et al., 2023a; Sun et al., 2023). However, they primarily target on problems involving less complicated reasoning like knowledge retrieval or simple tool utilization.

Regarding complicated scientific reasoning problems (Chen et al., 2023b; Wang et al., 2023b), questions are scattered among diverse topics and each topic additionally requires domain-specific knowledge. So annotating questions and their solutions domain by domain is much more laborconsuming. Most current benchmarks (Chen et al., 2023b; Wang et al., 2023b; Zhao et al., 2023) merely include hundreds of questions (in all; less for each single domain) from textbooks and provide no training samples. A concurrent work (Zhang et al., 2024) develop a large-scale scientific training corpus, but only focuses three common domains: math, physical and chemistry. Accordingly, the progress of reasoning tasks in these domains is slower than that in math domain: the most competitive approach only achieves 50% and 35% on TheoremQA and SciBench, respectively, not to mention methods built on open-source LLMs. Instead of developing an omniscient and proficient

LLMs on reasoning tasks across various scientific domains, we believe it is more practical to teach LLMs the ability to use domain-specific tools to facilitate their reasoning abilities in some domain when external functions (toolset) are attached.

# A.2 Tool Learning

LLMs, both proprietary ones and open-source ones, demonstrate promising capabilities leveraging external tools to solve problems beyond their limits (Qin et al., 2023a). Combined with specific tools, these *tool-augmented* LLMs achieve great success on various tasks such as machine learning (Wu et al., 2023; Shen et al., 2023; Patil et al., 2023; Yang et al., 2023; Liu et al., 2023), question answering (Peng et al., 2023; Gou et al., 2023a), daily assistance (Xu et al., 2023; Qin et al., 2023b; Song et al., 2023; Gao et al., 2023), *etc*.

Previous work usually pre-defines several tools, e.g., equation solver or calculator, to facilitate math reasoning tasks (Gou et al., 2023a; Lu et al., 2023; Hao et al., 2023; Chen et al., 2023c; Wang et al., 2023c; Xu et al., 2023b; Yin et al., 2023). Cai et al. (2023) generalize the concept of tools to Program functions. Following this concept, CRE-ATOR (Qian et al., 2023) scale up the function number towards thousand level. However, these ad-hoc, argument-free functions are more like solution wrapper rather than well-generalized tools. CRAFT (Yuan et al., 2023a) targetedly design an automatic pipeline to extract generalized functions for tool-use. Though leading to improvement, these functions are still not generalized enough and serve more as reference rather than as tools for direct calling. Ouyang et al. 2023 ask LLM to generate chemistry formulae as knowledge reference to assist the following reasoning and achieve enhanced performance on chemistry questions in SciBench. Similar as our attached toolset, Zhao et al. (2023) maintain a knowledge bank in which saves more than 900 financial definitions/equations/models as the format of functions for retrieval and use. To our best knowledge, our work is the first which (1) finetunes open-source, tool-augmented LLM agents for scientific reasoning tasks and (2) provides a benchmark covering multiple scientific domains to evaluate LLMs' tool-use abilities.

# **B** Training Details

### B.1 Retriever

1016

1019

1021

1024

1025

1026

1027

1030

1031

1032

1035

1037

1038

1039

1041

1042

1043

1045

1048

1049

1050

1051

1053

To fine-tune a retriever, we construct the training samples from MATHFUNC. We concatenate the question and its planning as the query, and view the generated functions as the keys. We finally collect a total of 8603 query-key pairs for training, and split 10% training samples as validation set.

$$ext{query} = [q; G_q] \ ext{key} = f \in ilde{F}_q$$

We follow DPR (Karpukhin et al., 2020) to train a dense retriever R. We use ROBERTA-BASE (Liu et al., 2019) as the backbone. We set the training step as 500, the batch size as 128 and the learning rate as 2e-5. We also set the temperature coefficient of the InfoNCE loss (van den Oord et al., 2019) as 0.07. We run this experiment on a single NVIDIA Quadro RTX8000 GPU. The whole training process lasts for about 20 minutes.

### **B.2** Planning and Action

We fine-tune CodeLlamA (Rozière et al., 2023), Mistral (Jiang et al., 2023) and DeepMath (Shao et al., 2024) on MATHFUNC to develop the planning and action modules in our tool-augmented agents SCIAGENT series, yielding SCIAGENT-CODER, SCIAGENT-MISTRAL and SCIAGENT-DEEPMATH respectively. We set the global batch size as 128. We use the learning rate as 2e-5 for CodeLlama, 2e-6 for Mistral, and 5e-6 for Deep-Math. We use a cosine scheduler with a 3% warmup period for 2 epochs. We train all models with *ZeRO Stage3* (Rajbhandari et al., 2021) on 8 V100 GPUs. The whole training process lasts for about 3 hours for 7B LLMs and 7.5 hours for 13B LLMs.

The planning and action modules share the same model but act differently with different input instructions. We detail the format of planning and action instructions as below:

1054**Planning.** Given a question q, we construct a plan-1055ning sample as  $(I_{\text{plan}}(q), G_q)$ , where  $I_{\text{plan}}(q)$  is the1056input instruction,  $G_q$  is the output, and  $I_{\text{plan}}(.)$  is1057the template for planning module. We provide an1058example of planning instruction as below:

Listing 1: An example of the planning sample. We separate the input instruction and output answer by the dashed line.

Read the following question and provide a high-level, step-by-step plan for	1060 1061
this problem.	1062
Question: Iwo complementary angles are	1063
in a ratio of \$3:2\$. What is the	1064
measure, in degrees, of the smaller	1065
angle?	1066
	1067
Answer:	1068
	1069
To solve this problem, we need to	1070
perform the following steps:	1071
- Set up an equation using the given	1072
ratio and the fact that the angles	1073
are complementary.	1074
- Solve the equation for the variable	1075
representing the ratio.	1076
- Calculate the measure of the smaller	1077
angle using the found variable value.	1078

1030

1082

1083

1084

1085

1088

1121

Action. Given a question q, we construct an action sample as  $(I_{action}(q, F_q), S_q)$ , where  $I_{action}(q, F_q)$ is the input instruction,  $S_q$  is the output answer, and  $I_{action}(.)$  is the template for action module.  $F_q$ are retrieved functions from the toolset in MATH-FUNC. We adopt the retriever in Appendix B.1 to retrieve three functions in the toolset. We provide an example of action instruction as below:

Listing 2: An example of the action sample. We separate the input instruction and output answer by the dashed line. We only show one retrieved function in this sample for visualization convenience.

Read the following mathematical question	1089
and answer it. Please note that you	1091
could (optionally, not required)	1092
call the following functions when	1093
writing the code:	1094
C C	1095
<pre>def angle_ratio_to_degrees(ratio):</pre>	1096
	1097
Computes the angle measures in	1098
degrees given their ratio.	1099
	1100
Parameters:	1101
- ratio (tuple): The ratio of the	1102
angles.	1103
	1104
Returns:	1105
- tuple: The angle measures in	1106
degrees.	1107
- 1	1108
Formula:	1109
- angle_1 = ratio[0] * x	1110
$- angle_2 = ratio[1] * x$	1111
- angle_3 = ratio[2] * x	1112
- where x is the common variable such	1115
that angle_1 + angle_2 + angle_3	1114
	1115
sum ratios - sum(ratio)	1110
v = 180 / sum ratios	1112
$x = 100$ / Sum_ratios	1110
ratio)	110
	112(

```
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1149
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1171
1172
1173
1174
1175
1176
1178
1179
```

1181

[...Two more retrieved functions...]

Question: Two complementary angles are in a ratio of \$3:2\$. What is the measure, in degrees, of the smaller angle?

Answer:

```
# Define the ratio of the angles.
ratio = (3, 2)
 Since the angles are complementary,
    they add up to 90 degrees.
# We can find the value of each part of
    the ratio by dividing 90 by the sum
   of the parts of the ratio.
value_of_each_part = 90 / sum(ratio)
# The measure of the smaller angle is
    the value of each part times the
    smaller part of the ratio.
smaller_angle = value_of_each_part * min
    (ratio)
# Print the result.
print(smaller_angle)
```

#### С **Evaluation Details**

We evaluate (1) the tool-use abilities of LLMs in scientific reasoning task, and (2) their performance when supplemented without any toolsets for comparison. We detail these two settings as below:

With Toolsets. We use the unified PoT-based prompt (Chen et al., 2023a) for all pretrainingbased models and our SCIAGENT series. The unified prompt consists of a short task description and two demonstrations. We show the prompt in Appendix F.4. For each question, we provide three retrieved functions and instruct LLMs to use them if (and only if) necessary. Note that we use the same retriever, *i.e.*, fine-tuned from MATHFUNC, for all LLMs. For MAmmoTH-Coder and ToRA-Coder which are fine-tuned on specific (tool-agnostic) instructions, we try to enable them to use retrieved tools while keeping the formats of their original instructions as much as possible. Specifically, we append a short *tool-augmented* description at the end of their original prompts:

[original prompt]

Please note that you could (optionally, not required) call the following functions when writing the program:

[retrieved functions]

Without Toolsets. Similar as above, we use the unified PoT-based prompt (Chen et al., 2023a) shown in Appendix F.5 for all pretraining-based models and our SCIAGENT series. And we follow the original instructions used for MAmmoTH-Coder and 1182 ToRA-Coder to evaluate their performance. 1183

#### D **Details of SCITOOLBENCH Annotation**

1184

1199

1200

1201

1202

1203

1204

1205

1206

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

We provide a more thorough description about SC-1185 ITOOLBENCH construction in this section. This 1186 semi-automatic annotation pipeline involves both 1187 GPT-4 and humans to balance the quality and cost. 1188 Specifically, we enlist two authors to serve as hu-1189 man annotators. Both of them are graduate students 1190 with proficiency in English. Additionally, they hold 1191 Bachelor of Science and/or Engineering degrees 1192 and have completed undergraduate-level courses 1193 in the five scientific domains corresponding to our 1194 benchmark. We detail the three subsequent sub-1195 modules in our annotation pipeline, *i.e.*, question 1196 filtering, positive function construction and nega-1197 tive function construction, as below. 1198

#### **D.1 Question Filtering**

We curate the questions from TheoremQA (Chen et al., 2023b) and SciBench (Wang et al., 2023b), both of which are available under the MIT License. Among 1495 questions in these original two datasets, we remove three kinds of questions. Image-required: There are 37 questions from TheoremQA which include images and necessitate visual understanding abilities. We remove these samples because our benchmark is text-oriented.

Reasoning-agnostic: There are some multi-choice questions from TheoremQA which merely requires the memorization of knowledge points but involves little reasoning process. For example:

Question: The open mapping theorem can be
proved by
(a) Baire category theorem.
(b) Cauchy integral theorem.
(c) Random graph theorem.
(d) None of the above.

We manually check each samples and remove 68 such kind of samples.

Over-difficult: Too hard questions confuse all models and weaken the discrimination of our benchmark. To balance the difficulty and discrimination, we employ 4 advanced proprietary models<sup>6</sup> to generate related functions and functionaugmented program solutions. We generate 6 so-

<sup>&</sup>lt;sup>6</sup>gpt-4, gpt4-32k, gpt-3.5-turbo, gpt-3.5turbo-16k

lutions for each model (one generated by greedy decoding and the other five by nucleus sampling with 0.6 temperature) and 24 solutions in all. We view questions that are answered incorrectly by all 24 solutions as *over-difficult* questions. We remove all *over-difficult* questions, and retain 73.5% questions in TheoremQA and 47.8% in SciBench.

By removing three kinds of samples mentioned above, there are a total of 865 questions in our SCITOOLBENCH benchmark.

## **D.2** Positive Function Construction

### **Function Generation**

1222

1223

1224

1225

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1242

1243

1244

1246

1247

1248

1249

1250

1251

1252

1253

1254

1255

1256 1257

1259

1261

In practice, we merge this sub-module to the process of over-difficult question identification. We randomly sample one set of functions which yield correct solutions for each question. As a result, we collect a total of 1216 candidates for the next verification sub-module. We additionally save other functions leading to correct solutions and use them as reference in the refinement sub-module.

### Function Verification

We verify the generated functions from both correctness and generalizations. We detail them separately as below.

1. Correctness: Since all candidate functions lead to correct solutions, we speculate that almost all of them are correct. We randomly sample 100 functions (20 per domain) and manually check their correctness. The results shown in Table 5 validate our speculation. Therefore, we assume all candidate functions are correct and retain them.

Table 5: The correctness of 100 randomly sampled functions across five domains.

	Correct	Partially Correct	Wrong	All
Math	18	2	0	20
Physics	19	1	0	20
Chemistry	20	0	0	20
Finance	19	0	1	20
EECS	17	3	0	20
All	93	6	1	100

2. Generalization: We encounter the similar problem as the function construction in MATHFUNC, *i.e.*, some of the auto-generated functions are not generalized enough. If ad-hoc functions were in the provided toolsets of our benchmark, they would cause a significant overestimation of LLMs' tooluse abilities. To mitigate it as much as possible, we manually check all candidate functions to ensure their generalization. Specifically, we design a binary classification task and assign each function a label in {Retained, Refined}. We label a function as refined if it had one of the problems listed below: (1) a pure solution wrapper. (2) merely defining a non-generalized expression (likely only occur in this question). (3) the argument names or document describing the special scenario of corresponding question and not being generalized/abstractive enough. (4) including adhoc constants or code snippets. The annotators firstly co-annotate 100 functions. We calculate Cohen's kappa value of their annotation results as 0.85, illustrating an ideal agreement. Therefore, the annotators separately annotate the remaining functions. It takes about 6 hours per annotator to classify about 650 functions. We show some Refined function cases in Figure 10, and the annotation interface in Figure 8.

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1296

1297

1298

1299

1300

1301

1302

1303

1304

1305

1306

1307

1309

As a result, we collect 1012 Retained and 206 Refined functions. We keep all Retained as the component of positive functions. We also feed the Refined functions to next refinement sub-module to modify them as much as possible.

# **Function Refinement**

This sub-module aims to rewrite 206 Refined functions to make them qualified. To this end, we associate each function with (1) the question from which it is derived, (2) the function-augmented solutions, and (3) the alternative functions from the generation sub-module (if have). Then we provide them to the annotators. The annotators are asked to rewrite the functions to improve their generalization as much as possible. If one function were successfully rewritten, we also require the annotator to write a solution involving the new function to the related question. The solution must yield correct answer to ensure the correctness of the rewritten function. We show some rewritten cases in Figure 10, and the screenshot of the annotation interface in Figure 9.

It takes approximately 12 hours per annotator to check each Refined function and, if applicable, rewrite it. As a consequence, we successfully rewrite 91 Refined functions and drop the remaining ones. We combine these 91 rewritten functions and the 1012 Retained functions to construct 1103 positive functions.

### **D.3** Negative Function Construction

The positive functions constructed above have sat-<br/>isfied the minimum requirements of the toolset in<br/>our benchmark. However, we find that such kind of1310<br/>1311

benchmark contains shortcuts for LLM to retrieve 1313 and use functions. Take a physical question about 1314 frequency-angular conversion as example, the pre-1315 vious modules construct a positive function named 1316 angular\_from\_frequency(...) to solve this 1317 question. Without any other similar functions, the 1318 LLMs could readily select and use the only func-1319 tion by superficial shortcuts. These shortcuts sig-1320 nificantly weaken the function-understanding and 1321 -use abilities evaluation of our benchmark. To miti-1322 gate this problem, we design an additional module 1323 to eliminate the shortcuts by constructing some 1324 (hard) negative functions for each positive func-1325 tion, like frequency\_from\_angular(...) and 1326 frequency\_from\_energy(...) in the above 1327 example. Among three similar functions, LLMs 1328 are forced to understand their usages and choose proper ones to use. In summary, we add negative 1330 functions into the toolset to simulate a more chal-1331 lenging scenario and better evaluate LLMs' tool-1332 use abilities. 1333

Listing 3: Prompt for constructing negative functions

1334 1335

1336

1337

1339

1340

1341

1342

1343 1344

1345

1346 1347 1348

1349

1350

1351 1352 1353

1354

1355 1359

1358

1359

1360

1362

1363

1364

1365 1366

```
Given a function about the {subfield}
   field, could you please write two
   more functions which satisfy:
 The functions should be in the same
   field with the provided function,
   while the knowledge point is not
   compulsorily the same.
- The functions should be similar, but
   not identical with the provided
    function
- The new written functions should be
   wrapped as the below format:
New function 1:
  `python
[new_written_function_1]
New function 2:
  `pvthon
[new_written_function_2]
```

Specifically, we employ GPT-4 for each positive function to generate two similar but not identical functions as the negative functions. The prompt used is shown as below. We do not validate the correctness of negative functions for simplicity, as they are not intended to be used for any question. We filter the duplicated functions and retain the other 1343 functions in all. By merging the 1103 positive functions and 1343 negative functions, we finally collect a total of 2446 functions in our toolset.

[Question 29]: Compute the double integrals over indicated rectangles $\min\{\lim_{x \to a} R \in [-5, 4] \in [0, 3]$
[Function 0]:
python def double_integral(function, x_limits, y_limits): """
Computes the double integral of a given function over a rectangular region.
Parameters: - function (callable): The function to be integrated. It should take two arguments (x, y). - x_limits (tuple): A tuple containing the lower and upper limits of integration for the x-variable. - y_limits (tuple): A tuple containing the lower and upper limits of integration for the y-variable.
Returns: - float: The value of the double integral.
Note: - This function uses the scipy library to compute the double integral. """
from scipy.integrate import nquad
# Define the limits of integration for the x and y variables. limits = [x_limits, y_limits]
<pre># Compute the double integral using the nquad function from scipy. result, _ = nquad(function, limits) return result</pre>
Let's consider the [Function 0]
Retained
Refined

Figure 8: The screenshot of our annotation interface to evaluate functions' generalization.

[Question 0]: Square ABCD. Rectangle AEFG. The degree of ∠AFG-20. Please find ∠AEB in terms of degree. Return the numeric value.
[Function 0]:
def calculate_angle_in_rectangle(angle1, angle2):
""" Calculates the angle in a rectangle given two other angles.
Parameters: - angle1 (float): The first angle in degrees. - angle2 (float): The second angle in degrees.
Returns: - flaat: The calculated angle in degrees. """
return angle1 - angle2
[Solution]: "pythom # Define the angles angle_AFS = 45 # in degrees angle_AFS = 04 = in degrees
# Calculate the angle ∠AEB angle_AE8 * calculate_angle_in_rectangle(angle_AF8, angle_AF6)
print(angle_AEB)
Let's consider to rewrite [Function 0]. You can optionally use one of the below alternative functions (if have) for substitution:
Next

Figure 9: The screenshot of our annotation interface to rewrite functions. We provide no alternative functions in this example for convenience of visualization.

# **Function before rewriting**

# **Function after rewriting**



\_ \_ \_ \_ \_ \_ \_ \_ \_ \_ \_

# Function before rewriting

def calculate\_emptying\_time(height, radius, side\_length, g=9.81): Calculates the time it takes for a cylindrical tank to go from full to empty. Parameters: - height (float): The height of the cylindrical tank. - radius (float): The radius of the cylindrical tank. - side\_length (float): The length of the side of the square hole in the bottom of the tank. - g (float): The acceleration due to gravity. Returns: - float: The time it takes for the tank to empty. from math import pi, sqrt
# Calculate the area of the tank and the hole
tank\_area = pi \* radius\*\*2
hole\_area = side\_length\*\*2 # Use Torricelli's law to calculate the time time = (2 \* height \* tank\_area) / (sqrt(2\*g\*height) \* hole\_area) return time

## Function after rewriting

def calculate\_drain\_time(volume, area, gravity=9.81): Calculates the time it takes for a cylindrical object to drain using Torricelli's Law.

arameters: - volume (float): The volume of the cylindrical object. - area (float): The area of the hole through which the object is draining. - gravity (float): The acceleration due to gravity.

------

Returns: - float: The time it takes for the object to drain.

"""
from math import sqrt
return volume / (area \* sqrt(2\*gravity))

1. Abstract the function description by changing "tank" to "object"

2. Decompose the area calculation and Torricelli's law

# Function before rewriting

Determines if the sumulative distribution function (CDE) of the standard Caussian distribution
Decentiones in the complative distribution function (CDF) of the standard gaussian distribution
is log-concave.
Returns: - int: 1 if the CDF is log-concave, 0 otherwise.
Note: • The second derivative of the natural logarithm of the CDF of the standard Gaussian distribution is always non-positive. • Therefore, the function is log-concave, and we can return 1 without performing any calculations.
return 1

Rewrite the specific function (and its variable) to an argument of the function

# Function after rewriting

	der is_log_concave(t, x):
-	Determines ≱ a given function `f` with respect to variable `x` is log-
	concave.
	Parameters: - f (sympy expression): The function for which the log-concavity is to be checked. - x (sympy symbol): The variable with respect to which log-concavity is to be checked.
	Returns: - bool: True if the function is log-concave, False otherwise. 
	from sympy import diff, log, simplify, solveset, S from sympy.calculus.util import continuous_domain
	# Compute the first derivative of the logarithm of the function first_derivative = diff(log(f), $x)$
	<pre># Compute the second derivative of the logarithm of the function second_derivative = diff(first_derivative, x)</pre>
	<pre># Simplify the second derivative second_derivative_simplified = simplify(second_derivative)</pre>
	# Check if the second derivative is nonpositive (<= 0) # We check if the second_derivative is nonpositive for all x in the domain of f
	# For a general solution, we attempt to prove the second derivative is always nonpositive
	<pre>return solveset(second_derivative_simplified&lt;=0, domain=S.Reals)==continuous_domain(second_derivative_simplified, x, S.Reals)</pre>

Figure 10: Three examples of Refined functions (before rewriting) and their rewritten functions (after rewriting). We also briefly describe the modifications for each example (color in red).

# **E** Examples

### Problem & Planing

Read the following mathematical question and answer it. Please note that you could (optionally, not required) call the following functions when writing the code:

Question:Suppose a stock has the following information. It is listed on the London stock exchange and operates throughout Europe. The yield on a UK 10 year treasury is 2.8%. The stock in question will earn 8.6% as per historical data. The Beta for the stock is 1.4, i.e., it is 140% volatile to the changes in the general stock market. What is the expected rate of return?

(1) Calculate the expected return based on the beta and the yield on a UK 10 year treasury.(2) Add the stock's historical return to the expected return.

### **Retrieved Functions**

```
def expected_return(rf, beta, rm):
    .....
    Computes the expected return using the Capital Asset Pricing Model (CAPM) formula.
   Parameters:
    - rf (float): The risk-free rate.
    - beta (float): The beta of the portfolio.
    - rm (float): The return on the market.
    Returns:
    - float: The expected return.
    .....
    return rf + beta * (rm - rf)
def expected_stock_return(true_prob, u, d):
    Calculates the expected return of the stock in a one-period binomial model.
    Parameters:
    - true_prob (float): The true probability of the stock price going up.
    - u (float): One plus the rate of capital gain on the stock if the price goes up.
    - d (float): One plus the rate of capital loss on the stock if the price goes down.
    Returns:
    - float: The expected return of the stock.
    .....
    return true_prob * u + (1 - true_prob) * d
. . .
```

#### Solution

```
# Given values.
rf = 0.028 # The yield on a UK 10 year treasury
beta = 1.4 # The stock is 140% volatile to the changes in the general stock market
rm = 0.086 # The stock in question will earn 8.6% as per historical data
# Calculate the expected return .
expected_return = expected_return(rf, beta, rm)
# Print the result.
print(expected_return)
```

Figure 11: An example of direct function call.

Problem & Planing

Read the following mathematical question and answer it. Please note that you could (optionally, not required) call the following functions when writing the code:

(1) Integrate the linear density function with respect to x from 0 to 8. (2) Divide the result by the length of the rod.

### **Retrieved Functions**

```
def average_value_of_function(f, a, b):
    Calculates the average value of a function `f` on the interval `[a, b]`.
    Parameters: - f (sympy expression): The function for which the average value is to be found.
    - a (float): The lower limit of the interval.
    - b (float): The upper limit of the interval.
    Returns: float: The average value of the function on the interval.
    from sympy import integrate, symbols
    x = symbols('x')
    integral = integrate(f, (x, a, b))
    average_value = integral / (b - a)
    return average_value
def avg_distance(vectors):
    Finds the average distance from the origin that can be reached by adding and subtracting the given
    vectors.
    Parameters: - vectors (list of numpy arrays): The vectors.
    Returns: float: The average distance.
    from itertools import product
    import numpy as np
    total_distance = 0
    count = 0
    for signs in product([-1, 1], repeat=len(vectors)):
        resultant = sum(v * s for v, s in zip(vectors, signs))
        distance = np.linalg.norm(resultant)
        total_distance += distance
       count += 1
    return total_distance / count
. . .
```

#### Solution

# The linear density is given by 12 / sqrt(x+1) kg/m. # The average density is the total mass of the rod divided by its length. # The total mass of the rod is the integral of the linear density times the length of the rod from 0 to 8. from sympy import symbols, integrate x = symbols('x') length = 8 # Define the linear density function. linear\_density = 12 / (x + 1)\*\*0.5 # Calculate the total mass of the rod. total\_mass = integrate(linear\_density, (x, 0, length)) # Calculate the average density by dividing the total mass by the length. average\_density = total\_mass / length print(average\_density)

Figure 12: An example of implicitly using functions as a reference.

### 1369 F Prompts

1373 1374

1375

1376

1377 1378

1379

1380

1381

1382 1383 1384

1385

1386 1387

1388

1389 1390

1391

1392 1393

1394

1395

1396 1397 1398

1399

1400

1401 1402

1403 1404

1405

1406 1407

1409

1410 1411 1412

1413

1414

1415 1416

1417

1418 1419 1420

1421

1422 1423

1424

1425

1426

1427 1428 1429

1430

1431

1432 1433 1434

1435

### 70 F.1 Tool/Function Collection

```
Listing 4: Prompt used for planning and function collection
Given the question and the reference solution, do the following things:
- Think about what math knowledge points are required to solve this problem step by
    step.
- write some python one or more functions to abstract the solution. Please note that
    the functions should be well-documented as much as possible and not too
    specific (for example, do not write the values in this problem within the
    functions. Pass them as the function arguments). We hope your written functions
   could be re-used in anywhere else.
-Instantiate these functions to solve the problem. The last line of your program
    should be a 'print' command to print the final answer
Here are some examples you may refer to:
Question: There are integers b,c for which both roots of the polynomial x^2-x-1
   are also roots of the polynomial x^5-bx-c. Determine the product bc.
Answer: Let r be a root of x^2-x-1. Then, rearranging, we have n
   $$Multiplying both sides by $r$ and substituting gives\n\\begin{align*}\nr^3 &=
    r^2+r \end{pmatrix} = (r+1)+r \end{pmatrix} = 2r+1.\n\end{align*}Repeating this process twice
    more, we have\n\\begin{align*}\nr^4 &= r(2r+1) \\\\\n&= 2r^2+r \\\\\n&= 2(r+1)+
    r \\\\\n&= 3r+2\n\\end{align*}and\n\\begin{align*}\nr^5 &= r(3r+2) \\\\\n&= 3r
    ^2+2r \\\\\n&= 3(r+1)+2r \\\\n&= 5r+3.\n\\end{align*}Thus, each root of $x^2-x
    -1$ is also a root of x^5-5x-3, which gives bc = 5\cdot 3 = \bcdot 3.
Think: To solve this question, we can follow the steps below: (1) Find the roots of
   the polynomial x^2-x-1. (2) Substitute them into the the polynomial x^5-bx-c and obtain two equations. (3) Solve the equations.
Functions:
  `function 1
def find_roots_of_polynomial(polynomial, variable):
   Finds the roots of a given polynomial using the sympy library.
   Parameters:
   - polynomial (sympy expression): The polynomial whose roots are to be found.
   - variable (sympy symbol): The variable of the polynomial.
   Returns:

    list: The roots of the polynomial.

   ......
   from sympy import solve
   roots = solve(polynomial, variable)
   return roots
```function 2
def substitute_roots_into_polynomial(roots, polynomial, variable):
   Substitutes the given roots into the polynomial and returns the resulting
      expressions.
   Parameters:
   - roots (list): The list of roots to be substituted into the polynomial.
    polynomial (sympy expression): The polynomial into which the roots are to be
       substituted.
   - variable (sympy symbol): The variable of the polynomial.
   Returns:
   - list: The resulting expressions after substituting the roots into the
      polynomial.
   return [polynomial.subs(variable, root) for root in roots]
```

```function 3 1437 def solve\_equations(equations, variables): 1438 1439 Solves a system of equations for the specified variables using the sympy library. 1440 1441 Parameters: 1442 - equations (list of sympy expressions or a single sympy expression): 1443 The equations to be solved. If solving a single equation, this can be a single 1444 expression. 1445 - variables (list of sympy symbols or a single sympy symbol): 1446 The variables for which the solution is to be found. If solving for a single 1447 variable, this can be a single symbol. 1448 1449 1450 Returns: - list of dictionaries: Each dictionary represents a solution, with keys being 1451 the variables and values being their corresponding values. 1452 If there's only one solution, the list will contain a single dictionary. 1453 ..... 1454 1455 from sympy import solve 1456 solution = solve(equations, variables, dict=True) 1457 return solution 1458 1459 1460 Solution: 1461 1462 `python # Import required functions and classes from sympy 1463 from sympy import symbols, Eq 1464 1465 # Define the variable and the polynomials 1466 x, b, c = symbols('x b c') 1467 polynomial1 = x \* \* 2 - x - 11468 polynomial2 = x\*\*5 - b\*x - c1469 1470 # Find the roots of the first polynomial 1471 roots = find\_roots\_of\_polynomial(polynomial1, x) 1472 1473 # Substitute the roots into the second polynomial 1474 resulting\_expressions = substitute\_roots\_into\_polynomial(roots, polynomial2, x) 1475 1476 # Set up the equations based on the resulting expressions 1477 equations = [Eq(expr, 0) for expr in resulting\_expressions] 1478 1479 # Solve the system of equations for b and c 1480 solutions = solve\_equations(equations, (b, c)) 1481 # This linear system has only one solution 1482 solution = solutions[0] 1483 1484 # Calculate the product bc 1485 product\_bc = solution[b] \* solution[c] 1486 print(product\_bc) 1487 1488 1489 \_\_\_ 1490 1491 Question: Medians \$\\overline{DP}\$ and \$\\overline{EQ}\$ of \$\\triangle DEF\$ are 1492 perpendicular. If \$DP= 18\$ and \$EQ = 24\$, then what is \${DE}\$? Answer: Point \$G\$ is the centroid of \$\\triangle DEF\$, so \$DG:GP = EG:GQ = 2:1\$. Therefore, \$DG = \\frac23(DP) = 12\$ and \$EG = \\frac23(EQ) =16\$, so applying the 1494 1495 Pythagorean Theorem to  $\left( EGD \right) = 0 = 0 = 0 = 0$ 1496 boxed $\{20\}$ \$. 1497 Think: Given two perpendicular medians in a triangle, we need to perform the 1498 following steps: (1) Identify the relationship between the segments of medians 1499 and the centroid. (2) Use the ratios provided to determine the lengths of the 1500 individual segments from the centroid to the vertices. (3) Use the Pythagorean 1501 theorem to determine the length of the side connecting the two vertices from 1502 which the medians originate. 1503 Functions: 1504 ``function 1 1505

```
def median_segments_length(median_length, ratio):
   Computes the lengths of the segments of a median split by the centroid.
   Parameters:
   - median_length (float): Total length of the median.
   - ratio (tuple): Ratio in which the centroid splits the median. Default is (2,1)
      for standard triangles.
   Returns:
   - tuple: Lengths of the two segments.
   Formula:
   - segment_1 = ratio[0]/sum(ratio) * median_length
   - segment_2 = ratio[1]/sum(ratio) * median_length
   .....
   segment_1 = ratio[0] / sum(ratio) * median_length
   segment_2 = ratio[1] / sum(ratio) * median_length
   return segment_1, segment_2
. . .
```function 2
def pythagorean_theorem(a, b):
   Computes the hypotenuse of a right triangle given two legs.
   Parameters:
   - a, b (float): Lengths of the two legs.
   Returns:
   - float: Length of the hypotenuse.
   Formula:
   - c = sqrt(a^{2} + b^{2})
   .....
   from sympy import sqrt
   return sqrt(a**2 + b**2)
Solution:
```python
# Given values
DP = 18
EQ = 24
# Point $G$ is the centroid.
ratio = (2,1)
# Determine the lengths of the segments split by the centroid
DG, GP = median_segments_length(DP, ratio)
EG, GQ = median_segments_length(EQ, ratio)
# Use Pythagoras theorem to determine DE
DE = pythagorean_theorem(EG, DG)
print(DE)
```

### F.2 Self-rectification

Listing 5: Prompt used for self-rectification

{
 "role": "system",

1506

1507

1508 1509 1510

1511

1512

1513 1514 1515

1516

1517 1518

1519

1520

1521 1522

1524

1525 1526

1527 1528

1531

1533

1534 1535

1536

1537 1538

1539

1540

1541 1542

1544 1545

1546

1547

1548 1549

1550

1553 1554

1555

1556

1557 1558

1559

1569

1562

1566

1567 1568

1569

1570 1571

<sup>&</sup>quot;content": "Suppose you are writing Python programs to solve mathematical questions. Your original program, however, can not run normally for some reasons. Given the error message returned from the interpreter, please debug your program carefully and generate the new program. Your generated correction MUST FOLLOW the following format:\n``function 1\n[the first created function]\n```\n...\n```function N\n[the N-th created function]\n```\n

```
```python\n[solution]\n```.\nD0 NOT only quote your output as ```python\n[
  1573
             solution]\n`
  1574
},
  1575
  1576
      "role": "user".
  1577
      "content": "Original Program:\n```function 1\ndef solve_quadratic(a, b, c):\n
  1578
            from sympy import symbols, Eq, solven x = symbols('x') n quadratic = Eq(a *
  1579
            x**2 + b * x + c, 0)\n solutions = solve(quadratic, x)\n return solutions\n
  1580
             ```\n\n```python\ndef find_d_values():\n a, b, c, d = symbols('a b c d')\n eq1
                                                                                                                                                                          1581
              = Eq(6 * a * d, a + d)\n eq2 = Eq((6 * d - 1) * (1 - a), a + d)\n eq3 = Eq(a,
4 * d - 1)\n solution = solve([eq1, eq2, eq3], (a, b, c, d), dict=True)\n
                                                                                                                                                                          1582
                                                                                                                                                                          1583
            d_values = [sol[d] for sol in solution]\n return d_values\nd_values =
                                                                                                                                                                          1584
             find_d_values()\nprint(d_values)\n```\n\nError Message:\nFile <string>:15\n\
                                                                                                                                                                          1585
            nFile <string>:9, in find_d_values()\n\nNameError: name 'Eq' is not defined"
                                                                                                                                                                          1586
},
{
                                                                                                                                                                          1587
                                                                                                                                                                          1588
      "role": "assistant",
                                                                                                                                                                          1589
      "content": "Analysis: The error is occurring because the Eq function from the
                                                                                                                                                                          1590
             sympy module is not recognized in the find_d_values function. This is
                                                                                                                                                                          1591
            happening because the ```from sympy import symbols, Eq, solve``` line is only
                                                                                                                                                                          1592
            inside the ```solve_quadratic``` function, so Eq and other imported functions
                                                                                                                                                                          1593
            and classes are only available within the scope of that function. To fix this,
                                                                                                                                                                          1594
            you should claim import statement within the ``find_d_values``` function as well.\n\ Program:\ function 1\ndef solve_quadratic(a, b, c):\n from sympy import symbols, Eq, solve\n x =  symbols('x')\n quadratic = Eq(a * x**2)
                                                                                                                                                                          1595
                                                                                                                                                                          1596
                                                                                                                                                                          1597
             + b * x + c, 0)\n solutions = solve(quadratic, x)\n return solutions\n```\n\n
                                                                                                                                                                          1598
            ``python\ndef find_d_values():\n from sympy import symbols, Eq, solve\n a, b,
c, d = symbols('a b c d')\n eq1 = Eq(6 * a * d, a + d)\n eq2 = Eq((6 * d -
                                                                                                                                                                          1599
            1) * (1 - a), a + d\n eq3 = Eq(a, 4 * d - 1)\n solution = solve([eq1, eq2, eq3], (a, b, c, d), dict=True)\n d_values = [sol[d] for sol in solution]\n
                                                                                                                                                                          1601
                                                                                                                                                                          1602
             return d_values\nd_values = find_d_values()\nprint(d_values)\n```"
                                                                                                                                                                          1603
},
                                                                                                                                                                          1604
                                                                                                                                                                          1605
      "role": "user",
                                                                                                                                                                          1606
      "content": "finished"
                                                                                                                                                                          1607
                                                                                                                                                                          1608
},
      "role": "user"
                                                                                                                                                                          1610
      "content": "Original Program:\n```function 1\ndef matrix_square(matrix):\n from
                                                                                                                                                                          1611
             sympy import Matrix\n if not isinstance(matrix, Matrix):\n matrix = Matrix(
                                                                                                                                                                          1612
            matrix)\n return matrix * matrix\n```\n\n```function 2\ndef solve_equations(
                                                                                                                                                                          1613
            equations, variables): \n from sympy import solve \n solution = solve(
                                                                                                                                                                          1614
            equations, variables, dict=True) \n return solution \n^{n} n^{n} n^{
                                                                                                                                                                          1615
                                                                                                                                                                          1616
            nfrom sympy import symbols, Eq \na, d = symbols('a d') \nmatrix = [[a, 3],
                                                                                                                                                                          1617
            [-8, d]] \nsquared_matrix = matrix_square(matrix) \nidentity_matrix = [[1, 0],
[0, 1]] \nequations = [Eq(squared_matrix[i][j], identity_matrix[i][j]) for i
                                                                                                                                                                          1618
                                                                                                                                                                          1619
               in range(2) for j in range(2)] \nsolution = solve_equations(equations, [a, d
                                                                                                                                                                          1620
            ]) \nans = count_solutions(solution)\nprint(ans)\n```\n\nError Message:\nFile
                                                                                                                                                                          1621
               <string>:17\n\nFile <string>:17, in <listcomp>(.0)\n\nTypeError: 'Add'
                                                                                                                                                                          1622
             object is not subscriptable"
                                                                                                                                                                          1623
},
{
                                                                                                                                                                          1624
                                                                                                                                                                          1625
      "role": "assistant",
                                                                                                                                                                          1626
      "content": "Analysis: The error message is occurring because when you are
                                                                                                                                                                          1627
            creating the list of equations with the list comprehension, squared_matrix[i
                                                                                                                                                                          1628
            ][j] is not working as expected. The squared_matrix is a SymPy Matrix object,
              and you should use the .row and .col methods.
\n
New Program:
_____function 1
                                                                                                                                                                          1630
             ndef matrix_square(matrix):\n from sympy import Matrix\n if not isinstance(
                                                                                                                                                                          1631
            matrix, Matrix):\n matrix = Matrix(matrix)\n return matrix * matrix\n``\n\n
             ```function 2\ndef solve_equations(equations, variables): \n from sympy
  1633
            import solve \n solution = solve(equations, variables, dict=True) \n return
solution \n```\n\n``function 3\ndef count_solutions(solutions):\n return len(
solutions)\n```\n\n```python\nfrom sympy import symbols, Eq \na, d = symbols('
  1634
  1635
  1636
            a d') \nmatrix = [[a, 3], [-8, d]] \nsquared_matrix = matrix_square(matrix) \
nidentity_matrix = [[1, 0], [0, 1]] \nequations = [Eq(squared_matrix.row(i)[j], identity_matrix[i][j]) for i in range(2) for j in range(2)] \nsolution =
  1637
  1638
  1639
             solve_equations(equations, [a, d]) \nans = count_solutions(solution)\nprint(
             ans)\n```"
  1641
},
  1642
```

```
{
     "role": "user",
     "content": "finished"
},
     "role": "user"
     "content": "Original Program:\n```function 1\ndef sum_to_product_cos(alpha, beta)
           :\n return (cmath.cos(alpha - beta) + cmath.cos(alpha + beta)) / 2\n
   function 2\ndef simplify_expression(expression): \n from sympy import
           simplify \n return simplify(expression) n^{n} python/nfrom sympy import symbols, cos\n\nx = symbols('x')\nexpr = cos(x) + cos(3*x) + cos(7*x) + cos(7*x)
           (9*x)\n\n# Apply sum-to-product formula\nexpr1 = 2 * sum_to_product_cos(5*x,
           4*x)\nexpr2 = 2 * sum_to_product_cos(5*x, 2*x)\n\n# Simplify the expression\
           nsimplified_expr = simplify_expression(expr1 + expr2)\n\n# Extract the values
           of a, b, c, and d and sum them up\na, b, c, d = 4, 1, 3, 5\nans = a + b + c
+ d\nprint(ans)\n```\n\nError Message:\nFile <string>:12\n\nFile <string>:2,
           in sum_to_product_cos(alpha, beta)\n\nNameError: name 'cmath' is not defined"
},
{
     "role": "assistant",
     "content": "Analysis: The error message is occurring because the cmath module is
           being used in the sum_to_product_cos function but it has not been imported.
           You can fix this by importing the cmath module in this function. New Program
           :\n```function 1\ndef sum_to_product_cos(alpha, beta):\n import cmath\n
           return (cmath.cos(alpha - beta) + cmath.cos(alpha + beta)) / 2\n```\n\n```
           function 2\ndef simplify_expression(expression): \n from sympy import
           simplify \n return simplify(expression)\n``` \n\n```python\nfrom sympy import
           symbols, cos\n x = symbols('x')\nexpr = cos(x) + cos(3*x) + cos(7*x) + cos(7*x)
           (9*x)\n\n# Apply sum-to-product formula\nexpr1 = 2 * sum_to_product_cos(5*x,
           4*x\\nexpr2 = 2 * sum_to_product_cos(5*x, 2*x)\n\n# Simplify the expression\
           nsimplified_expr = simplify_expression(expr1 + expr2)\n\
            of a, b, c, and d and sum them up\na, b, c, d = 4, 1, 3, 5 nans = a + b + c
           + d\nprint(ans)\n```"
},
     "role": "user",
     "content": "Error message:\nTraceback (most recent call last):\n File, line 9, in
             <module>\n exec(\n File '<string>', line 13, in <module>\n File '<string>',
           line 3, in sum_to_product_cos\n File, line 356, in __complex__\n return
           complex(float(re), float(im))\n File, line 351, in __float__\n raise
TypeError('Cannot convert expression to float')\nTypeError: Cannot convert
           expression to float\n"
},
     "role": "assistant",
     "content": "Analysis: The error occurs because you are trying to use cmath.cos
           function (which expects float or complex numbers as arguments) with SymPy
           symbolic expressions (5x and 4x, 5x and 2x are symbolic expressions). To fix
           this, you should use SymPy's cos function instead of cmath.cos function to
           work with symbolic expressions.\n\nNew Program:\ndef sum_to_product_cos(alpha,
    beta):\n from sympy import cos\n return (cos(alpha - beta) + cos(alpha +
           beta)) / 2\ndef simplify_expression(expression): \n from sympy import
           simplify \n return simplify(expression) \nfrom sympy import symbols, cos\n\nx
            = symbols('x')\nexpr = cos(x) + cos(3*x) + cos(7*x) + cos(9*x) \ln \# Apply
           sum-to-product formula\nexpr1 = 2 * sum_to_product_cos(5*x, 4*x)\nexpr2 = 2 *
             sum_to_product_cos(5*x, 2*x)\n\n# Simplify the expression\nsimplified_expr =
             simplify_expression(expr1 + expr2)\n\ Extract the values of a, b, c, and d
             and sum them up\na, b, c, d = 4, 1, 3, 5 = a + b + c + d print(ans)
},
     "role": "user",
     "content": "finished"
},
     "role": "user".
     "content": "Original Program:\n```function 1\ndef tan_double_angle(tan_theta):\n
           return 2 * tan_theta / (1 - tan_theta**2)\n```\n\n```function 2\ndef
           tan_triple_angle(tan_theta):\n return (3 * tan_theta - tan_theta**3) / (1 - 3
 * tan_theta**2)\n```\n\n```function 3\ndef solve_polynomial(polynomial, x):\n
             from sympy import solve
\n return solve(polynomial, x)
\n```\n
\n```python
\nfrom
```

1644

1645

1646 1647

1648

1649

1650

1651

1652 1653

1654

1655

1656

1658

1659

1660 1661

1662

1663 1664

1665

1668

1669

1670

1672 1673

1674

1675 1676

1678 1679

1681

1682 1683 1684

1685

1686

1687

1690

1692 1693

1694

1695

1696

1697

1699

1701

1702

1703

1704 1705

1706

1707

1708

1709

1710 1711

```
26
```

	sympy import symbols\ntan_theta = symbols('tan_theta')\ntan_2_theta =	1713
	tan double angle(tan theta)\ntan 3 theta = tan triple angle(tan theta)\	1714
	neguation = tan theta + tan 2 theta + tan 3 theta\nsimplified equation =	1715
	simplify expression(equation) the tas $solutions = solve polynomial($	1716
	simplified equation, tan theta) $\ln \pi$ Find the solution in the given range (0	1717
	<pre>&lt; tan theta &lt; 1\\ntan theta value = [solution for solution in</pre>	1718
	tan theta solutions if $0 < solution < 11[0] norint(tan theta value) (n:::)n)$	1710
	nFror Message \nFile <tring>13\n\nMameFror name 'simplify expression' is</tring>	1720
	not defined"	1720
ι	hot der med	1721
, נ ג		1722
ι	"rolo", "assistant"	1723
	"content", "Analysis, The error message is accurring because there is no function	1724
	content. Analysis, the error message is occur ing because there is no information	1725
	define such a function that uses Sympler complete function to simplify	1720
	evented a function that uses sympt is simplify function to simplify	1720
	expressions. (n(new program: (n - thection 1(neet simplify_expression(	1720
	expression): (n from sympy import simpting (n return simplify(expression))n	1729
	(i) in the the set of a set of the set of th	1730
	/ (1 - tan_theta**2) \n \n\n Iunction 3(hdef tan_tripte_angle(tan_theta): \n	1731
	return $(3 \times \tan \tan \tan 2 - \tan \tan 2 - \tan 2 $	1732
	tunction 4 nder solve_polynomial(polynomial, x): (n from sympy import solve(n	1733
	return solve(polynomial, x) \n \n\n python\ntrom sympy import symbols\	1734
	ntan_theta = symbols('tan_theta')(htan_2_theta = tan_double_angle(tan_theta))	1735
	ntan_3_theta = tan_triple_angle(tan_theta)\nequation = tan_theta +	1736
	tan_2_theta + tan_3_theta\nsimplified_equation = simplify_expression(equation	1737
	)\ntan_theta_solutions = solve_polynomial(simplified_equation, tan_theta)\n\n	1738
	# Find the solution in the given range $(0 < tan_theta < 1)$ (ntan_theta_value =	1739
	[solution for solution in tan_theta_solutions if 0 < solution < 1][0]\nprint	1740
	(tan_theta_value)\n``"	1741
},		1742
ł		1743
	"role": "user",	1744
	"content": "finished"	1745
}		1746
<b>F.3</b>	Function-augmented Solutions	1748
	Listing 6: Prompt used for the generation of function-augmented solutions (cross-retrieval strategy)	
You will encounter a mathematical problem and are required to write a piece of		

You will encounter a mathematical problem and are required to write a piece of	
Python code to solve this problem.	1751
	1752
Now we have a suite of wrapped functions. Take note:	
- The newly provided wrapped functions have NOT been verified. They may be	1754
irrelevant or potentially flawed.	1755
- It's essential that the solution doesn't overly depend on wrapped functions.	
You're welcome to utilize one or more functions from the new set in your solution	
but only after you've determined:	1758
(1) Their accuracy.	1759
(2) Their inclusion significantly streamlines the problem-solving approach.	1760
	1761
Additionally take note that	
(1) The last line of your written code shall be a 'print' command to print the	1763
final answer.	1764
(2) The wrapped functions should not be duplicated within your code. Instead,	1765
call them directly if needed.	1766
(3) Should you need to create custom functions, do so without adding	1767
documentation comments for the sake of brevity.	1768
(4) Write simple but clear annotations interleaving your code solution.	1769
	1770
	1771
Retrieved functions:	
LIST OT CALLED TUNCTION NAMES TROM THE NEW SET]	1773

Your Written Python Code.]

```
1780
            For example:
1781
1782
            Question: What is the 100th digit to the right of the decimal point in the decimal
1783
                representation of \frac{13}{90}?
1784
1785
            New provided functions:
             ``New Function 0
1786
            def decimal_representation(numerator, denominator, max_digits=1000):
1787
1788
1789
               Computes the decimal representation of a fraction.
1790
1791
               Parameters:
1792
               - numerator (int): The numerator of the fraction.
1793
               - denominator (int): The denominator of the fraction.
1794
               - max_digits (int): The maximum number of decimal digits to compute.
1795
1796
               Returns:
               - str: The decimal representation of the fraction as a string.
1797
               .....
1798
1799
               result = ""
1800
1801
               remainder = numerator % denominator
               for _ in range(max_digits):
1802
                  remainder *= 10
1803
                  result += str(remainder // denominator)
1804
                  remainder %= denominator
1805
                  if remainder == 0:
1806
1807
                     break
               return result
            . . .
1809
1810
            ```New Function 1
1811
1812
            def decimal_to_scientific(decimal_number):
               from sympy import log, floor
1813
               exponent = -floor(log(decimal_number, 10))
1814
1815
               coefficient = decimal_number * 10**(-exponent)
               return coefficient, exponent
1816
1817
1818
            ```New Function 2
1819
1820
            def repeating_decimal_representation(numerator, denominator):
1821
1822
               Computes the repeating decimal representation of a fraction.
1823
1824
               Parameters:
                - numerator (int): The numerator of the fraction.
               - denominator (int): The denominator of the fraction.
1827
               Returns:
1829
               - str: The repeating decimal representation of the fraction as a string.
               .....
1830
1831
1832
               # Initialize the result string and a dictionary to store remainders.
1833
               result = ""
               remainders = {}
1834
1835
               # Perform long division to find the decimal representation.
1836
               while numerator != 0:
1838
                  # If the remainder has been seen before, we found the repeating block.
                  if numerator in remainders:
1840
                      start = remainders[numerator]
                      return result[:start] + "(" + result[start:] + ")"
1841
1842
                  # Otherwise, store the remainder and continue the division.
1843
                  remainders[numerator] = len(result)
1844
                  numerator *= 10
1845
                  result += str(numerator // denominator)
1846
                  numerator %= denominator
1847
            return result
1849
```

```
28
```

```
```New Function 3
                                                                                                 1851
def nth_digit_of_decimal_representation(numerator, denominator, n):
                                                                                                 1852
                                                                                                 1853
   Computes the nth digit after the decimal point of the decimal representation of a
                                                                                                 1854
                                                                                                 1855
        fraction.
   Parameters:
                                                                                                 1857
   - numerator (int): The numerator of the fraction.
                                                                                                 1858
   - denominator (int): The denominator of the fraction.
                                                                                                 1859
   - n (int): The position of the digit after the decimal point.
                                                                                                 1860
                                                                                                 1861
                                                                                                 1862
   Returns:
   - int: The nth digit after the decimal point of the decimal representation of the
                                                                                                 1863
                                                                                                 1864
        fraction.
   .....
                                                                                                 1865
                                                                                                 1866
   # Get the repeating decimal representation of the fraction.
   decimal_representation = repeating_decimal_representation(numerator, denominator)
                                                                                                 1868
                                                                                                 1869
   # Remove the parentheses from the repeating block.
                                                                                                 1870
   decimal_representation = decimal_representation.replace("(", "").replace(")", "")
                                                                                                 1871
                                                                                                 1872
   # Calculate the nth digit using the repeating block.
                                                                                                 1873
   return int(decimal_representation[(n - 1) % len(decimal_representation)])
                                                                                                 1875
                                                                                                 1876
Retrieved functions:
                                                                                                 1877
[decimal_representation, nth_digit_of_decimal_representation]
                                                                                                 1879
```python
  1880
# Use the nth_digit_of_decimal_representation function to find the 100th digit
   1881
numerator = 13
   1882
denominator = 90
   1883
n = 100
   1885
# Call the function and print the result
   1886
result = nth_digit_of_decimal_representation(numerator, denominator, n)
print(result)
   1888
  1889
   1890
   1891
   1892
Question: The square root of x is greater than 3 and less than 4. How many integer
   1893
     values of $x$ satisfy this condition?
   1894
   1895
New provided functions:
 ``New Function 0
   1897
def solve_square_root_equation(a, b, c):
   .....
   1899
   Solves a square root equation of the form sqrt(ax - b) = c.
   1900
  1901
   Parameters:
  1902
   - a (float): Coefficient of x inside the square root.
  1903
   - b (float): Constant term inside the square root.
   1904
   - c (float): Constant term on the right side of the equation.
  1905
   1906
   Returns:
  1907
   - float: The value of x that satisfies the equation.
   1908
  1910
   Formula:
   -x = (c^2 + b) / a
   1911
   .....
   1912
  return (c**2 + b) / a
   1913
  1914
  1915
```New Function 1
                                                                                                1916
def find_integer_square_less_than_double():
                                                                                                 1917
   .....
                                                                                                1918
   Finds the only integer whose square is less than its double.
                                                                                                 1919
```

```
1921
               Returns:
                - int: The integer that satisfies the condition.
1922
1923
1924
               Method:
                 Iterate through integers starting from 1, and check if the square of the
1925
                    integer is less than its double.
                - If the condition is satisfied, return the integer.
1927
1928
                - If the condition is not satisfied for any integer up to a certain limit, return
1929
                    None.
                .....
1930
               limit = 100
1931
1932
                for x in range(1, limit):
1933
                   if x**2 < 2*x:
1934
                      return x
1935
               return None
1936
1938
            ```New Function 2
1939
            def solve_equation():
               .....
1940
1941
               Solves the equation (x-2)^{(25-x^2)} = 1 for integer solutions.
1942
1943
               Returns:
                - list: A list of integer solutions for x.
1944
                .....
1945
               solutions = []
1947
                # Case 1: Exponent is 0 (25 - x^2 = 0)
1949
               x1 = 5
1950
                x2 = -5
               solutions.extend([x1, x2])
1951
1952
1953
                # Case 2: Base is 1 (x - 2 = 1)
1954
               x3 = 3
               solutions.append(x3)
1956
                # Case 3: Base is -1 and exponent is even (x - 2 = -1 \text{ and } 25 - x^2 = 2n \text{ for some}
1958
                   integer n)
1959
                x4 = 1
1960
               solutions.append(x4)
1961
1962
               return solutions
            . . .
1963
1964
            ```New Function 3
1965
            def count_integers_in_range(lower_bound, upper_bound, exclude_zero=True):
1967
               Counts the number of integers within a given range.
1969
1970
               Parameters:
                - lower_bound (int): The lower bound of the range.
1971
1972
                - upper_bound (int): The upper bound of the range.
                - exclude_zero (bool): Whether to exclude 0 from the count. Default is True.
1973
1974
               Returns:
1976
                - int: The number of integers within the range.
                .....
1978
                count = upper_bound - lower_bound + 1
                if exclude_zero and lower_bound <= 0 and upper_bound >= 0:
                   count -= 1
1980
1981
               return count
            . . .
1982
1983
            Retrieved functions:
1984
1985
            ٢٦
            ```python
1987
            # The lower and upper bounds of x for which sqrt(x) > 3 and sqrt(x) < 4
1989
            lower_bound = 9
```

upper\_bound = 16
# Counting the number of integers between 9 (exclusive) and 16 (exclusive)
num\_integers = len([x for x in range(lower\_bound + 1, upper\_bound)])
# Printing the result
print(num\_integers)

1990 1991

1992

1993 1994

1995

1996 1998

1999

# F.4 Evaluation with Toolsets

Listing 7: Prompt used for evaluation (setting with toolsets) Read the following questions and answer them. For each question, you are required to write a Python program to solve it. Please note that we provide you several functions for each question. You could ( 2003 optionally, not required) call the functions to help you to solve the question 2004 if necessary. Note that the last line of your program should be a 'print' command to print the 2006 final answer 2007 2008 \_\_\_\_\_ Question: 2010 What is the 100th digit to the right of the decimal point in the decimal 2011 representation of \$\\frac{13}{90}\$? Functions: 2014 def repeating\_decimal\_representation(numerator, denominator): ..... 2016 Computes the repeating decimal representation of a fraction. 2017 2018 Parameters: 2019 - numerator (int): The numerator of the fraction. 2020 - denominator (int): The denominator of the fraction. 2021 2022 Returns: 2023 - str: The repeating decimal representation of the fraction as a string. ..... 2025 2026 # Initialize the result string and a dictionary to store remainders. 2027 result = "" remainders = {} 2030 # Perform long division to find the decimal representation. while numerator != 0: 2032 # If the remainder has been seen before, we found the repeating block. 2033 if numerator in remainders: 2034 start = remainders[numerator] return result[:start] + "(" + result[start:] + ")" # Otherwise, store the remainder and continue the division. 2037 remainders[numerator] = len(result) numerator \*= 10 result += str(numerator // denominator) numerator %= denominator 2041 2042 return result 2043 def nth\_digit\_of\_decimal\_representation(numerator, denominator, n): 2046 ..... 2047 Computes the nth digit after the decimal point of the decimal representation of a fraction. Parameters: - numerator (int): The numerator of the fraction. 2052 - denominator (int): The denominator of the fraction. - n (int): The position of the digit after the decimal point. 2054 Returns: 2056

```
- int: The nth digit after the decimal point of the decimal representation of the
       fraction.
   .....
   # Get the repeating decimal representation of the fraction.
   decimal_representation = repeating_decimal_representation(numerator, denominator)
   # Remove the parentheses from the repeating block.
   decimal_representation = decimal_representation.replace("(", "").replace(")", "")
   # Calculate the nth digit using the repeating block.
   return int(decimal_representation[(n - 1) % len(decimal_representation)])
def decimal_representation(numerator, denominator, max_digits=1000):
   .....
   Computes the decimal representation of a fraction.
   Parameters:
   - numerator (int): The numerator of the fraction.
   - denominator (int): The denominator of the fraction.
   - max_digits (int): The maximum number of decimal digits to compute.
   Returns:
   - str: The decimal representation of the fraction as a string.
   .....
   result = ""
   remainder = numerator % denominator
   for _ in range(max_digits):
      remainder *= 10
      result += str(remainder // denominator)
      remainder %= denominator
      if remainder == 0:
        break
   return result
Solution:
# find the 100th digit.
numerator = 13
denominator = 90
n = 100
# Call the function and print the result.
result = nth_digit_of_decimal_representation(numerator, denominator, n)
print(result)
 _____
Ouestion:
The square root of $x$ is greater than 3 and less than 4. How many integer values of
    $x$ satisfy this condition?
Functions:
def count_integers_in_range(lower_bound, upper_bound, exclude_zero=True):
   .....
   Counts the number of integers within a given range.
   Parameters:
   - lower bound (int): The lower bound of the range.
   - upper_bound (int): The upper bound of the range.
   - exclude_zero (bool): Whether to exclude 0 from the count. Default is True.
   Returns:
   - int: The number of integers within the range.
   .....
   count = upper_bound - lower_bound + 1
   if exclude_zero and lower_bound <= 0 and upper_bound >= 0:
      count -= 1
```

2058

2059

2061

2062

2063

2064

2067

2070 2071

2072 2073

2074 2075

2076 2077

2078

2079 2080

2081

2084

2085 2086

2087

2088

2090

2091

2094

2095 2096

2097

2098

2099 2100 2101

2102

2103

2104 2105

2106 2107

2108 2109

2110 2111

2112

2113 2114

2115 2116

2117 2118

2119

2120 2121

2122

2123 2124

```
return count
```

```
2128
  2129
def find_integer_square_less_than_double():
  2130
   .....
  2131
   Finds the only integer whose square is less than its double.
  2132
  2133
   Returns:
  2134
   - int: The integer that satisfies the condition.
  2135
  2136
   Method:
  2137
   - Iterate through integers starting from 1, and check if the square of the
  2138
       integer is less than its double.
  2139
   - If the condition is satisfied, return the integer.
  2140
   - If the condition is not satisfied for any integer up to a certain limit, return
  2141
  2142
        None.
   .....
  2143
   limit = 100
  2144
   for x in range(1, limit):
  2145
      if x**2 < 2*x:
  2146
         return x
  2147
   return None
  2148
  2149
  2150
def solve_square_root_equation(a, b, c):
  2151
  2152
   Solves a square root equation of the form sqrt(ax - b) = c.
  2153
  2154
   Parameters:
  2155
   - a (float): Coefficient of x inside the square root.
  2156
   - b (float): Constant term inside the square root.
  2157
   - c (float): Constant term on the right side of the equation.
  2158
  2159
   Returns:
  2160
   - float: The value of x that satisfies the equation.
  2161
  2162
   Formula:
  2163
   -x = (c^2 + b) / a
  2164
   .....
  2165
   return (c**2 + b) / a
  2166
  2167
  2168
Solution:
  2169
# We need to find the integer values of x for which sqrt(x) > 3 and sqrt(x) < 4. To
  2170
    this end, we can count the number of integers in the range (9, 16) directly.
  2171
result = 0
  2172
for x in range(10, 16):
  2173
   if 9 < x < 16:
  2174
      result += 1
  2175
print(result)
  2179
F.5 Evaluation without Toolsets
  2178
```

Listing 8: Prompt used for evaluation (setting without toolsets)

2179 2180 Read the following questions and answer them. For each question, you are required to write a Python program to solve it. 2181 Please note that we provide you several functions for each question. You could ( 2182 optionally, not required) call the functions to help you to solve the question 2183 if necessary. 2184 Note that the last line of your program should be a 'print' command to print the 2185 final answer 2186 2187 \_\_\_\_\_ 2188 Question: 2189 What is the 100th digit to the right of the decimal point in the decimal 2190 representation of \$\\frac{13}{90}\$? 2191 2192 2193

```
Solution:
```

```
2194
           from decimal import Decimal, getcontext
2195
           # Set the precision to 101 (100 digits after decimal + 1 digit before decimal)
2196
2197
           getcontext().prec = 101
2198
2199
           # Calculate the decimal representation of 13/90
2200
           dec = Decimal(13) / Decimal(90)
2201
2202
           # Convert the decimal to a string
2203
           dec_str = str(dec)
2204
           # Get the 100th digit to the right of the decimal point
2205
2206
           digit_100th = dec_str[101]
2207
2208
           print(digit_100th)
2209
2210
           _____
           Ouestion:
2211
2212
           The square root of x is greater than 3 and less than 4. How many integer values of
2213
                $x$ satisfy this condition?
2214
2215
           Solution:
2216
           # We need to find the integer values of x for which sqrt(x) > 3 and sqrt(x) < 4. To
               this end, we can count the number of integers in the range (9, 16) directly.
2217
2218
           result = 0
2219
           for x in range(10, 16):
              if 9 < x < 16:
2220
                 result += 1
2221
           print(result)
2223
```