## Learning to Ask: When LLM Agents Meet Unclear Instruction

**Anonymous ACL submission** 

#### Abstract

Equipped with the capability to call functions, modern LLM agents can leverage external tools for addressing a range of tasks unattainable through language skills alone. However, the effective execution of these tools relies heavily not just on the advanced capabilities of LLM agents but also on precise user instructions, which often cannot be ensured in the real world. To evaluate the performance of LLM agents tool-use under imperfect instructions, we meticulously examine the real-world instructions queried from users, analyze the error patterns, and build a challenging tool-use benchmark called Noisy ToolBench (Noisy-ToolBench). We find that due to the next-token prediction training objective, LLM agents tend to arbitrarily generate the missed argument, which may lead to hallucinations and risks. To address this issue, we propose a novel framework, Ask-when-Needed (AwN), which prompts LLM agents to ask questions to users whenever they encounter obstacles due to unclear instructions. Moreover, to reduce the manual labor involved in user-LLM interaction and assess LLM agents' performance in tool utilization from both accuracy and efficiency perspectives, we design an automated evaluation tool named ToolEvaluator. Our experiments demonstrate that the AwN significantly outperforms existing frameworks for tool learning in the NoisyToolBench. We will release all related code and datasets to support future research.

#### 1 Introduction

005

011

012

015

017

022

035

040

042

043

LLMs have undergone remarkable development since OpenAI introduced ChatGPT-3.5 (Bang et al., 2023). This model demonstrates a significant advancement in solving multiple tasks, including code generation (Dong et al., 2023; Sakib et al., 2023; Feng et al., 2023), machine translation (Jiao et al., 2023; Peng et al., 2023), even game playing (Wu et al., 2024). However, despite their impressive capabilities, LLMs often struggle with complex computations and delivering accurate, timely information (Qu et al., 2024). Tool learning emerges as a promising solution to mitigate these limitations of LLMs by enabling dynamic interaction with external tools (Schick et al., 2024). 044

045

046

047

051

055

058

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

081

The incorporation of tool usage capabilities marks a pivotal step towards enhancing the intelligence of LLMs, pushing them closer to exhibiting human-like intelligence. The integration of tool usage allows LLMs to perform a broader array of complex and varied tasks, including managing emails, designing presentations, and browsing the web to gather real-time information. For example, LLMs can perform complex calculations using a calculator tool, access real-time weather updates through weather APIs, and execute programming code via interpreters (Qin et al., 2023a; Schick et al., 2024; Mialon et al., 2023; Yang et al., 2023a). Toolformer (Schick et al., 2024) is a pioneering work in empowering language models with self-learning capabilities for tool usage. Then, significant research efforts have been directed toward accessing a wider variety of tools or using multiple tools simultaneously to resolve a single query, such as Gorilla(Patil et al., 2023), RestGPT (Song et al., 2023) and ToolLLM (Qin et al., 2023b).

Despite the significant strides made, existing frameworks and benchmarks often operate under the assumption that user instructions are always explicit and unambiguous, a premise that diverges from real-world scenarios (Qin et al., 2023a; Song et al., 2023; Patil et al., 2023). Due to the feature of API calls, it requires precise user instructions since the arguments for the function call can hardly tolerate ambiguity. We find that due to the next-token prediction training objective, LLMs tend to arbitrarily generate the missed argument, which may lead to hallucinations and risks (as the example shown in Figure 1a). Furthermore, as the tasks assigned to LLMs grow in complexity, multiple and sequential API calls are needed to resolve a single task. This



(a) The execution process of previous frameworks.

(b) The execution process of our framework.

Figure 1: The motivating example of our Ask-when-Needed (AwN) framework.

complexity amplifies the challenge, as any error in the sequence of API calls can culminate in an outcome that strays from the user's original intention.Hence, LLMs tool-use under unclear instruction is an important but rarely investigated direction.

To address this oversight, we conduct a systematic analysis of actual user instructions, identifying and categorizing potential issues into several key areas. These include instructions lacking essential information, instructions with ambiguous references, instructions containing inaccuracies, and instructions that are unfeasible for LLMs to execute due to the limitations of the tools available. Building on this observation, we meticulously design a noisy instruction benchmark, named NoisyTool-Bench, which is pimarily used for assessing the capability of LLMs to detect ambiguities in user queries and to pose relevant questions for clarification accordingly. Specifically, NoisyToolBench includes a collection of provided APIs, ambiguous queries, anticipated questions for clarification, and the corresponding responses.

094

100

101

102

103

104

105

106

107

109

110

111

112

113

114

115

116

117

118

119

120

121

122

To improve the performance of LLMs tool-use under unclear instructions, we propose a novel framework called Ask-when-Needed (AwN). Our key insight is encouraging LLMs to proactively ask questions to seek clarifications from users when uncertainties arise during instruction execution. By facilitating dialogue throughout the process, our method aims to ensure the accurate invocation of functions (See Figure 1b)

To evaluate the performance of LLMs tool-use under unclear instruction, we design several innovative metrics from the accuracy and efficiency perspectives. For accuracy, we measure the LLMs' proficiency in asking appropriate clarifying questions, their ability to execute the correct function calls, and their success in delivering final responses that meet the users' needs. For efficiency, we calculate the average number of redundant asked questions and the average number of actions required to complete the instruction. An ideal LLM should achieve higher accuracy with fewer number of queries. Recognizing the labour-intensive nature of manually communicating with LLMs and verifying all execution results, we also innovatively design an automatic evaluation system, ToolEvaluator, to streamline the whole process. ToolEvaluator leverages the advanced problem-solving capabilities of GPT-40 to communicate with LLMs and automatically evaluate the performance of LLMs' tool-using under unclear instruction. Our experiments on 6 LLMs and 2 tool-using frameworks demonstrate that the AwN significantly outperforms existing baseline methods for tool learning in the Noisy-ToolBench.

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

152

153

154

155

156

157

158

159

The key contributions of this research are summarized as follows:

- We conduct a systematic study on real-world user instruction for tool utilization and categorize the prevalent issues into four distinct categories.
- We create and release a novel benchmark, Noisy-ToolBench, which can be used to evaluate the performance of LLMs' tool-using under imperfect user instruction.
- We design five evaluation metrics from both accuracy and efficiency perspectives and introduce an automatic evaluation system, ToolEvaluator, that can proxy users to interact and assess LLMs.
- We introduce a novel framework, named AwN method, to prompt LLMs to actively ask questions to request clarifications from users when facing uncertainties. Experimental results show that AwN can significantly improve the LLMs' tool-using under unclear instructions.

#### 161 162

165

166

167

168

169

170

171

173

174

176

177

178

179

181

182

183

186

190

191

194

195

196

2 **Related Works** 

Tool Learning for LLMs. LLMs have recently made significant advancements, with ChatGPT being recognized as a major step towards achieving AGI (Wu et al., 2023; Lund and Wang, 2023; Jiao et al., 2023). These LLMs possess strong reasoning capabilities, enabling them to perform increasingly complex tasks (Liu et al., 2023). However, to progress further towards AGI, it is crucial for LLMs to master the utilization of tools. Toolformer is the first innovative AI model designed to use several specialized tools, such as a web browser, a code interpreter, and a language translator, within a single framework (Schick et al., 2023). The model's ability to seamlessly switch between these tools and apply them contextually represents a significant ad-175 vancement in AI capabilities. Recent studies like RestGPT (Song et al., 2023) and ToolLLM (Qin et al., 2023b), have connected LLMs with real-life Application Programming Interfaces (APIs), allowing them to sequentially employ multiple external tools to solve user queries. The tool-augmented approach empowers LLMs to use various kinds of tools to do more sophisticated tasks, showcasing an enhanced level of capability compared to pure LLMs. Besides, API-Bank (Li et al., 2023), ToolAlpaca (Tang et al., 2023), ToolBench (Qin et al., 2023b), ToolQA (Zhuang et al., 2023) and RestBench (Song et al., 2023) are exemplary benchmarks to systematically evaluate the performance of tool-augmented LLMs performance in response to user's queries. However, current models often ignore the situations in which users might not give exact instructions, which can result in the tools not 193 working properly. Thus, our study aims to tackle this specific challenge by developing a new benchmark specifically for ambiguous instructions.

Prompting LLMs for Decision Making. In cer-197 tain situations, addressing user queries may require more than a single API call. This necessi-199 tates the effective division of the overarching task into smaller, more manageable components, which 201 presents a significant challenge. Prior research has 202 focused extensively on enhancing LLMs's ability to effectively plan and execute complex tasks. The 'Chain of Thought' prompting approach facilitates 206 advanced reasoning by introducing intermediate steps in the reasoning process (Wei et al., 2022). The ReAct methodology improves the integration of reasoning and action, enabling LLMs to take informed actions based on environmental feedback 210

(Yao et al., 2022). Meanwhile, Reflexion is designed to reduce errors in the reasoning process by revisiting and learning from previous mistakes (Shinn et al., 2023). DFSDT expands upon Reflexion, allowing LLMs to evaluate various options and choose the most viable path (Qin et al., 2023b). In our work, we innovatively involve users in the process of executing instructions. Our approach, referred to as AwN, motivates LLMs to consider the necessity of requesting further information from users during each tool invocation round. This strategy aims at clarifying users' ambiguous instructions to help execute the tasks in alignment with the users' intentions.

211

212

213

214

215

216

217

218

219

220

221

222

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

257

258

259

260

Learning to Ask. Since user queries may not always be clear, and the execution of LLMs may encounter uncertainties and ambiguities, learning to ask questions has emerged as a challenging yet crucial research area (Rao and Daumé III, 2018; Kuhn et al., 2022; Andukuri et al., 2024). For example, some researchers introduce a learning framework that empowers an embodied visual navigation agent to proactively seek assistance(Zhang et al., 2023). Recently, similar ideas have been adopted in the software engineering, leveraging a communicator to enhance the reliability and quality of generated code (Wu, 2023). Our work focuses on the tool-learning scenario, which is more sensitive to the user's unclear query. A concurrent study (Qian et al., 2024) also focuses on the reliability of tool-learning systems under unclear instruction. However, they did not systematically examine realworld user behavior, leading to the limited and biased nature of their dataset that doesn't accurately capture user errors. Our research addresses this shortfall by starting with a user analysis. Additionally, Qian's methodology depends significantly on human manual interaction and assessment of LLM performances, which is time-consuming and hard to reproduce. In contrast, we introduce an automated evaluation method that can proxy humans to communicate with and automatically evaluate the performance of LLMs.

#### 3 **Noisy ToolBench**

Several tool-learning benchmarks have been introduced to assess LLMs' ability in tool utilization. However, these benchmarks overlook the potential ambiguity in users' instruction, which might hinder LLMs from executing tasks as intended by the user. For instance, as depicted in Figure 1a, if a

user inquires, "How is today's weather" without specifying the location, LLMs cannot accurately activate the APIs to fetch the correct weather information. This scenario underscores the critical role of interaction between users and LLMs in executing instructions accurately. However, previous tool-learning benchmarks only contain perfect user instruction in a one-query-one-execution manner.

261

262

263

267

270

272

274

275

278

279

281

287

289

290

291

296

297

298

301

303

307

308

311

To create a realistic benchmark for ambiguous instructions, the initial step involves a systematic examination of the common errors in user instructions that could complicate correct execution by LLMs. Therefore, we first collect real-world user instructions that are problematic. Then, we classify these instructions into various categories based on their characteristics. Lastly, we manually create our dataset, ensuring it reflects the distribution of errors found in the real-world user instructions.

#### 3.1 User Instruction Analysis

To analyze the issues in real-world user instruction, we recruit human annotators to write user queries according to the API provided. Firstly, we select 100 APIs from the ToolBench (Qin et al., 2023b), containing real-world RESTful APIs spanning 49 categories, ranging from sports to finance. Secondly, we hire 10 volunteers, who have a Bachelor's degree, are proficient in English, and have experience using LLMs. We provide them with the 100 APIs, and then ask them to write down an instruction to prompt LLMs to call each API, ending up with 1000 user queries. Finally, we manually identify the problematic user queries and categorized them as follows.

• Instructions Missing Key Information (IMKI): These are user instructions that omit crucial details necessary for the successful execution of a function. An example of IMKI would be, "Set an alarm to wake me up" without providing a specific time. Asking for more information is needed when encountering this issue.

• Instructions with Multiple References (IMR): These user instructions include elements that can be interpreted in several ways, potentially leading to confusion for LLMs in understanding the user's actual intent. For example, an IMR instance is "I want to know the director of the movie 'The Matrix'," where the ambiguity arises because there are multiple versions of 'The Matrix', each possibly having a different director. This issue is similar to IMKI but is more subtle and difficult to detect. Pointing out potential ref-

Type of Issue	Ratio
Instructions Missing Key Information (IMKI)	56.0%
Instructions with Multiple References (IMR)	11.3%
Instructions with Errors (IwE)	17.3%
Instructions Beyond Tool Capabilities (IBTC)	15.3%

Table 1: Distribution of problematic instructions.

erences and asking for clarification are needed when encountering this issue. 312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

351

352

354

355

- Instructions with Errors (IwE): This category consists of user instructions that contain the necessary information for executing a function, but the information is incorrect. An example of IWE is, "Please help me to log in to my Twitter. My user account is 'abcde@gmail.com' and the password is '123456'," where the user might have provided the wrong account details or password due to typographical errors. Asking for the correct information is needed when encountering this issue.
- Instructions Beyond Tool Capabilities (IBTC): These are user instructions that request actions or answers beyond what LLMs can achieve with the available APIs. In such cases, the existing toolaugmented LLM frameworks might randomly choose an available API, leading to an incorrect function call. This scenario highlights the need for LLMs to recognize their limitations in tool usage. Telling the user that the query is beyond the capabilities and refusing to generate API calls are needed when encountering this issue.

Table 1 shows the ratio of the four issues, where the most common issue in the instructions is "Instructions Missing Key Information", with a significant 56.0% of all errors. This issue is a clear indication that users often do not provide adequate information to effectively use the APIs. Additionally, issues such as "Instructions with Errors" and "Instructions Beyond Tool Capabilities" were identified at rates of 17.3% and 15.3%, respectively.

#### 3.2 Data Construction

Our user instruction analysis reveals that there are four kinds of instruction issues that may lead to LLMs' tool utilization failures: Instructions Missing Key Information (IMKI), Instructions with Multiple References (IMR), Instructions with Errors (IwE), and Instructions Beyond Tool Capabilities (IBTC). So, we build our benchmark with the four issues by intentionally modifying the problem-free instructions from well-established datasets to problematic ones. We first select 200 data with problem-



Figure 2: The comparison of our QwN prompting compared with original CoT/ReAct Prompting

free instruction from ToolBench and then manually modify the user instructions to make them suffer from the four kinds of instruction issues. Then we annotate the expected questions that LLMs should ask when facing each imperfect user query, which will be used to measure whether LLMs can ask the right questions, as well as the answer to the question, which will be used to proxy the human responses. We conduct a two-round crossverification to ensure the quality of the annotation. Each data is annotated and verified by different people and any disagreement data will be re-annotated until reach a consensus. Finally, each data entry in NoisyToolBench has the following five components: the imperfect user query, the available APIs, the questions that LLMs should ideally ask, the answers to these questions, and the expected function calls along with their respective arguments.

357

361

366

371

373

375

376

377

384

388

#### 4 Ask-when-Needed Prompting

Previous approaches to tool-using often overlooked the importance of user engagement during the reasoning and planning stages. To address this oversight, we introduce a new prompting strategy named Ask-when-Needed (AwN). The key insight is prompting LLMs to detect the potential flaws in user instructions and proactively seek clarifications by asking questions before generating the API call.

AwN is built upon widely-used tool-using methods, such as CoT and ReAct. As in Figure 2, we introduce an additional step before the generation of API calls. This step involves presenting all available information to the LLMs, including the user query and API guideline, and prompting them to determine the adequacy and correctness of user instruction. If LLMs identify any missing argument needed for function execution based on the API's requirements, they are encouraged to ask questions to the user for this information. AwN prompts LLMs not to generate any API call until obtaining all the necessary information. In other words, only if no further information is needed, they can bypass the query step and directly initiate the API call. We also provide various kinds of specific instructions and demonstration examples for different kinds of instruction issues.

You are AutoGPT, tasked with processing user requests through a variety of APIs you have access to. Sometimes, the information provided unclear, incomplete, by users may be or incorrect. Your main responsibility is to if user's instructions are determine the sufficiently clear and detailed for effective use of the APIs. Here's your strategy: If user instructions are missing crucial 1. details for the APIs, pose a question to obtain the necessary information. If the user's instructions appear to be 2. incorrect, delve deeper by asking questions to clarify and rectify the details. If the user's request falls outside the 3. capabilities of your current APIs, notify them that vou're unable to meet the request by stating: "Due to the limitation of the toolset, I cannot solve the question".

#### **5** Experiments

In this section, we evaluate the performance of our Ask-when-Needed (AwN) prompting technique on the NoisyToolBench dataset. We first introduce the evaluation metrics, where we specify the criteria used to assess the effectiveness of AwN. Then, we describe the evaluation pipeline, detailing the step-by-step process employed to measure AwN's performance. Lastly, we discuss the main experiments, presenting the results and findings from our comprehensive testing of the AwN technique.

#### 5.1 Evaluation Metrics

We evaluate the performance of LLMs' tool-using under unclear instructions from two perspectives: accuracy and efficiency. The accuracy assessment aims to measure the LLMs' capability to make correct decisions during the instruction execution phase and to generate accurate final answers. In contrast, the efficiency assessment focuses on the number of redundant decisions made by the LLMs, considering that unnecessary communication could lead to a waste of processing time. Specifically, we 402

389

390

391

392

393

394

395

396

397

398

399

400

401

- 403 404
- 405 406 407

408

409

410 411

412 413

414 415

416

417

418

419

420

421

422



Figure 3: Illustration of the Auto-Interaction module.

design the following five metrics:

424

425

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

453

454

455

456

457

458

459

460

461

- Accuracy 1 (A1). A1 evaluates the capability of LLMs to ask the anticipated questions that pinpoint the ambiguous elements in user instructions. A1 is considered a success if the LLMs manage to ask the correct questions at any point. Conversely, it is deemed a failure if they do not.
  - Accuracy 2 (A2). A2 assesses the ability of LLMs to use all available information to invoke the correct API calls. It is deemed a success if the LLMs call all the anticipated APIs with the correct arguments. If they fail to do so, it is considered a failure.
  - Accuracy 3 (A3). A3 measures the ability of LLMs to extract the anticipated information from previous API calls to fulfill the user's instructions. This is achieved and considered a success if the user's instructions are successfully executed. If not, it is regarded as a failure.
- Average Redundant Asked questions (Re). This metric evaluates the quantity of irrelevant or redundant questions asked by LLMs during the instruction process. Irrelevant questions are those that do not meet the initial expectations of the query, and redundant questions include those that are repetitive or have previously been asked. This metric is crucial for assessing the LLMs' ability to precisely identify the ambiguous aspects of user instructions and to formulate appropriate questions to clarify these uncertainties. The larger the value, the worse the performance.
- **Steps.** Steps quantifies the average number of actions required to complete an instruction, including inference generation, asking questions, and conducting API calls. A smaller number indicates fewer unnecessary steps in the instruction execution process, signifying a more efficient and direct approach to accomplishing the task.

#### 5.2 Auto-Evaluation Pipeline

To assess how LLMs perform under unclear instructions, interacting with LLMs and making assessments are needed. Previous work employs individuals to interact with and evaluate LLMs throughout the entire evaluation process, which is inefficient and not reproducible. To address this, we design an automated evaluation method named ToolEvaluator to proxy this process. ToolEvaluator can automatically interact with LLMs and assess their performances. 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

493

494

495

496

497

498

**Auto-Interaction.** ToolEvaluator can proxy the user's communication with LLMs. When LLMs post a question, ToolEvaluator calculates the semantic similarity between the asked question and the expected question by the sentencetransformer (Reimers and Gurevych, 2019). If the similarity is higher than a threshold, ToolEvaluator replies with the predefined answer to the LLMs. Otherwise, this question is treated as an irrelevant question and ToolEvaluator replies with a standard reply of "Sorry, I cannot provide additional information about this.". This approach streamlines the evaluation process by reducing the need for human interaction with LLMs, as illustrated in Figure 3.

**Auto-Assessment.** ToolEvaluator can also automatically assess how well LLMs perform under ambiguous instructions according to the five metrics introduced above. A1 measures whether LLMs can ask the right question. ToolEvaluator calculates the semantic similarity between the LLMs-asked question and the expected question to asses A1. A2 measures whether LLMs can conduct correct API calls. Following the previous works (Yang et al., 2023b; Chiang and yi Lee, 2023; Wang et al., 2023; Yuan et al., 2023), ToolEvaluator adopts GPT-4o as a judge to identify whether the generated API

Model	Framework	IMKI			IMR			IwE			IBTC
		A1(%)	A2(%)	A3(%)	A1(%)	A2(%)	A3(%)	A1(%)	A2(%)	A3(%)	A1(%)
gpt-3.5	СоТ	0.74	0.36	0.22	0.20	0.24	0.12	0.5	0.24	0.16	0.38
	+ AwN	0.74	0.44	0.24	0.86	0.46	0.20	0.74	0.48	0.28	0.48
	DFSDT	0.64	0.16	0.12	0.60	0.18	0.16	0.48	0.14	0.14	0.46
	+ AwN	0.88	0.52	0.46	0.88	0.56	0.48	0.72	0.42	0.36	0.64
gpt-4	СоТ	0.74	0.48	0.32	0.72	0.52	0.36	0.52	0.26	0.24	0.34
	+ AwN	0.94	0.62	0.50	0.76	0.44	0.38	0.48	0.34	0.34	0.94
	DFSDT	0.82	0.16	0.16	0.70	0.28	0.26	0.54	0.12	0.10	0.54
	+ AwN	0.80	0.56	0.48	0.80	0.50	0.44	0.52	0.38	0.36	0.94
gpt-4o	СоТ	0.52	0.48	0.34	0.18	0.28	0.16	0.12	0.12	0.10	0.10
	+ AwN	0.90	0.58	0.36	0.80	0.46	0.30	0.60	0.44	0.32	0.92
	DFSDT	0.58	0.20	0.18	0.26	0.18	0.16	0.18	0.06	0.04	0.08
	+ AwN	0.88	0.60	0.46	0.90	0.52	0.36	0.64	0.46	0.38	0.94
deepseek-v3	СоТ	0.44	0.40	0.20	0.24	0.28	0.24	0.10	0.14	0.14	0.30
	+ AwN	0.70	0.52	0.36	0.70	0.54	0.46	0.40	0.30	0.26	0.98
	DFSDT	0.42	0.30	0.26	0.60	0.20	0.18	0.22	0.12	0.12	0.48
	+ AwN	0.72	0.52	0.42	0.82	0.52	0.48	0.54	0.38	0.36	0.98
gemini-1.5	СоТ	0.22	0.18	0.10	0.22	0.10	0.02	0.08	0.12	0.06	0.52
	+ AwN	0.86	0.40	0.18	0.74	0.24	0.08	0.58	0.28	0.22	0.68
	DFSDT	0.62	0.02	0.02	0.6	0.08	0.04	0.36	0.06	0.02	0.48
	+ AwN	0.82	0.40	0.12	0.76	0.28	0.04	0.66	0.36	0.26	0.70
claude-3.5	CoT	0.24	0.26	0.20	0.12	0.28	0.24	0.08	0.26	0.24	0.30
	+ AwN	0.54	0.5	0.5	0.32	0.30	0.24	0.34	0.34	0.26	0.88
	DFSDT	0.26	0.18	0.14	0.12	0.18	0.18	0.12	0.20	0.18	0.62
	+ AwN	0.52	0.44	0.42	0.32	0.30	0.18	0.36	0.36	0.30	0.86

Table 2: Assessing the accuracy of various LLMs using different prompting methods in our benchmark.

calls are the same as the expected API calls. A3 measures whether LLMs can correctly generate the final answer. ToolEvaluator adopts GPT-4o as a judge to identify whether the final answer aligns with the user's intent. For measuring the efficiency, ToolEvaluator counts the number of generated irrelevant questions as Re and counts the total number of actions during the process as Steps. All the details can be found in the Appendix due to the space limitation.

#### 509 5.3 The Effectiveness of ToolEvaluator

499

500

505

506

507

508

Since ToolEvaluator is an automatic evaluation 510 511 method, the evaluation can be inaccurate due to the imperfect nature of AI techniques, such as sen-512 tence transformer or GPT-40 as the judge. In this 513 section, we conduct a human annotation to validate 514 the effectiveness of ToolEvaluator. Specifically, 515 516 we randomly select 50 cases and ask annotators to assess the accuracy and efficiency, according to 517 the evaluation metrics mentioned above. Then we 518 compare the assessment results from ToolEvalua-519 tor and human annotators. ToolEvaluator achieves 520

90% accuracy, indicating its effectiveness.

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

#### 5.4 Experimental Setup

We evaluated the performance of AwN against two baseline methods, chain-of-thought (CoT) (Wei et al., 2022) and depth-first search-based decision tree (DFSDT) (Qin et al., 2023b), which are two widely-used tool-learning methods. All the experiments are conducted with several LLMs as engines, gpt-3.5-turbo-0125, gpt-4-turbo-2024-04-09, gpt-4o-2024-11-20, deepseek-v3, gemini-1.5flash-latest and claude-3-5-haiku-20241022, using the default setting. Since an ideal reaction under Instructions Beyond Tool Capabilities (IBTC) is telling the user that the query is beyond the capabilities and refusing to generate API calls, its performance in A2 and A3 are measured neither.

#### 5.5 Main Result

We evaluate the performance of AwN as well as the baseline methods on our NoisyToolBench dataset. The accuracy-related results are shown in Table 2 and the efficiency-related results are in Table 3.

Model	FrWork	IMKI		I	MR	I	wE	IBTC		
liouer		Re	Steps	Re	Steps	Re	Steps	Re	Steps	
gpt-3.5	CoT	-	4.46	-	4.02	-	3.90	-	2.10	
	+ AwN	0.66	5.36	1.1	5.98	0.76	5.08	-	2.40	
	DFSDT	-	12.82	-	12.80	-	13.82	-	5.50	
	+ AwN	1.44	16.94	0.98	11.24	0.94	11.68	-	3.94	
gpt-4	CoT	-	4.00	-	3.98	-	3.34	-	2.04	
	+ AwN	0.16	3.94	0.20	3.94	0.36	3.46	-	1.16	
	DFSDT	-	83.96	-	21.04	-	22.40	-	4.06	
	+ AwN	0.48	9.82	0.74	13.08	0.62	9.42	-	2.10	
gpt-40	CoT	-	3.00	-	2.98	-	2.48	-	1.28	
	+ AwN	0.62	3.86	0.70	3.96	0.46	3.18	-	1.10	
	DFSDT	-	5.98	-	9.58	-	5.78	-	8.98	
	+ AwN	0.86	6.70	1.18	7.68	0.88	8.56	-	1.14	
deepseek-v3	CoT	-	4.20	-	3.52	-	3.12	-	1.18	
	+ AwN	0.20	3.88	0.06	3.60	0.04	2.92	-	1.10	
	DFSDT	-	59.08	-	41.70	-	24.24	-	11.64	
	+ AwN	1.16	15.86	1.80	24.60	1.20	11.82	-	1.32	
gemini-1.5	CoT	-	4.00	-	4.12	-	2.86	-	4.52	
	+ AwN	0.42	6.44	0.68	6.36	0.48	4.54	-	1.46	
	DFSDT	-	750.80	-	685.00	-	725.14	-	559.78	
	+ AwN	5.34	445.16	9.08	532.56	1.94	411.18	-	1.46	
claude-3.5	CoT	-	2.64	-	3.40	-	3.04	-	1.90	
	+ AwN	0.18	3.74	0.33	1.03	0.36	3.76	-	1.68	
	DFSDT	-	3.34	-	9.64	-	5.98	-	4.26	
	+ AwN	0.76	6.74	0.80	17.46	1.04	13.08	-	2.76	

Table 3: Assessing the efficiency of various LLMs using different prompting methods in our benchmark.

543

544

545

546

549

550

553

554

555

556

562

563

565

566

AwN enhances the capability of LLM Agents to ask pertinent questions across different issues. For example, as is shown in Table 2, AwN improved the A1 scores from 0.52 to 0.90, from 0.18 to 0.80, from 0.12 to 0.60, and from 0.10 to 0.92 for gpt-4o-based CoT as well as from 0.58 to 0.88, from 0.26 to 0.90, from 0.18 to 0.64 and from 0.08 to 0.94 for gpt-4o-based DFSDT.

Asking the right question leads to the better generation and execution of API calls. Besides the significant improvements on A1, AwN also achieves considerable performance in generating correct API calls (A2) and returning the expected final answer (A3). For example, AwN improved the A2 scores from 0.48 to 0.58, from 0.28 to 0.46, from 0.12 to 0.44 for gpt-4o-based CoT as well as from 0.20 to 0.60, from 0.18 to 0.52, from 0.06 to 0.46 for gpt-4o-based DFSDT.

AwN can improve most of the LLM agents without generating excessive unnecessary questions. As is shown in Table 3, AwN only leads to 0.16, 0.20, and 0.36 redundant questions for gpt-4-based-CoT, as well as 0.48, 0.74, and 0.62 redundant questions for gpt-4-based-DFSDT.

However, a few LLM agents tend to ask more

irrelevant or redundant questions, as indicated by the higher Re scores in Table 3. For example, in Gemini-1.5-based DFSDT, where the average number of redundant questions is 5.34, 9.08, and 1.94. This suggests that while the AwN aids in identifying and addressing ambiguities in user instructions, it also leads to a less efficient querying process. 567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

AwN can reduce the average cost of LLM's tool-using. The average number of steps measures the cost of LLMs' tool-using. As is shown in Table 3, adopting AwN can reduce the number of actions, especially for gpt-4-based DFSDT, deepseek-v3-based-DFSDT and gemini-1.5-based-DFSDT. Although AwN can lead to a higher cost for a few LLM agents, such as claude-3.5, considering the significant performance improvements achieved, the moderate increase in cost is justifiable and worthwhile.

### 6 Conclusion

This study explores how unclear user instructions hinder modern LLMs' tool usage. To investigating the common error patterns in real-world instructions, we propose: (1) Noisy ToolBench (NoisyToolBench), a novel benchmark for evaluating LLM performance under ambiguous instructions; (2) Ask-when-Needed (AwN), an innovative approach enabling LLMs to request clarification when uncertain; and (3) an automated evaluator (ToolEvaluator) to assess accuracy and efficiency. Experimental results show that the AwN algorithm markedly surpasses existing methods in the Noisy-ToolBench dataset and significantly improves the performance of LLMs' tool-using under unclear user instructions.

#### Limitations

This paper has two limitations:

- 1. Although AwN can improve the performance, there is still a big gap to perfect. We hope that this work can serve as the first stepping stone, inspiring future researchers to delve deeper into this field of study.
- 2. The automatic evaluation process is not 100% accurate, leading to some potential false negatives and false positives. In the future, more efforts are needed to build a more reliable auto-evaluation method.

615

616

617

618

619

621

622

624

626

627

628

630

631

632

633

634

637

645

651

654

657

References

- Chinmaya Andukuri, Jan-Philipp Fränken, Tobias Gerstenberg, and Noah D Goodman. 2024. Star-gate: Teaching language models to ask clarifying questions. arXiv preprint arXiv:2403.19154.
  - Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, et al. 2023. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. *arXiv preprint arXiv:2302.04023*.
  - Cheng-Han Chiang and Hung yi Lee. 2023. Can large language models be an alternative to human evaluations? In Annual Meeting of the Association for Computational Linguistics.
  - Yihong Dong, Xue Jiang, Zhi Jin, and Ge Li. 2023. Self-collaboration code generation via chatgpt. *arXiv preprint arXiv:2304.07590*.
  - Yunhe Feng, Sreecharan Vanam, Manasa Cherukupally, Weijian Zheng, Meikang Qiu, and Haihua Chen.
    2023. Investigating code generation performance of chatgpt with crowdsourcing social data. In 2023 IEEE 47th Annual Computers, Software, and Applications Conference (COMPSAC), pages 876–885. IEEE.
  - Wenxiang Jiao, Wenxuan Wang, Jen-tse Huang, Xing Wang, Shuming Shi, and Zhaopeng Tu. 2023. Is chatgpt a good translator? yes with gpt-4 as the engine. *arXiv preprint arXiv:2301.08745*.
  - Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2022. Clam: Selective clarification for ambiguous questions with generative language models. *arXiv preprint arXiv:2212.07769*.
  - Minghao Li, Feifan Song, Bowen Yu, Haiyang Yu, Zhoujun Li, Fei Huang, and Yongbin Li. 2023. Apibank: A benchmark for tool-augmented llms. *arXiv preprint arXiv:2304.08244*.
  - Hanmeng Liu, Ruoxi Ning, Zhiyang Teng, Jian Liu, Qiji Zhou, and Yue Zhang. 2023. Evaluating the logical reasoning ability of chatgpt and gpt-4. *arXiv preprint arXiv:2304.03439*.
  - Brady D Lund and Ting Wang. 2023. Chatting about chatgpt: how may ai and gpt impact academia and libraries? *Library Hi Tech News*, 40(3):26–29.
  - Grégoire Mialon, Roberto Dessì, Maria Lomeli, Christoforos Nalmpantis, Ram Pasunuru, Roberta Raileanu, Baptiste Rozière, Timo Schick, Jane Dwivedi-Yu, Asli Celikyilmaz, Edouard Grave, Yann LeCun, and Thomas Scialom. 2023. Augmented language models: a survey. *Preprint*, arXiv:2302.07842.
  - Shishir G. Patil, Tianjun Zhang, Xin Wang, and Joseph E. Gonzalez. 2023. Gorilla: Large language model connected with massive apis. *Preprint*, arXiv:2305.15334.

Keqin Peng, Liang Ding, Qihuang Zhong, Li Shen, Xuebo Liu, Min Zhang, Yuanxin Ouyang, and Dacheng Tao. 2023. Towards making the most of chatgpt for machine translation. *arXiv preprint arXiv:2303.13780*. 667

668

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

708

709

710

711

712

714

715

716

717

718

719

- Cheng Qian, Bingxiang He, Zhong Zhuang, Jia Deng, Yujia Qin, Xin Cong, Zhong Zhang, Jie Zhou, Yankai Lin, Zhiyuan Liu, and Maosong Sun. 2024. Tell me more! towards implicit user intention understanding of language model driven agents. *Preprint*, arXiv:2402.09205.
- 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. *Preprint*, arXiv:2304.08354.
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. 2023b. Toolllm: Facilitating large language models to master 16000+ real-world apis. *arXiv preprint arXiv:2307.16789*.
- Changle Qu, Sunhao Dai, Xiaochi Wei, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, Jun Xu, and Ji-Rong Wen. 2024. Tool learning with large language models: A survey. *arXiv preprint arXiv:2405.17935*.
- Sudha Rao and Hal Daumé III. 2018. Learning to ask good questions: Ranking clarification questions using neural expected value of perfect information. *arXiv preprint arXiv:1805.04655*.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. *EMNLP*.
- Fardin Ahsan Sakib, Saadat Hasan Khan, and AHM Karim. 2023. Extending the frontier of chatgpt: Code generation and debugging. *arXiv preprint arXiv:2307.08260*.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2024. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. *Preprint*, arXiv:2302.04761.

- 721 722
- 72 72
- 7: 7:
- 729 730 731
- 732 733 734 735
- 7 7
- 739 740
- 741 742 743
- 744 745 746
- 747
- 748 749
- 750 751
- 752 753
- 755

- 757 758 759
- 760
- 761
- 7
- 765
- 766

- 769 770
- 771 772

- Noah Shinn, Beck Labash, and Ashwin Gopinath. 2023. Reflexion: an autonomous agent with dynamic memory and self-reflection. *arXiv preprint arXiv:2303.11366*.
- 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 realworld restful apis. *Preprint*, arXiv:2306.06624.
- Qiaoyu Tang, Ziliang Deng, Hongyu Lin, Xianpei Han, Qiao Liang, and Le Sun. 2023. Toolalpaca: Generalized tool learning for language models with 3000 simulated cases. *arXiv preprint arXiv:2306.05301*.
- Wenxuan Wang, Zhaopeng Tu, Chang Chen, Youliang Yuan, Jen-tse Huang, Wenxiang Jiao, and Michael R Lyu. 2023. All languages matter: On the multilingual safety of large language models. *arXiv preprint arXiv:2310.00905*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.
- Jie JW Wu. 2023. Does asking clarifying questions increases confidence in generated code? on the communication skills of large language models. *Preprint*, arXiv:2308.13507.
- Tianyu Wu, Shizhu He, Jingping Liu, Siqi Sun, Kang Liu, Qing-Long Han, and Yang Tang. 2023. A brief overview of chatgpt: The history, status quo and potential future development. *IEEE/CAA Journal of Automatica Sinica*, 10(5):1122–1136.
- Yue Wu, Xuan Tang, Tom M. Mitchell, and Yuanzhi Li. 2024. Smartplay: A benchmark for llms as intelligent agents. *Preprint*, arXiv:2310.01557.
- Sherry Yang, Ofir Nachum, Yilun Du, Jason Wei, Pieter Abbeel, and Dale Schuurmans. 2023a. Foundation models for decision making: Problems, methods, and opportunities. *Preprint*, arXiv:2303.04129.
- Xianjun Yang, Xiao Wang, Qi Zhang, Linda Ruth Petzold, William Yang Wang, Xun Zhao, and Dahua Lin. 2023b. Shadow alignment: The ease of subverting safely-aligned language models. *ArXiv*, abs/2310.02949.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2022.
  React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*.
- Youliang Yuan, Wenxiang Jiao, Wenxuan Wang, Jen-tse Huang, Pinjia He, Shuming Shi, and Zhaopeng Tu. 2023. Gpt-4 is too smart to be safe: Stealthy chat with llms via cipher. *arXiv preprint arXiv:2308.06463*.

- Jenny Zhang, Samson Yu, Jiafei Duan, and Cheston Tan. 2023. Good time to ask: A learning framework for asking for help in embodied visual navigation. *Preprint*, arXiv:2206.10606.
- Yuchen Zhuang, Yue Yu, Kuan Wang, Haotian Sun, and Chao Zhang. 2023. Toolqa: A dataset for llm question answering with external tools. *Preprint*, arXiv:2306.13304.

774

780