

# 000 001 002 003 004 005 ADVANCING AND BENCHMARKING PERSONALIZED 006 TOOL INVOCATION FOR LLMs 007 008 009

010 **Anonymous authors**  
011 Paper under double-blind review  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053

## ABSTRACT

Tool invocation is a crucial mechanism for extending the capabilities of Large Language Models (LLMs) and has recently garnered significant attention. It enables LLMs to solve complex problems through tool calls while accessing up-to-date world knowledge. However, existing work primarily focuses on the fundamental ability of LLMs to invoke tools for problem-solving, without considering personalized constraints in tool invocation. In this work, we introduce the concept of Personalized Tool Invocation and define two key tasks: Tool Personalization and Parameter Personalization. Tool Personalization addresses user preferences when selecting among functionally similar tools, while Parameter Personalization considers cases where a user query lacks certain tool parameters, requiring the model to infer them from the user profile. To tackle these challenges, we propose **PTool**, a data synthesis framework designed for personalized tool invocation. Additionally, we construct **PTBench**, a benchmark to evaluate personalized tool invocation. We then fine-tune various open-source models, demonstrating the effectiveness of our framework and providing valuable insights. Our model, training data, and the benchmark will be publicly released upon acceptance.

## 1 INTRODUCTION

Recently, large language models (LLMs) have demonstrated remarkable capabilities in natural language processing tasks, particularly in human-computer interaction, where they can effectively comprehend user queries and provide reasonable responses (Zhao et al., 2023). However, the knowledge embedded within LLMs is not inherently up-to-date, as updating these models requires extensive retraining with large-scale data, which incurs significant time and economic costs. To equip LLMs with the ability to solve complex problems and access the latest information, tool invocation capabilities are essential. For instance, LLMs can leverage mathematical tools to decompose and solve intricate mathematical problems or utilize internet APIs (Liu et al., 2025; Qin et al., 2024) and search engines (Schick et al., 2024; Nakano et al., 2021) to retrieve the most recent knowledge.

Existing research on enhancing LLMs's tool invocation abilities primarily focuses on improving fundamental capabilities (Qin et al., 2024; Yan et al., 2024; Lin et al., 2024), such as ensuring adherence to the required syntax, comprehending tool functionalities, interpreting explicit user instructions, and extracting tool parameters. However, in real-world applications, user intents are often implicit rather than explicitly stated, requiring models to infer based on personalized profiles. Recent studies have begun to explore personalized tool invocation (Xu et al., 2025; Cheng et al., 2025), yet a systematic definition is still lacking. Existing work typically relies on idealized assumptions about user profiles—assuming that users' implicit preferences are observable to the model. However, such preferences are embedded in behavioral histories and are rarely provided explicitly in practice.

In this work, we first provide a systematic formulation of personalized tool invocation and highlight two common concepts that exemplify this challenge: (1) **Tool Personalization**. When multiple tools offer similar functionalities, users often exhibit specific preferences. For example, in online shopping, users may choose different platforms depending on their preferences for particular product categories. Some users may prioritize platforms with superior maintenance services when purchasing high-value electronic products, despite the higher cost, while preferring platforms with faster delivery when buying inexpensive daily necessities. Inferring such preferences necessitates reasoning from user attributes, such as age, interests, and purchasing behavior. (2) **Parameter Per-**

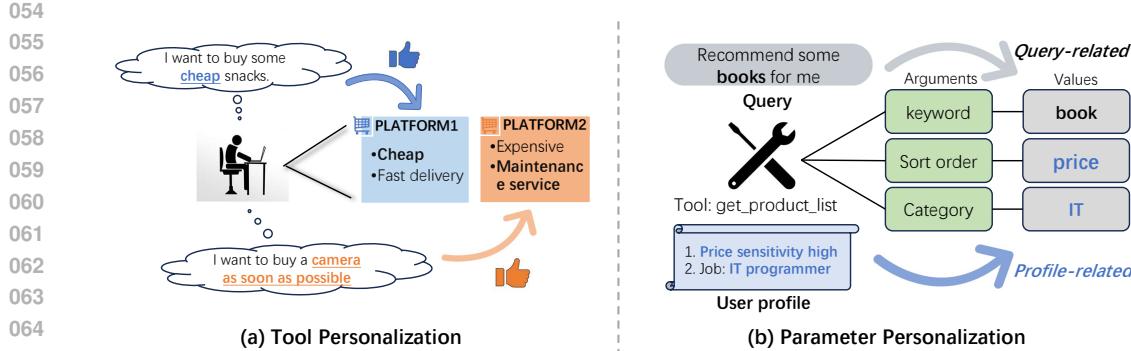


Figure 1: Examples of Personalized Tool Invocation. (a) Tool Personalization: Users may prefer different tools with similar functionalities. (b) Parameter Personalization: Certain tool parameters may be missing from the user’s query and must be inferred from the user’s profile.

**sonalization.** In everyday scenarios, users tend to express their needs concisely and omit crucial details. For instance, a user might simply request, “Order me a hamburger from KFC”, without specifying essential information such as the delivery address, recipient contact details, or preferred delivery time. This requires the model to infer the missing information from the user profile, such as the user’s work location, current time, and phone number, ensuring a seamless and accurate tool invocation process. To enhance and systematically evaluate a model’s ability in personalized tool invocation, we further introduce an automated data synthesis framework for this task, termed **PTool**, which consists of three key stages: tool generation, user profile construction, and user behavior simulation. Firstly, we consider several commonly used real-world scenarios, where each scenario contains multiple functionally similar platforms organized in a hierarchical tree structure. We then leverage an advanced LLM to recursively decompose platform functionalities using a depth-first expansion approach, progressively refining them until distinct tools are defined for each functional category. Secondly, we abstract and summarize platform features and API parameters to extract both basic user attributes and personalized characteristics, including psychological traits and behavioral tendencies. To construct a diverse set of user profiles, we employ a bottom-up clustering approach for feature induction and a top-down assignment strategy for attribute allocation. Finally, we exploit the role-playing capabilities of LLMs to simulate user behaviors based on the assigned user profiles, generating both historical interactions and potential user queries. To establish reliable ground-truth labels, we further integrate a multi-agent framework that conditions query generation on user profiles. Following manual review and annotation, we construct **Personalized ToolBench (PTBench)**, a benchmark designed to evaluate large models’ ability in personalized tool invocation, consisting of 1,301 high-quality annotated data samples. Key contributions are summarized as follows:

- We propose a role-playing-based paradigm for personalized tool invocation towards real-world applications, incorporating both user tool personalization and parameter personalization.
- We develop a systematic personalized data synthesis framework and construct PTBench, a benchmark for personalized tool invocation, enabling a comprehensive evaluation of the ability to invoke tools based on user information.
- We demonstrate that training open-source models on our synthesized dataset significantly improves personalized tool invocation capabilities, while also enhancing general tool invocation without compromising other general abilities.

## 2 RELATED WORK

### 2.1 TOOL INVOCATION

Tool invocation (also termed tool calling) involves tool selection from candidate tools and parameter extraction from queries. Existing works can be categorized into two tuning-free and tuning-based methods (Qu et al., 2025; Liu et al.). Tuning-free methods mainly rely on the prompt strategy

108 with few-shot learning, involving encouraging LLM to reason by providing examples (Yao et al.,  
 109 2022), rewriting tool documentation with LLMs to enhance the comprehension (Yuan et al., 2024),  
 110 summarizing tool description with more concise and precise sentence (Xu et al., 2024), leveraging  
 111 multi-agent collaboration to decompose the tool-calling task (Shi et al., 2024). Tuning-based meth-  
 112 ods leverage tool-learning samples to train existing LLMs, where the research problems comprise  
 113 data collection and training strategy. Toolformer (Schick et al., 2024) and ToolkenGPT (Hao et al.,  
 114 2024) add a special tool-related token into the vocabulary, switching the decoding process into tool  
 115 selection and calling. Some works leverage advanced LLM to synthesize tool-calling samples to  
 116 improve the tool-invocation ability of lightweight models, demonstrating the efficiency of the distil-  
 117 lation from advanced models (Qin et al., 2024; Yang et al., 2023b; Liu et al., 2025).

## 118 2.2 PERSONALIZED LLMs

120 Personalized LLMs represent LLMs that have been adapted to align with user preferences and char-  
 121 acteristics (Zhang et al., 2024c). Existing works mainly focus on the generation of personalized texts  
 122 or applications in information systems. LLMs are customized as personal conversational AI assis-  
 123 tants for various domains, including education (Kasneci et al., 2023; Dan et al., 2023; Park et al.,  
 124 2024), healthcare (Belyaeva et al., 2023; Abbasian et al., 2024; Jin et al., 2024), finance (Liu et al.,  
 125 2023; Lakkaraju et al., 2023), legal (Nguyen, 2023), and etc. User profiles are provided via prompts  
 126 or hidden representation, leading the model to generate personalized text in the dialog. Personalized  
 127 LLMs have been extensively applied in information systems such as recommender systems (Wu  
 128 et al., 2023; Chen et al., 2024). LLMs are leveraged as an augmentation module for traditional rec-  
 129 ommender systems, serving as the content interpreter (Bao et al., 2023; Li et al., 2023; Yang et al.,  
 130 2023a), the knowledge base (Xi et al., 2024; Wei et al., 2024), or the explainer (Lei et al., 2024;  
 131 Wang et al., 2023). Also, many works directly deploy LLMs as the direct recommenders via prompt  
 132 techniques (Lyu et al., 2024; Hou et al., 2024) or fine-tuning (Zhang et al.). **Recent work has begun**  
 133 **to investigate personalized tool invocation by treating user characteristics as explicit profile features**  
 134 **for selecting appropriate tools (Xu et al., 2025; Cheng et al., 2025).** However, these approaches over-  
 135 look a key challenge: in practical settings, users' implicit preferences are not directly observable by  
 136 the model and instead must be inferred from their behavioral histories.

## 137 3 PERSONALIZED TOOL INVOCATION

139 We innovatively consider a practical and high-demand scenario in LLM tool invocation: **personal-**  
 140 **ized tool invocation.** This scenario requires the model to leverage user-specific information when  
 141 selecting and configuring tools to address user needs. In this chapter, we formally define the task of  
 142 personalized tool invocation.

143 Given an LLM with model parameters  $\theta$ , the general tool invocation task requires the model, when  
 144 provided with a query  $q$  and a set of candidate tools  $T$ , to select the appropriate tool  $t^i$  and populate  
 145 its corresponding parameters  $a_1^i, \dots, a_m^i$ , forming the solution  $A = [(t^i, a_1^i, \dots, a_m^i), \dots]$ .

146 In conventional formulations of this task, correctness is typically determined by whether the selected  
 147 tool successfully resolves the query. However, this setting overlooks the fact that multiple tools may  
 148 solve one problem (e.g., APIs from different platforms with similar capabilities), and that users often  
 149 have preferences for certain tools—a concept we refer to as **tool personalization**, defined as follows:

150 **Definition 3.1. (Tool Personalization)** *User  $u$  prefers  $t^1$  for query  $q_1$  and  $t^2$  for query  $q_2$ , where  
 151  $q_1, q_2$  can be solved by both  $t^1$  and  $t^2$ :*

$$t^1 \succ_{(u, q_1)} t^2; \quad t^2 \succ_{(u, q_2)} t^1 \quad (1)$$

154 Moreover, in  $A$ , both tool selection and parameter values are determined solely based on the in-  
 155 formation contained in the query. For instance, consider the query: "Book me a flight from Los  
 156 Angeles to New York at 8:45 AM tomorrow". However, in real-world scenarios, users often do not  
 157 provide such detailed query information. Instead, they may omit certain essential details required  
 158 for tool invocation, meaning that the model cannot extract all necessary parameters from the query  
 159 alone. We refer to this personalized scenario as an **Parameter Personalization**, defined as follows:

160 **Definition 3.2. (Parameter Personalization)** *Given the profile of the user  $u$  as  $P_u$ , the query  $q$  and  
 161 the solution  $A$ , there exists value  $\alpha \in A$ ,  $\alpha \in P_u$  and  $\alpha \notin q$ . The phenomenon is called parameter  
 162 personalization, and the query  $q$  is called a profile-dependent query.*

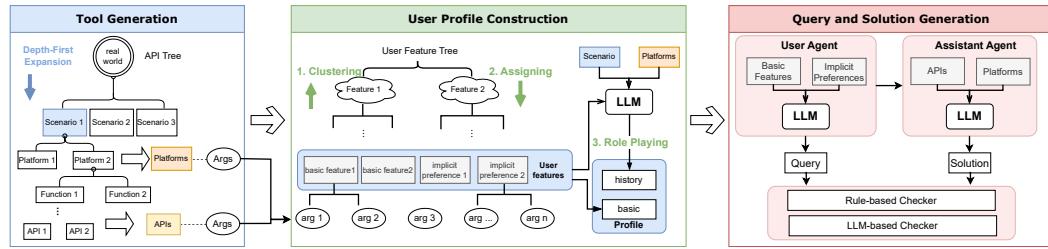


Figure 2: Framework of our personalized tool invocation data synthesis: PTool. The pipeline comprises three stages: Tool Generation, User Profile Generation and Query and Answer Generation.

## 4 PERSONALIZED TOOL INVOCATION DATA SYNTHESIS

To address the two challenges in personalized tool invocation mentioned above, we propose an automated data synthesis framework, PTool, for generating high-quality training and evaluation data for personalized tool invocation. The framework consists of three key stages: **Tool Generation**, **User Profile Construction**, and **Query and Solution Generation**, as illustrated in Figure 2. The detailed processes of each stage are described in the subsequent parts of this section.

### 4.1 TOOL GENERATION

To cover the majority of scenarios encountered in daily life, we first constructed a diversified tool library across multiple contexts. Inspired by existing work, we employed an advanced Large Language Model (LLM)-based data synthesis method to generate APIs. Similar to ToolACE, we also developed a structure akin to an API Tree, which allows for the generation of diverse tools.

Specifically, we initially define several demand scenarios from everyday life (e.g., shopping, food delivery, office) as the first-level nodes of the tree. Then, using a depth-first expansion approach, we iteratively refine the functionality at each node until we derive specific API descriptions as the leaf nodes. Notably, in order to generate data that enhances the model’s Tool Personalization capability, tools with similar functionalities are required. However, this API Tree expansion approach alone cannot achieve this. Therefore, at the second level of the tree expansion, we introduce the concept of platforms. For each scenario, we generated multiple platforms with distinct characteristics. For example, in the video entertainment scenario, platforms such as YouTube and TikTok were included, where YouTube focuses on long-form videos and TikTok emphasizes short, lifestyle-oriented clips. This enables us to obtain multiple tools with functionally interchangeable capabilities.

### 4.2 USER PROFILE CONSTRUCTION

Personalization requires constructing diverse and realistic user profiles. This process involves three key challenges: (1) defining feature sets relevant to tool invocation, ensuring a structured linkage between user traits and tool selection; (2) maintaining sufficient diversity across profiles to enable generalization to unseen users; and (3) ensuring that profiles contain only observable basic and behavioral information, without incorporating detailed psychological attributes.

**Bottom-up Feature Tree Construction.** To systematically define user profile features, we adopt a tool-driven hierarchical clustering approach. We construct a feature tree, where platform characteristics and tool parameters serve as leaf nodes. Using LLM-based clustering, we recursively merge semantically related parameters, summarizing them into higher-level features until the number of parent nodes at each level falls within a predefined threshold. Notably, we categorize features during initial clustering: explicit basic features (e.g., age, gender) are directly observable, while implicit preferences (e.g., shopping preferences) remain latent and are used in user behavior generation.

**Top-down Characteristic Assignment.** Once the user feature tree is constructed, we encounter the second issue: how to diversify the assignment of values to these features to generate distinct user profiles. When using an advanced LLM to assign  $N$  different user features, two options typically arise: one is to assign all features for a single user at a time and repeat this process  $N$  times; the other is to assign all features for  $N$  users in one pass. The first method incurs higher inference costs and

216 makes it challenging to avoid repetition across multiple generations, while the second is constrained  
 217 by the model’s context length limitation, especially when  $N$  or the number of features is large.  
 218 Therefore, we adopt a top-down hierarchical assignment based on the tree structure. Specifically,  
 219 for nodes at the  $l$ -th layer, we assign  $k_l$  different values simultaneously, and for the  $(l + 1)$ -th layer  
 220 nodes, the model generates  $k_{l+1}$  different values for each parent node’s feature value. Thus, for a  
 221 user feature tree with depth  $L$ , we can ultimately obtain  $N = \prod_{l=0}^L k_l$  distinct user profiles. It’s  
 222 important to note that each time the LLM generates  $k_l$ , this number can be much smaller than  $N$ ,  
 223 allowing the LLM to generate diverse features in one pass.

224 **User Behavior Generation.** Once user profiles are assigned, they include both explicit basic fea-  
 225 tures (e.g., occupation, gender, location) and implicit preferences (e.g., price sensitivity, product  
 226 affinity). However, in real-world scenarios, user preferences are typically inferred through behav-  
 227 ior patterns rather than explicitly stated. To simulate authentic behavioral traits, we employ an  
 228 LLM-based role-playing approach, where the model generates user actions on various platforms  
 229 based on their profile and platform characteristics. For instance, given a user’s preference for budget-  
 230 conscious shopping, the model may generate interactions such as “searches for hiking backpacks on  
 231 Amazon” or “purchases coffee from Walmart for \$30.” While implicit preferences remain unob-  
 232 servable to the model during task execution, they are embedded in prompts when generating tool  
 233 invocation solutions, ensuring accurate and contextually appropriate tool selection.

### 234 4.3 QUERY AND SOLUTION GENERATION

236 For generating query-solution pairs, we adopt a multi-agent collaborative approach, involving two  
 237 agents: the user agent and the assistant agent. The user agent generates queries by role-playing  
 238 based on the user profile, while the assistant agent generates tool invocation solutions. The user  
 239 agent’s role information includes both basic and implicit features, as these provide a more accurate  
 240 user representation than explicit behavioral features.

241 Given that a user’s platform preferences may vary across queries, we explicitly incorporate platform  
 242 information into the user agent’s prompt. This enables the agent to generate queries aligned with  
 243 the user’s platform preferences. Additionally, we instruct the user agent to avoid revealing profile  
 244 information in the queries, ensuring the generation of profile-dependent queries as well.

245 To ensure the correctness of tool invocations, we employ a two-tier verification strategy: rule-based  
 246 validation and model-based verification. Rule-based validation checks the format of tool invocations  
 247 to prevent issues such as unresolvable results or hallucinated tools and parameters. Model-based  
 248 verification inputs the user profile, query, and solution triples into the LLM to verify parameter  
 249 correctness, detect hallucinations, and assess whether the solution effectively resolves the query.  
 250 Furthermore, to ensure evaluation accuracy, we manually inspect tool invocation parameters. These  
 251 parameters are annotated as profile-related or query-related, indicating whether they originate from  
 252 the user profile or the query, facilitating more precise error feedback during evaluation.

## 254 5 EXPERIMENTS

### 256 5.1 EXPERIMENTAL SETTINGS

258 **Dataset Details.** We leverage GPT-4-turbo to synthesize the personalized tool invocation dataset  
 259 via our proposed framework. The overall dataset consists of a total of 80 users and 8,197 queries  
 260 under 5 scenarios, including shopping, takeout, entertainment, work, and travel. Under each sce-  
 261 nario, there are 3 platforms and 24 APIs in each platform as tools. We separate the dataset into  
 262 training and test sets, randomly selecting all queries of 6 users and about 6% queries of another 74  
 263 users to form the test set PTBench. The 6 users will not be visible to models in the training process,  
 264 termed as untrained. To ensure the quality of the test set, we manually verify each sample. Addi-  
 265 tionally, we construct a dataset comprising 116 samples from two unseen scenarios—finance and  
 266 lifestyle—which were not exposed to the models during training, to evaluate their generalization  
 267 capability. The statistics are illustrated in Appendix A.2.

268 **Evaluation.** We first evaluate the format accuracy by checking if the model’s output can give for-  
 269 matted output, verifying instruction following ability. The solution of each sample comprises two  
 parts: platform and tool invocation. The models are required to select the correct user-preferred plat-

270 Table 1: Comparison with baseline models on PTBench in terms of accuracy. **Bold** and underline  
 271 represent the best and the 2nd best results. **Tool-P** denotes the tool personalization.  $T^*$  denotes the  
 272 correctness \* in tool invocation. **DS-R1-Dis** is the abbreviation of DeepSeek-R1-Distill.

Type	Model	Format	Tool-P	Param Value		Tool Invocation			Overall		
				Query	Profile	T-name	T-param	T-value	Trained	Untrained	Overall
API	GPT-4-turbo	<b>97.78</b>	54.84	<b>81.23</b>	<u>68.32</u>	<u>91.78</u>	<u>77.09</u>	<b>35.18</b>	<u>18.34</u>	<u>18.56</u>	<b>18.47</b>
	GPT-4o	90.12	44.84	71.44	61.04	82.83	69.91	28.69	13.50	17.08	15.51
	Deepseek-v3	90.95	52.80	73.09	64.16	84.60	75.30	30.85	17.08	17.57	17.36
	Deepseek-r1	81.99	48.19	63.04	58.06	73.76	62.94	26.24	14.77	14.94	14.86
	Qwen-max	76.92	49.46	60.94	54.40	70.91	58.43	23.48	14.56	17.07	15.97
	Claude-3.5-sonnet	96.86	58.26	78.24	65.04	71.10	64.45	23.26	13.29	13.95	13.67
OSS	DS-R1-Llama-8B	64.27	30.19	38.23	30.12	50.80	38.02	9.81	4.85	3.94	4.34
	DS-R1-Qwen-7B	60.95	14.69	23.41	10.39	36.56	21.13	2.21	0.42	0.66	0.55
	Qwen2.5-7B-Inst	78.58	37.95	61.32	41.65	68.33	54.30	18.37	7.17	7.55	7.38
	Llama-3.1-8B-Inst	88.65	40.53	66.48	51.41	79.97	62.52	21.33	9.29	9.85	9.60
	Mistral-7B-v0.3	85.87	39.03	55.98	37.23	66.12	35.72	14.50	6.74	5.59	6.09
	Hammer2.1-7b	96.49	36.38	72.96	52.59	84.02	63.16	22.62	7.39	6.89	7.11
	ToolACE-8B	40.35	16.81	32.89	20.49	38.87	26.31	9.06	3.38	3.78	3.60
	Watt-tool-8B	37.49	22.81	27.16	19.90	34.08	22.18	8.26	5.91	4.11	4.89
	xLAM-7b-r	95.29	32.85	67.94	49.68	86.88	59.34	22.17	6.96	7.71	7.38
	Ours	95.75	<b>73.74</b>	<u>79.33</u>	<b>73.41</b>	<b>92.42</b>	<b>82.90</b>	<u>34.17</u>	<b>27.01</b>	<b>26.60</b>	<b>26.78</b>

291  
 292  
 293 form and then generate suitable tool invocations. Platform accuracy demonstrates the ability of tool  
 294 preference understanding. The tool invocation consists of three parts: tool name, parameters, and  
 295 parameter values, where the parameter values comprise query-related and profile-related parameters.  
 296 Profile-related parameters require the model to infer from the user profile, evaluating the ability to  
 297 handle profile-dependent query. We calculate the accuracy of the function name, function parameter,  
 298 and function value, respectively. The calculations of accuracy are detailed in Appendix A.1.

299 **Baselines.** We compare the latest open-source models and API-based models, as well as fine-tuned  
 300 tool-calling models. Open-source models include DeepSeek-R1-Distill-Llama-8B(DeepSeek-AI,  
 301 2025), DeepSeek-R1-Distill-Qwen-7B(DeepSeek-AI, 2025), Qwen2.5-7B-Instruct(Team, 2024a;b),  
 302 Llama-3.1-8B-Instruct (AI@Meta, 2024) and Mistral-7B-Instruct-v0.3(Jiang et al., 2023). API-  
 303 based models include GPT-4-turbo<sup>1</sup>, GPT-4o<sup>1</sup>, Deepseek-v3(DeepSeek-AI, 2024), Deepseek-  
 304 r1(DeepSeek-AI, 2025), Qwen-max(Team, 2024b) and Claude-3.5-sonnet<sup>2</sup>. Models fine-tuned for  
 305 tool-calling include Hammer2.1-7b(Lin et al., 2024), ToolACE-8B(Liu et al., 2025), watt-tool-8B<sup>3</sup>  
 306 and xLAM-7b-r(Zhang et al., 2024b; Liu et al., 2024; Zhang et al., 2024a).

307 **Implementation Details.** To validate the effectiveness of our model, we conducted various experiments  
 308 by training LLMs with the synthesized dataset. We train the open-source LLM, Qwen2.5-7B-  
 309 Instruct(Team, 2024a;b), in the supervised fine-tuning (SFT) manner. Due to limited resources, we  
 310 adopt the parameter-efficient LoRA(Hu et al., 2022) training strategy to fine-tune the model. As for  
 311 the hyper-parameters setting, we set the rank as 8, alpha as 16 learning rate as  $10^{-4}$ , LR scheduler  
 312 as cosine, WarmUp Ratio as 0.1 and epoch as 1 for all modules in the model.

## 313 5.2 MAIN RESULTS

314  
 315 The overall results are illustrated in Table 1. The detailed results of trained and untrained users are  
 316 presented in Appendix A.2. We have the following findings according to the results:

317 *Finding 1:* API-based large models significantly outperform smaller OSS models across various  
 318 dimensions, including format compliance, tool preference capabilities, and tool invocation abilities.  
 319 This aligns with the findings of most benchmarks, primarily attributed to the enhanced capabilities  
 320 enabled by the larger scale of model parameters.

321<sup>1</sup><https://chatgpt.com>

322<sup>2</sup><https://www.anthropic.com>

323<sup>3</sup><https://ollama.com>

324 Table 2: Ablation of user profile on PTBench. The models are trained with various variants. The  
 325 input in evaluation remains consistent with the training input.

Data	Untrained	Trained	Overall
All	<b>26.60</b>	<b>27.01</b>	<b>26.78</b>
All w/o Basic	9.69	24.26	16.06
All w/o History	24.63	25.31	24.93
All w/o Basic&History	5.91	7.81	6.74

326  
 327  
 328  
 329  
 330  
 331  
 332  
 333  
 334  
 335 *Finding 2:* Most models fall short on the tool-preference task, including the state-of-the-art model-  
 336 GPT-4-turbo, indicating the high complexity of selecting a suitable one from several similar tools  
 337 according to the user profile. Our model outperforms nearly all models in all aspects by a consider-  
 338 able improvement, presenting the necessity of personalized tool-invocation enhancement.

339 *Finding 3:* Our model demonstrates a significant improvement in its performance across various  
 340 tasks on PTBench. Notably, the enhancement in the Tool Preference task is particularly pronounced  
 341 when compared to the pre-trained Qwen2.5-7B-Instruct model. This also indicates that, even with-  
 342 out additional manual verification of the training data, the model achieves a high accuracy, demon-  
 343 strating the effectiveness of the proposed synthesis framework. Additionally, our model shows a  
 344 significant improvement on untrained users, presenting the generalization of the model.

345 *Finding 4:* All models exhibit lower accuracy on profile-dependent parameter values compared to  
 346 query-dependent parameters, indicating that inferring parameters from the profile presents a greater  
 347 challenge. While our trained model does not surpass GPT-4-turbo in accuracy on query-dependent  
 348 parameters, it outperforms larger models on profile-dependent parameters. Furthermore, the im-  
 349 provement over the pre-trained Qwen2.5-7B-Instruct model is more substantial, demonstrating the  
 350 effectiveness of our data generation framework in handling the query-dependent query tasks.

### 351 352 353 ABLATION STUDY

354 To investigate the importance of various parts in our synthesized user profile, we conduct the ablation  
 355 study on the user profile, including 4 variants on the user profile:

- 356 • **All.** All information in the user profile is used, including basic features and behavioral history.
- 357 • **All w/o Basic.** Basic features are omitted.
- 358 • **All w/o History.** The behavioral history is given.
- 359 • **All w/o Basic&History.** Both basic features and behavioral history are omitted.

360 First, We use the four dataset variants to train and then evaluate the model with the consistent input.  
 361 The results are reported in Table 2. From the result, we can observe that the existence of user history  
 362 and basic features hold contributions to the overall performance of the model to an extent.

363 Additionally, we conduct experiments under two settings: (1) train the model with the All variant  
 364 and evaluate the model with the four variants, illustrated in Figure 3a; (2) train the model with  
 365 the four variants and evaluate the models with the All variant, illustrated in Figure 3b. The results  
 366 exhibit that the model shows poor performance in the tool preference task when lacking user history  
 367 information in training or evaluation. On the other hand, the accuracy of tool invocation suffers  
 368 when basic features are absent, led by the challenging profile-dependent query task.

369 To further confirm that the curated instructions can only be completed with personalized informa-  
 370 tion, we conducted an additional experiment where all personalized information was removed from  
 371 the instructions. As shown in Table 3, model performance decreases in all settings compared to main  
 372 results, with the most pronounced decline observed in the precision of tool values. These results con-  
 373 firm that personalized information is crucial and indispensable for achieving optimal performance.

378  
379  
380 Table 3: Evaluating Model Performance Without Personalized Information.  
381  
382  
383  
384  
385  
386

Model	Format	Platform	T-name	T-param	T-value	Overall
GPT-4o	92.80	48.29	87.26	64.54	8.59	5.35
Deepseek-v3	98.34	51.25	91.69	77.28	10.16	5.72
Qwen2.5-7B-Instruct	90.58	44.32	82.64	64.64	8.96	4.16
Llama-3.1-8B-Instruct	95.29	43.31	87.35	69.07	9.23	3.97
Ours	95.57	52.34	96.51	87.55	6.18	5.91

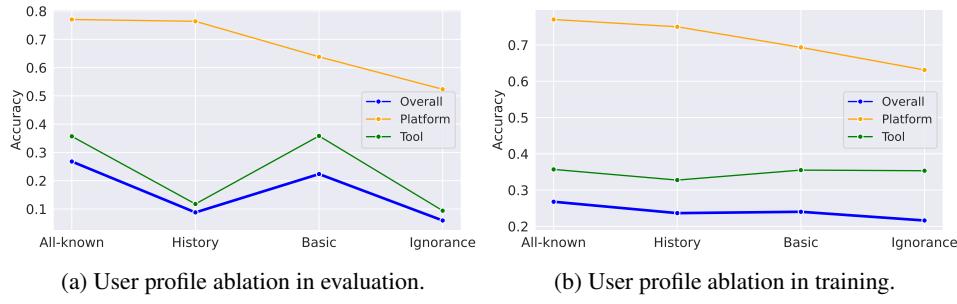
387  
388  
389  
390  
391  
392  
393  
394  
395  
396 (a) User profile ablation in evaluation.  
397 (b) User profile ablation in training.  
398  
399  
400  
401

Figure 3: Ablation study on user profile in evaluation and training, respectively.

402  
403  
404  
405  
406  
407  
408  
409  
410  
5.4 ERROR ANALYSIS

411 To gain deeper insights into the types of errors made by the models during the evaluation, we conduct  
412 investigations into the error types on our model, GPT-4-turbo, and Qwen2.5-7B-Instruct. We only  
413 analyze solutions with the correct format.

414 We analyze the function errors generally and divide them into 6 categories: wrong tools, missing  
415 tools, excessive tools, missing parameters, excessive parameters, and wrong parameters. The results  
416 are shown in Figure 4. From the pie chart, it is evident that filling the correct parameters is more  
417 challenging than the selection of the correct tools. After training with our synthesized data, the  
418 model is more familiar with the candidate tools, demonstrating less error percentage in tool selection.

419  
420  
421  
422  
423  
424  
425  
426  
427  
5.5 FURTHER ANALYSIS

428  
429  
430  
431  
432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
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  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
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  
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  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
7010  
7011  
7012  
7013  
7014  
7015  
7016  
7017  
7018  
7019  
7020  
7021  
7022  
7023  
7024  
7025  
7026  
7027  
7028  
7029  
7030  
7031  
7032  
7033  
7034  
7035  
7036  
7037  
7038  
7039  
7040  
7041  
7042  
7043  
7044  
7045  
7046  
7047  
7048  
7049  
7050  
7051  
7052  
7053  
7054  
7055  
7056  
7057  
7058  
7059  
7060  
7061  
7062  
7063  
7064  
7065  
7066  
7067  
7068  
7069  
7070  
7071  
7072  
7073  
7074  
7075  
7076  
7077  
7078  
7079  
7080  
7081  
7082  
7083  
7084  
7085  
7086  
7087  
7088  
7089  
7090  
7091  
7092  
7093  
7094  
7095  
7096  
7097  
7098  
7099  
70100  
70101  
70102  
70103  
70104  
70105  
70106  
70107  
70108  
70109  
70110  
70111  
70112  
70113  
70114  
70115  
70116  
70117  
70118  
70119  
70120  
70121  
70122  
70123  
70124  
70125  
70126  
70127  
70128  
70129  
70130  
70131  
70132  
70133  
70134  
70135  
70136  
70137  
70138  
70139  
70140  
70141  
70142  
70143  
70144  
70145  
70146  
70147  
70148  
70149  
70150  
70151  
70152  
70153  
70154  
70155  
70156  
70157  
70158  
70159  
70160  
70161  
70162  
70163  
70164  
70165  
70166  
70167  
70168  
70169  
70170  
70171  
70172  
70173  
70174  
70175  
70176  
70177  
70178  
70179  
70180  
70181  
70182  
70183  
70184  
70185  
70186  
70187  
70188  
70189  
70190  
70191  
70192  
70193  
70194  
70195  
70196  
70197  
70198  
70199  
70200  
70201  
70202  
70203  
70204  
70205  
70206  
70207  
70208  
70209  
70210  
70211  
70212  
70213  
70214  
70215  
70216  
70217  
70218  
70219  
70220  
70221  
70222  
70223  
70224  
70225  
70226  
70227  
70228  
70229  
70230  
70231  
70232  
70233  
70234  
70235  
70236  
70237  
70238  
70239  
70240  
70241  
70242  
70243  
70244  
70245  
70246  
70247  
70248  
70249  
70250  
70251  
70252  
70253  
70254  
70255  
70256  
70257  
70258  
70259  
70260  
70261  
70262  
70263  
70264  
70265  
70266  
70267  
70268  
70269  
70270  
70271  
70272  
70273  
70274  
70275  
70276  
70277  
70278  
70279  
70280  
70281  
70282  
70283  
70284  
70285  
70286  
70287  
70288  
70289  
70290  
70291  
70292  
70293  
70294  
70295  
70296  
70297  
70298  
70299  
70300  
70301  
70302  
70303  
70304  
70305  
70306  
70307  
70308  
70309  
70310  
70311  
70312  
70313  
70314  
70315  
70316  
70317  
70318  
70319  
70320  
70321  
70322  
70323  
70324  
70325  
70326  
70327  
70328  
70329  
70330  
70331  
70332  
70333  
70334  
70335  
70336  
70337  
70338  
70339  
70340  
70341  
70342  
70343  
70344  
70345  
70346  
70347  
70348  
70349  
70350  
70351  
70352  
70353  
70354  
70355  
70356  
70357  
70358  
70359  
70360  
70361  
70362  
70363  
70364  
70365  
70366  
70367  
70368  
70369  
70370  
70371  
70372  
70373  
70374  
70375  
70376  
70377  
70378  
70379  
70380  
70381  
70382  
70383  
70384  
70385  
70386  
70387  
70388  
70389  
70390  
70391  
70392  
70393  
70394  
70395  
70396  
70397  
70398  
70399  
70400  
70401  
70402  
70403  
70404  
70405  
70406  
70407  
70408  
70409  
70410  
70411  
70412  
70413  
70414  
70415  
70416  
70417  
70418  
70419  
70420  
70421  
70422  
70423  
70424  
70425  
70426  
70427  
70428  
70429  
70430  
70431  
70432  
70433  
70434  
70435  
70436  
70437  
70438  
70439  
70440  
70441  
70442  
70443  
70444  
70445  
70446  
70447  
70448  
70449  
70450  
70451  
70452  
70453  
70454  
70455  
70456  
70457  
70458  
70459  
70460  
70461  
70462  
70463  
70464  
70465  
70466  
70467  
70468  
70469  
70470  
70471  
70472  
70473  
70474  
70475  
70476  
70477  
70478  
70479  
70480  
70481  
70482  
70483  
70484  
70485  
70486  
70487  
70488  
70489  
70490  
70491  
70492  
70493  
70494  
70495  
70496  
70497  
70498  
70499  
70500  
70501  
70502  
70503  
70504  
70505  
70506  
70507  
70508  
70509  
70510  
70511  
70512  
70513  
70514  
70515  
70516  
70517  
70518  
70519  
70520  
70521  
70522  
70523  
70524  
70525  
70526  
70527  
70528  
70529  
70530  
70531  
70532  
70533  
70534  
70535  
70536  
70537  
70538  
70539  
70540  
70541  
70542  
70543  
70544  
70545  
70546  
70547  
70548  
70549  
70550  
70551  
70552  
70553  
70554  
70555  
70556  
70557  
70558  
70559  
70560  
70561  
70562  
70563  
70564  
70565  
70566  
70567  
70568  
70569  
70570  
70571  
70572  
70573  
70574  
70575  
70576  
70577  
70578  
70579  
70580  
70581  
70582  
70583  
70584  
70585  
70586  
70587  
70588  
70589  
70590  
70591  
70592  
70593  
70594  
70595  
70596  
70597  
70598  
70599  
70600  
70601  
70602  
70603  
70604  
70605  
70606  
70607  
70608  
70609  
70610  
70611  
70612  
70613  
70614  
70615  
70616  
70617  
70618  
70619  
70620  
70621  
70622  
70623  
70624  
70625  
70626  
70627  
70628  
70629  
70630  
70631  
70632  
70633  
70634  
70635  
70636  
70637  
70638  
70639  
70640  
70641  
70642  
70643  
70644  
70645  
70646  
70647  
70648  
70649  
70650  
70651  
70652  
70653  
70654  
70655  
70656  
70657  
70658  
70659  
70660  
70661  
70662  
70663  
70664  
70665  
70666  
70667  
70668  
70669  
70670  
70671  
70672  
70673  
70674  
70675  
70676  
70677  
70678  
70679  
70680  
70681  
70682  
70683  
70684  
70685  
70686  
70687  
70688  
70689  
70690  
70691  
70692  
70693  
70694  
70695  
70696  
70697  
70698  
70699  
70700  
70701  
70702  
70703  
70704  
70705  
70706  
70707  
70708  
70709  
70710  
70711  
70712  
70713  
70714  
70715  
70716  
70717  
70718  
70719  
70720  
70721  
70722  
70723  
70724  
70725  
70726  
70727  
70728  
70729  
70730  
70731  
70732  
70733  
70734  
70735  
70736  
70737  
70738  
70739  
70740  
70741  
70742  
70743  
70744  
70745  
70746  
70747  
70748  
70749  
70750  
70751  
70752  
70753  
70754  
70755  
70756  
70757  
70758  
70759  
70760  
70761  
70762  
70763  
70764  
70765  
70766  
70767  
70768  
70769  
70770  
70771  
70772  
70773  
70774  
70775  
70776  
70777  
70778  
70779  
70780  
70781  
70782  
70783  
70784  
70785  
70786  
70787  
70788  
70789  
70790  
70791  
70792  
70793  
70794  
70795  
70796  
70797  
70798  
70799  
70800  
70801  
70802  
70803  
70804  
70805  
70806  
70807  
70808  
70809  
70810  
70811  
70812  
70813  
70814  
70815  
70816  
70817  
70818  
70819  
70820  
70821  
70822  
70823  
70824  
70825  
70826  
70827  
70828  
70829  
70830  
70831  
70832  
70833  
70834  
70835  
70836  
70837  
70838  
70839  
70840  
70841  
70842  
70843  
70844  
70845  
70846  
70847  
70848  
70849  
70850  
70851  
70852  
70853  
70854  
70855  
70856  
70857  
70858  
70859  
70860  
70861  
70862  
70863  
70864  
70865  
70866  
70867  
70868  
70869  
70870  
70871  
70872  
70873  
70874  
70875  
70876  
70877  
70878  
70879  
70880  
70881  
70882  
70883  
70884  
70885  
70886  
70887  
70888  
70889  
70890  
70891  
70892  
70893  
70894  
70895  
70896  
70897  
70898  
70899  
70900  
70901  
70902  
70903  
70904  
70905  
70906  
70907  
70908  
70909  
70910  
70911  
70912  
70913  
70914  
70915  
70916  
70917  
70918  
70919  
70920  
70921  
70922  
70923  
70924  
70925  
70926  
70927  
70928  
70929  
70930  
70931  
70932  
70933  
70934  
70935  
70936  
70937  
70938  
70939  
70940  
70941  
70942  
70943  
70944  
70945  
70946  
70947  
70948  
70949  
70950  
70951  
70952  
70953  
70954  
70955  
70956  
70957  
70958  
70959  
70960  
70961  
70962  
70963  
70964  
70965  
70966  
70967  
70968  
70969  
70970  
70971  
70972  
70973  
70974  
70975  
70976  
70977  
70978  
70979  
70980  
70981  
70982  
70983  
70984  
70985  
70986  
70987  
70988  
70989  
70990  
70991  
70992  
70993  
70994  
70995  
70996  
70997  
70998  
70999  
70100  
70101  
70102  
70103  
70104  
70105  
70106  
70107  
70108  
70109  
70110  
70111  
70112  
70113  
70114  
70115  
70116  
70117  
70118  
70119  
70120  
70121  
70122  
70123  
70124  
70125  
70126  
70127  
70128  
70129  
70130  
70131  
70132  
70133  
70134  
70135  
70136  
70137  
70138  
70139  
70140  
70141  
70142  
70143  
70144  
70145  
70146  
70147  
70148  
70149  
70150  
70151  
70152  
70153  
70154  
70155  
70156  
70157  
70158  
70159  
70160  
70161  
70162  
70163  
70164  
70165  
70166  
70167  
70168  
70169  
70170  
70171  
70172  
70173  
70174  
70175  
70176  
70177  
70178  
70179  
70180  
70181  
70182  
70183  
70184  
70185  
70186  
70187  
70188  
70189  
70190  
70191  
70192  
70193  
70194  
70195  
70196  
70197  
70198  
70199  
70200  
70201  
70202  
70203  
70204  
70205  
70206  
70207  
70208  
70209  
70210  
70211  
70212  
70213  
70214  
70215  
70216  
70217  
70218  
70219  
70220  
70221  
70222  
70223  
70224  
70225  
70226  
70227  
70228  
70229  
70230  
70231  
70232  
70233  
70234  
70235  
70236  
70237  
70238  
70239  
70240  
70241  
70242  
70243  
70244  
70245  
70246  
70247  
70248  
70249  
70250  
70251  
70252  
70253  
70254  
70255  
70256  
70257  
70258  
70259  
70260  
70261  
70262  
70263  
70264  
70265  
70266  
70267  
70268  
70269  
70270  
70271  
70272  
70273  
70274  
70275  
70276  
70277  
70278  
70279  
70280  
70281  
70282  
70283  
70284  
70285  
70286  
70287  
70288  
70289  
70290  
70291  
70292  
70293  
70294  
70295  
70296  
70297  
70298  
70299  
702100  
702101  
702102  
702103  
702104  
702105  
70

Table 4: Models performance results on new scenarios. Bold represents the best result.

Model	Format	Platform	T-name	T-param	T-value	Overall
GPT-4-turbo	91.80	38.37	76.32	55.30	16.92	5.48
Deepseek-v3	94.10	47.05	86.80	77.19	20.10	7.77
Qwen2.5-7B-Instruct	83.48	28.44	60.09	24.31	4.13	1.84
Llama-3.1-8B-Instruct	94.50	30.73	71.10	52.30	9.63	2.29
Hammer2.1-7b	94.04	33.03	75.23	37.16	9.63	4.59
xLAM-7b-r	98.62	26.61	83.49	41.28	10.55	3.67
<b>Ours</b>	<b>100.00</b>	<b>73.39</b>	<b>88.99</b>	<b>75.69</b>	<b>26.15</b>	<b>18.35</b>

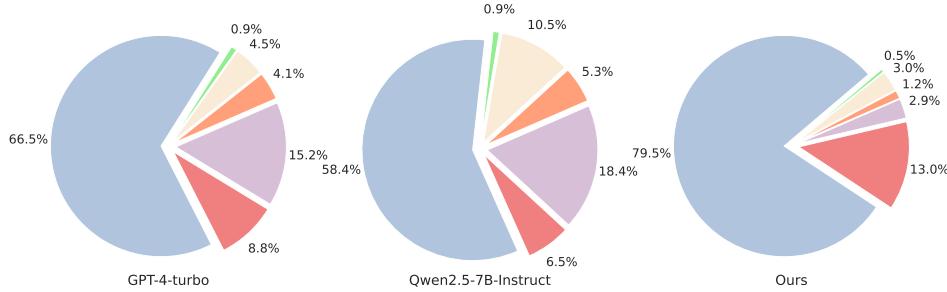


Figure 4: Error Analysis on PTBench. T-wrong, T-missing, and T-excessive represent wrong tools, missing tools and excessive tools, respectively. P-missing, P-excessive and P-error represent missing parameters, excessive parameters and wrong parameters, respectively.

**Generalization to Unseen Scenarios.** To further examine the generalizability of our model beyond the five common scenarios, we conduct additional evaluation on unseen domains. Specifically, we synthesized 218 samples covering three new scenarios: finance(42 samples), lifestyle(74 samples), and knowledge(102 samples). As shown in Table 4, our model consistently outperforms the baselines under all domains, demonstrating strong robustness and adaptability. These results provide evidence that the synthesized framework can be extended to a broader range of real-world scenarios. Further derailed results are shown in Appendix A.2.

**General Capabilities.** In order to validate that our synthesized data does not introduce negative effects on the model’s general capabilities, we employ a diverse set of benchmarks to assess the performance from different perspectives, Including general ability (MMLU (Hendrycks et al., 2021a;b)), coding (HumanEval (Chen et al., 2021)), math (GSM8K (Cobbe et al., 2021)), reasoning (CommonSenseQA (Talmor et al., 2019)), abstract reasoning (ARC (Chollet, 2019)), and basic function-calling (tool-invocation) ability (BFCL non-live (Yan et al., 2024)). xLAM-7B-r, LLaMA-3-8B-Instruct, Raw Qwen2.5-7B-Instruct serve as baselines. The results are shown in Figure 6. From the figure, it is evident that there is no significance deterioration on abilities of our model compared to the raw model Qwen2.5-7B-Instruct. Nonetheless, our model gains a notable improvement on BFCL non-live, These findings suggest that our approach effectively enhances personalized functional calling capabilities without compromising the underlying LLM’s other abilities.

## 6 CONCLUSION

In this work, we introduce the concept of personalized tool invocation, which encompasses two primary tasks: tool preference and profile-dependent queries. These tasks require the model’s ability to understand the user’s profile, select preferred tools based on historical behavior, and extract tool parameters from user information. To enhance and evaluate the model’s personalized tool invocation capabilities, we propose a data synthesis framework and create a benchmark, PTBench, by manually inspecting a subset of the generated data. Extensive experimental evaluations assess the personalized tool invocation abilities of existing models, confirming the effectiveness of our synthesized data and its harmlessness to other model capabilities.

486 REFERENCES  
487

488 Mahyar Abbasian, Zhongqi Yang, Elahe Khatibi, Pengfei Zhang, Nitish Nagesh, Iman Azimi,  
489 Ramesh Jain, and Amir M Rahmani. Knowledge-infused llm-powered conversational health  
490 agent: A case study for diabetes patients. *arXiv preprint arXiv:2402.10153*, 2024.

491 AI@Meta. Llama 3 model card. 2024. URL [https://github.com/meta-llama/llama3/blob/main/MODEL\\_CARD.md](https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md).

492

493 Keqin Bao, Jizhi Zhang, Yang Zhang, Wenjie Wang, Fuli Feng, and Xiangnan He. Tallrec: An  
494 effective and efficient tuning framework to align large language model with recommendation. In  
495 *Proceedings of the 17th ACM Conference on Recommender Systems*, pp. 1007–1014, 2023.

496

497 Anastasiya Belyaeva, Justin Cosentino, Farhad Hormozdiari, Krish Eswaran, Shravya Shetty, Greg  
498 Corrado, Andrew Carroll, Cory Y McLean, and Nicholas A Furlotte. Multimodal llms for health  
499 grounded in individual-specific data. In *Workshop on Machine Learning for Multimodal Health-  
500 care Data*, pp. 86–102. Springer, 2023.

501

502 Chen Chen, Xinlong Hao, Weiwen Liu, Xu Huang, Xingshan Zeng, Shuai Yu, Dexun Li, Shuai  
503 Wang, Weinan Gan, Yuefeng Huang, et al. Acebench: Who wins the match point in tool learning?  
504 *arXiv preprint arXiv:2501.12851*, 2025.

505

506 Jin Chen, Zheng Liu, Xu Huang, Chenwang Wu, Qi Liu, Gangwei Jiang, Yuanhao Pu, Yuxuan  
507 Lei, Xiaolong Chen, Xingmei Wang, et al. When large language models meet personalization:  
508 Perspectives of challenges and opportunities. *World Wide Web*, 27(4):42, 2024.

509

510 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, and et al.  
511 Evaluating large language models trained on code. 2021.

512

513 Zihao Cheng, Hongru Wang, Zeming Liu, Yuhang Guo, Yuanfang Guo, Yunhong Wang, and  
514 Haifeng Wang. ToolSpectrum: Towards personalized tool utilization for large language models.  
515 In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.),  
516 *Findings of the Association for Computational Linguistics: ACL 2025*, pp. 20679–20699, Vi-  
517 enna, Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-89176-256-  
518 5. doi: 10.18653/v1/2025.findings-acl.1063. URL <https://aclanthology.org/2025.findings-acl.1063/>.

519

520 François Chollet. On the measure of intelligence, 2019. URL <https://arxiv.org/abs/1911.01547>.

521

522 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
523 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John  
524 Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*,  
525 2021.

526

527 Yuhao Dan, Zhikai Lei, Yiyang Gu, Yong Li, Jianghao Yin, Jiaju Lin, Linhao Ye, Zhiyan Tie,  
528 Yougen Zhou, Yilei Wang, et al. Educhat: A large-scale language model-based chatbot system  
529 for intelligent education. *arXiv preprint arXiv:2308.02773*, 2023.

530

531 DeepSeek-AI. Deepseek-v3 technical report, 2024. URL <https://arxiv.org/abs/2412.19437>.

532

533 DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning,  
534 2025. URL <https://arxiv.org/abs/2501.12948>.

535

536 Shibo Hao, Tianyang Liu, Zhen Wang, and Zhitong Hu. Toolkengpt: Augmenting frozen language  
537 models with massive tools via tool embeddings. *Advances in neural information processing sys-  
538 tems*, 36, 2024.

539

540 Yupu Hao, Pengfei Cao, Zhuoran Jin, Huanxuan Liao, Yubo Chen, Kang Liu, and Jun Zhao. Eval-  
541 uating personalized tool-augmented LLMs from the perspectives of personalization and proac-  
542 tivity. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar  
543 (eds.), *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics*

540 (Volume 1: Long Papers), pp. 21897–21935, Vienna, Austria, July 2025. Association for Com-  
 541 putational Linguistics. ISBN 979-8-89176-251-0. doi: 10.18653/v1/2025.acl-long.1064. URL  
 542 <https://aclanthology.org/2025.acl-long.1064/>.

543

544 Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob  
 545 Steinhardt. Aligning ai with shared human values. *Proceedings of the International Conference*  
 546 *on Learning Representations (ICLR)*, 2021a.

547 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob  
 548 Steinhardt. Measuring massive multitask language understanding. *Proceedings of the Interna-*  
 549 *tional Conference on Learning Representations (ICLR)*, 2021b.

550

551 Yupeng Hou, Junjie Zhang, Zihan Lin, Hongyu Lu, Ruobing Xie, Julian McAuley, and Wayne Xin  
 552 Zhao. Large language models are zero-shot rankers for recommender systems. In *European*  
 553 *Conference on Information Retrieval*, pp. 364–381. Springer, 2024.

554

555 Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,  
 556 and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Con-*  
 557 *ference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=nZeVKeFYf9>.

558

559 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chap-  
 560 lot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier,  
 561 Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril,  
 562 Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023. URL <https://arxiv.org/abs/2310.06825>.

563

564 Mingyu Jin, Qinkai Yu, Dong Shu, Chong Zhang, Lizhou Fan, Wenyue Hua, Suiyuan Zhu, Yanda  
 565 Meng, Zhenting Wang, Mengnan Du, et al. Health-llm: Personalized retrieval-augmented disease  
 566 prediction system. *arXiv preprint arXiv:2402.00746*, 2024.

567

568 Enkelejda Kasneci, Kathrin Seßler, Stefan Küchemann, Maria Bannert, Daryna Dementieva, Frank  
 569 Fischer, Urs Gasser, Georg Groh, Stephan Günnemann, Eyke Hüllermeier, et al. Chatgpt for  
 570 good? on opportunities and challenges of large language models for education. *Learning and*  
*individual differences*, 103:102274, 2023.

571

572 Kausik Lakkaraju, Sai Krishna Revanth Vuruma, Vishal Pallagani, Bharath Muppasani, and Biplav  
 573 Srivastava. Can llms be good financial advisors?: An initial study in personal decision making  
 574 for optimized outcomes. *arXiv preprint arXiv:2307.07422*, 2023.

575

576 Yuxuan Lei, Jianxun Lian, Jing Yao, Xu Huang, Defu Lian, and Xing Xie. Recexplainer: Aligning  
 577 large language models for explaining recommendation models. In *Proceedings of the 30th ACM*  
*SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 1530–1541, 2024.

578

579 Ruyu Li, Wenhao Deng, Yu Cheng, Zheng Yuan, Jiaqi Zhang, and Fajie Yuan. Exploring the upper  
 580 limits of text-based collaborative filtering using large language models: Discoveries and insights.  
*arXiv preprint arXiv:2305.11700*, 2023.

581

582 Qiqiang Lin, Muning Wen, Qiuying Peng, Guanyu Nie, Junwei Liao, Jun Wang, Xiaoyun Mo, Jiamu  
 583 Zhou, Cheng Cheng, Yin Zhao, Jun Wang, and Weinan Zhang. Hammer: Robust function-calling  
 584 for on-device language models via function masking, 2024. URL <https://arxiv.org/abs/2410.04587>.

585

586 Weiwen Liu, Xingshan Zeng, Xu Huang, xinlong hao, Shuai Yu, Dexun Li, Shuai Wang, Weinan  
 587 Gan, Zhengying Liu, Yuanqing Yu, Zehong WANG, Yuxian Wang, Wu Ning, Yutai Hou, Bin  
 588 Wang, Chuhan Wu, Wang Xinzhi, Yong Liu, Yasheng Wang, Duyu Tang, Dandan Tu, Lifeng  
 589 Shang, Xin Jiang, Ruiming Tang, Defu Lian, Qun Liu, and Enhong Chen. ToolACE: Enhancing  
 590 function calling with accuracy, complexity, and diversity. In *The Thirteenth International Confer-*  
 591 *ence on Learning Representations*, 2025. URL <https://openreview.net/forum?id=8EB8k6DdCU>.

592

593 Xiao-Yang Liu, Guoxuan Wang, Hongyang Yang, and Daochen Zha. Fingpt: Democratizing  
 594 internet-scale data for financial large language models. *arXiv preprint arXiv:2307.10485*, 2023.

594 Z Liu, Z Lai, Z Gao, E Cui, Z Li, X Zhu, L Lu, Q Chen, Y Qiao, J Dai, et al. Controlllm: augment  
 595 language models with tools by searching on graphs (2023). *arXiv preprint arXiv:2310.17796*.  
 596

597 Zuxin Liu, Thai Hoang, Jianguo Zhang, Ming Zhu, Tian Lan, Shirley Kokane, Juntao Tan, Weiran  
 598 Yao, Zhiwei Liu, Yihao Feng, et al. Apigen: Automated pipeline for generating verifiable and  
 599 diverse function-calling datasets. *arXiv preprint arXiv:2406.18518*, 2024.

600 Hanjia Lyu, Song Jiang, Hanqing Zeng, Yinglong Xia, Qifan Wang, Si Zhang, Ren Chen, Chris  
 601 Leung, Jiajie Tang, and Jiebo Luo. LLM-rec: Personalized recommendation via prompting large  
 602 language models. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), *Findings of the*  
 603 *Association for Computational Linguistics: NAACL 2024*, pp. 583–612, Mexico City, Mexico,  
 604 June 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-naacl.39.  
 605 URL <https://aclanthology.org/2024.findings-naacl.39/>.

606 Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christo-  
 607 pher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted  
 608 question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.

609

610 Ha-Thanh Nguyen. A brief report on lawgpt 1.0: A virtual legal assistant based on gpt-3. *arXiv*  
 611 *preprint arXiv:2302.05729*, 2023.

612 Minju Park, Sojung Kim, Seunghyun Lee, Soonwoo Kwon, and Kyuseok Kim. Empowering person-  
 613 alized learning through a conversation-based tutoring system with student modeling. In *Extended*  
 614 *Abstracts of the CHI Conference on Human Factors in Computing Systems*, pp. 1–10, 2024.

615

616 Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru  
 617 Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein,  
 618 dahai li, Zhiyuan Liu, and Maosong Sun. ToolLLM: Facilitating large language models to master  
 619 16000+ real-world APIs. In *The Twelfth International Conference on Learning Representations*,  
 620 2024. URL <https://openreview.net/forum?id=dHng200Jjr>.

621

622 Changle Qu, Sunhao Dai, Xiaochi Wei, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, Jun Xu, and Ji-  
 623 Rong Wen. Tool learning with large language models: A survey. *Frontiers of Computer Science*,  
 624 19(8):198343, 2025.

625

626 Timo Schick, Jane Dwivedi-Yu, Roberto Dessim, Roberta Raileanu, Maria Lomeli, Eric Hambro,  
 627 Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can  
 628 teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36, 2024.

629

630 Zhengliang Shi, Shen Gao, Xiuyi Chen, Yue Feng, Lingyong Yan, Haibo Shi, Dawei Yin, Pengjie  
 631 Ren, Suzan Verberne, and Zhaochun Ren. Learning to use tools via cooperative and interac-  
 632 tive agents. pp. 10642–10657, Miami, Florida, USA, November 2024. doi: 10.18653/v1/2024.  
 633 findings-emnlp.624. URL 2024.findings-emnlp.624/.

634

635 Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. CommonsenseQA: A ques-  
 636 tion answering challenge targeting commonsense knowledge. In Jill Burstein, Christy Doran, and  
 637 Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of*  
 638 *the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long*  
 639 *and Short Papers)*, pp. 4149–4158, Minneapolis, Minnesota, June 2019. Association for Com-  
 640 *Computational Linguistics*. doi: 10.18653/v1/N19-1421. URL <https://aclanthology.org/N19-1421/>.

641

642 Qwen Team. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024a.

643

644 Qwen Team. Qwen2.5: A party of foundation models, September 2024b. URL <https://qwenlm.github.io/blog/qwen2.5/>.

645

646 Lei Wang, Songheng Zhang, Yun Wang, Ee-Peng Lim, and Yong Wang. LLM4Vis: Explainable  
 647 visualization recommendation using ChatGPT. In Mingxuan Wang and Imed Zitouni (eds.),  
 648 *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing:*  
 649 *Industry Track*, pp. 675–692, Singapore, December 2023. Association for Computational Lin-  
 650 *guistics*. doi: 10.18653/v1/2023.emnlp-industry.64. URL <https://aclanthology.org/2023.emnlp-industry.64/>.

648 Wei Wei, Xubin Ren, Jabin Tang, Qinyong Wang, Lixin Su, Suqi Cheng, Junfeng Wang, Dawei Yin,  
 649 and Chao Huang. Llmrec: Large language models with graph augmentation for recommendation.  
 650 In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*, pp.  
 651 806–815, 2024.

652 Likang Wu, Zhi Zheng, Zhaopeng Qiu, Hao Wang, Hongchao Gu, Tingjia Shen, Chuan Qin, Chen  
 653 Zhu, Hengshu Zhu, Qi Liu, Hui Xiong, and Enhong Chen. A survey on large language models  
 654 for recommendation. *CoRR*, abs/2305.19860, 2023.

655 Yunjia Xi, Weiwen Liu, Jianghao Lin, Xiaoling Cai, Hong Zhu, Jieming Zhu, Bo Chen, Ruim-  
 656 ing Tang, Weinan Zhang, and Yong Yu. Towards open-world recommendation with knowledge  
 657 augmentation from large language models. In *Proceedings of the 18th ACM Conference on Rec-  
 658 ommender Systems*, pp. 12–22, 2024.

659 Qiancheng Xu, Yongqi Li, Heming Xia, Fan Liu, Min Yang, and Wenjie Li. PEToolLLM: Towards  
 660 personalized tool learning in large language models. In Wanxiang Che, Joyce Nabende, Ekaterina  
 661 Shutova, and Mohammad Taher Pilehvar (eds.), *Findings of the Association for Computational  
 662 Linguistics: ACL 2025*, pp. 21488–21503, Vienna, Austria, July 2025. Association for Compu-  
 663 tational Linguistics. ISBN 979-8-89176-256-5. doi: 10.18653/v1/2025.findings-acl.1107. URL  
 664 <https://aclanthology.org/2025.findings-acl.1107/>.

665 Yang Xu, Yunlong Feng, Honglin Mu, Yutai Hou, Yitong Li, Xinghao Wang, Wanjun Zhong,  
 666 Zhongyang Li, Dandan Tu, Qingfu Zhu, Min Zhang, and Wanxiang Che. Concise and precise  
 667 context compression for tool-using language models. pp. 16430–16441, Bangkok, Thailand, Au-  
 668 gust 2024. doi: 10.18653/v1/2024.findings-acl.974. URL 2024.findings-acl.974/.

669 Fanjia Yan, Huanzhi Mao, Charlie Cheng-Jie Ji, Tianjun Zhang, Shishir G. Patil, Ion Stoica,  
 670 and Joseph E. Gonzalez. Berkeley function calling leaderboard. [https://gorilla.cs.berkeley.edu/blogs/8\\_berkeley\\_function\\_calling\\_leaderboard.html](https://gorilla.cs.berkeley.edu/blogs/8_berkeley_function_calling_leaderboard.html),  
 671 2024.

672 Fan Yang, Zheng Chen, Ziyan Jiang, Eunah Cho, Xiaojiang Huang, and Yanbin Lu. Palr: Personal-  
 673 ization aware llms for recommendation. *arXiv preprint arXiv:2305.07622*, 2023a.

674 Rui Yang, Lin Song, Yanwei Li, Sijie Zhao, Yixiao Ge, Xiu Li, and Ying Shan. GPT4tools: Teaching  
 675 large language model to use tools via self-instruction. In *Thirty-seventh Conference on Neural  
 676 Information Processing Systems*, 2023b. URL <https://openreview.net/forum?id=cwjh81qmOL>.

677 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.  
 678 React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*,  
 679 2022.

680 Siyu Yuan, Kaitao Song, Jiangjie Chen, Xu Tan, Yongliang Shen, Ren Kan, Dongsheng Li, and  
 681 Deqing Yang. Easytool: Enhancing llm-based agents with concise tool instruction. *arXiv preprint  
 682 arXiv:2401.06201*, 2024.

683 Jianguo Zhang, Tian Lan, Rithesh Murthy, Zhiwei Liu, Weiran Yao, Juntao Tan, Thai Hoang, Liang-  
 684 wei Yang, Yihao Feng, Zuxin Liu, et al. Agentohana: Design unified data and training pipeline  
 685 for effective agent learning. *arXiv preprint arXiv:2402.15506*, 2024a.

686 Jianguo Zhang, Tian Lan, Ming Zhu, Zuxin Liu, Thai Hoang, Shirley Kokane, Weiran Yao, Juntao  
 687 Tan, Akshara Prabhakar, Haolin Chen, et al. xlam: A family of large action models to empower  
 688 ai agent systems. *arXiv preprint arXiv:2409.03215*, 2024b.

689 Junjie Zhang, Ruobing Xie, Yupeng Hou, Xin Zhao, Leyu Lin, and Ji-Rong Wen. Recommendation  
 690 as instruction following: A large language model empowered recommendation approach. *ACM  
 691 Transactions on Information Systems*.

692 Zhehao Zhang, Ryan A Rossi, Branislav Kveton, Yijia Shao, Diyi Yang, Hamed Zamani, Franck  
 693 Dernoncourt, Joe Barrow, Tong Yu, Sungchul Kim, et al. Personalization of large language mod-  
 694 els: A survey. *arXiv preprint arXiv:2411.00027*, 2024c.

702 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min,  
 703 Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv*  
 704 *preprint arXiv:2303.18223*, 2023.

## 706 A EXPERIMENTS

### 709 A.1 EVALUATION METRICS

710 We categorize the PTBench metrics according to the key competencies required in personalized  
 711 tool-use scenarios as follows:

- 713 • **General Tool Use Capabilities.** This includes assessing format adherence through *Format Accu-*  
 714 *racy*, selecting correct tool through *Tool Name Accuracy*, specifying valid parameter names  
 715 through *Tool Parameter Accuracy*, and filling correct values through *Tool Parameter-Value Accu-*  
 716 *racy*.
- 717 • **Personalized Tool Use Capabilities.** This includes (i) tool preference personalization, measured  
 718 by *Platform Accuracy*, and (ii) profile-dependent query personalization, measured by the *Profile-*  
 719 *related Parameter-Value Accuracy*.

720 The calculation of various metrics in PTBench are formulated as follows:

- 722 • **Format Accuracy** indicates the instruction-following ability.

$$724 \text{format\_acc} = \frac{\#\text{parsable samples}}{\#\text{total}} \quad (2)$$

- 726 • **Platform Accuracy** indicates the tool preference recognition ability.

$$727 \text{platform\_acc} = \frac{\#\text{correct platform samples}}{\#\text{total}} \quad (3)$$

- 729 • **Query-related Parameter-Value Accuracy** indicates the ability to extract values from query.

$$731 \text{query\_param\_acc} = \frac{\#\text{correct query params}}{\#\text{total query params}} \quad (4)$$

- 733 • **Profile-related Parameter-Value Accuracy** indicates the ability to extract values from profile.

$$734 \text{profile\_param\_acc} = \frac{\#\text{correct profile params}}{\#\text{total profile params}} \quad (5)$$

- 736 • **Tool Name Accuracy** indicates the tool selection ability.

$$738 \text{tool\_name\_acc} = \frac{\#\text{correct name samples}}{\#\text{total}} \quad (6)$$

- 740 • **Tool Parameter Accuracy** indicates the tool comprehension ability.

$$741 \text{tool\_param\_acc} = \frac{\#\text{correct param samples}}{\#\text{total}} \quad (7)$$

- 743 • **Tool Parameter-Value Accuracy** indicate the value extraction on context ability.

$$745 \text{tool\_value\_acc} = \frac{\#\text{correct value samples}}{\#\text{total}} \quad (8)$$

- 747 • **Overall Accuracy on Trained Users** indicate the personalized tool ability on trained users.

$$748 \text{trained\_overall\_acc} = \frac{\#\text{correct trained samples}}{\#\text{trained total}} \quad (9)$$

- 750 • **Overall Accuracy on Untrained Users** indicate the personalized tool selection ability on trained  
 751 users.

$$752 \text{untrained\_overall\_acc} = \frac{\#\text{correct untrained samples}}{\#\text{untrained total}} \quad (10)$$

- 754 • **Overall Accuracy** indicate the overall personalized tool selection ability.

$$755 \text{overall\_acc} = \frac{\#\text{correct samples}}{\#\text{total}} \quad (11)$$

756  
 757 Table 5: Statistics of our synthesized dataset. The samples in the test set are verified by human  
 758 annotators. Trained and untrained represent the user profiles present and absent in the training set,  
 759 respectively. Unseen scenario represents additional data used in generalization study.

Dataset	#Scenario	#Platform	#API	#User	#Query
<b>Train</b>	5	15	360	74	7,096
<b>Test(PTBench)</b>	8	24	576	95	1,301
-Trained	5	15	360	74	474
-Untrained	5	15	360	6	609
-Unseen Scenarios	3	9	216	15	218
<b>Total</b>	8	21	576	95	8,397

## A.2 DETAILED RESULTS

771 **Dataset statistics.** The statistics of our training and test sets are illustrated in Table 5.

772 **Detailed results on trained and untrained user sets.** The detailed results of the trained and un-  
 773 trained user set on PTBench are illustrated in Table 11 and Table 12, respectively.

774 **Detailed results on three unseen scenarios.** The detailed results of unseen scenarios (lifestyle,  
 775 finance, and knowledge) are illustrated in Table 6, Table 7 and Table 8, respectively.

776 Table 6: Models performance results on lifestyle scenario. Bold represents the best result.

Model	Format	Platform	T-name	T-param	T-value	Overall
GPT-4-turbo	98.65	41.90	94.60	70.27	20.27	4.05
Deepseek-v3	100.00	47.30	97.30	<b>83.78</b>	21.62	6.76
Qwen2.5-7B-Instruct	81.08	20.27	72.97	21.62	4.05	1.35
Llama-3.1-8B-Instruct	97.30	39.19	93.24	64.87	10.81	2.70
Hammer2.1-7b	91.89	25.67	87.84	36.49	8.10	2.70
xLAM-7b-r	97.30	28.38	95.95	44.59	8.11	2.70
Ours	<b>100.00</b>	<b>63.51</b>	<b>97.30</b>	82.43	<b>27.03</b>	<b>14.86</b>

787 Table 7: Models performance results on finance scenario. Bold represents the best result.

Model	Format	Platform	T-name	T-param	T-value	Overall
GPT-4-turbo	90.48	40.48	76.19	54.76	21.43	7.14
Deepseek-v3	100.00	54.76	90.48	80.95	<b>23.81</b>	11.90
Qwen2.5-7B-Instruct	83.33	30.95	54.76	21.43	7.14	4.76
Llama-3.1-8B-Instruct	90.48	33.33	59.52	47.62	4.76	2.38
Hammer2.1-7b	92.86	38.10	80.95	38.10	7.14	7.14
xLAM-7b-r	100.00	26.19	90.48	35.71	11.90	7.14
Ours	<b>100.00</b>	<b>73.81</b>	<b>90.48</b>	<b>80.95</b>	21.43	<b>16.67</b>

## A.3 ADDITIONAL ANALYSIS ON RELATED WORKS.

801 **Comparison with Existing Personalized Tool Benchmarks.** We compare our benchmark with  
 802 existing personalized tool-use benchmarks, such as ETAPP (Hao et al., 2025) and ToolSpec-  
 803 trum (Cheng et al., 2025). Despite the concurrency, our benchmark maintains several substantive  
 804 advantages in terms of scale, personalization coverage, and profile design , which is shown in 9:

805

- 806 **• Broader coverage:** We include substantially more domains and tools.
- 807 **• Comprehensive personalization definition:** We jointly define personalization in both tool selec-  
 808 tion and parameter completion.

810 Table 8: Models performance results on knowledge scenario. Bold represents the best result.  
811

812 Model	813 Format	814 Platform	815 T-name	816 T-param	817 T-value	818 Overall
GPT-4-turbo	87.38	34.95	63.11	44.66	12.62	5.83
Deepseek-v3	87.38	43.69	77.67	70.87	17.48	6.80
Qwen2.5-7B-Instruct	85.29	33.33	52.94	27.45	2.94	0.98
Llama-3.1-8B-Instruct	94.12	23.53	59.80	45.10	10.78	1.96
Hammer2.1-7b	96.08	36.27	63.73	37.25	11.76	4.90
xLAM-7b-r	99.02	25.49	71.57	41.18	11.76	2.94
Ours	<b>100.00</b>	<b>80.39</b>	<b>82.35</b>	<b>68.63</b>	<b>27.45</b>	<b>21.57</b>

820

821

822 • **More realistic profile design:** Instead of using basic features + implicit preferences (where im-  
823 plicit preferences are directly provided), we structure profiles as basic features + user history. This  
824 design is closer to real-world settings, where users' implicit preferences (e.g., price sensitivity) are  
825 not explicitly written but must be inferred from behavioral history. This strengthens the practical  
826 value of our data synthesis pipeline.

827

828 Table 9: Comparison among personalized tool-use benchmarks.  
829

830 Benchmarks	831 #Tools	832 #Users	833 #Samples	834 Tool P.	835 Param P.	836 User Traj
ETAPP	35	16	800	✓	✗	✗
ToolSpectrum	42	158	1000	✓	✓	✗
Ours	<b>360</b>	85	<b>1301</b>	✓	✓	✓

837 **Generalization on other personalized tool benchmarks.** To further investigate whether our  
838 method generalizes beyond our proposed benchmark PTBench, we additionally evaluate the models  
839 on two other personalized tool-use benchmarks: ACEBench (Chen et al., 2025) and ToolSpectrum  
840 (Cheng et al., 2025). In addition to general-purpose model Qwen2.5-7B-Instruct, we also include  
841 PEToolLLM (Xu et al., 2025) as strong personalized tool-use baselines, ensuring a fair and com-  
842 prehensive comparison. The results are presented in Table 10. From the results, it is clear that our  
843 model achieves consistent and robust improvements across all three benchmarks, demonstrating not  
844 only its capability in general personalized tool invocation but also its effectiveness in handling di-  
845 verse tool-use patterns. These observations collectively show that our synthesized data and training  
846 strategy enable the model to generalize well across multiple personalized tool-use benchmarks.

847

848 Table 10: Performance on various personalized tool benchmarks. Bold represents the best result.  
849

850 Model	851 ACEBench	852 ToolSpectrum	853 PTBench
Qwen2.5-7B-Instruct	0.58	0.1759	0.0738
PETool-sft	0.34	0.1648	0.0433
PETool-sft-dpo	0.10	0.0133	0.0130
Ours	<b>0.66</b>	<b>0.1782</b>	<b>0.2678</b>

854  
855 

## B HUMAN-IN-THE-LOOP VERIFICATION

856 In the human verification stage, we adopt a systematic evaluation and refinement protocol consisting  
857 of three main steps.

861 

### B.1 DESIGNING EVALUATION CRITERIA

862 • **Query Reasonableness:** Ensures that queries include all required parameters, align with user  
863 profiles, and exclude meaningless characters.

864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

- **Platform Consistency:** Checks whether the platform preference implied in the query is consistent with the answer. If no explicit platform is specified, historical preferences from the user profile are used for verification.
- **Tool Invocation Accuracy:** Verifies that the invoked tool appropriately addresses the query and that its parameters are correctly specified.

## B.2 HUMAN ANNOTATION AND REFINEMENT

A human annotator reviews queries, answers, and tool invocations against the above criteria, making necessary corrections to ensure overall data quality.

A second annotator categorizes tool invocation parameters into two groups:

- **Query-dependent Parameters:** Explicitly provided in the user query.
- **Profile-dependent Parameters:** Not directly mentioned in the query but inferable from the user profile.

This classification enables a fine-grained evaluation of accuracy on different parameters.

## C EXAMPLES

To enhance the understanding of the proposed personalized tool invocation, we illustrate an example in Figure 7.

918 Table 11: Comparison with baseline models on trained users in PTBench. **Bold** and underline  
 919 represent the best and the 2nd best results.  
 920

921 Type	922 Model	923 Format	924 Preference		925 Param Value		926 Tool Invocation			927 Overall
			928 Platform	929 Query	930 Profile	931 T-name	932 T-param	933 T-value		
928 API	GPT-4-turbo	929 <b>0.9831</b>	0.5569	930 <b>0.7927</b>	0.7080	931 0.9325	0.7869	932 <b>0.3502</b>	0.1834	
	GPT-4o	0.8840	0.4157	0.6520	0.6164	0.8143	0.6941	0.2637	0.1350	
	Deepseek-v3	0.8903	0.5043	0.6868	0.6508	0.8376	0.7617	0.3059	0.1708	
	Deepseek-r1	0.8376	0.4958	0.6112	0.6317	0.7637	0.6604	0.2574	0.1477	
	Qwen-max	0.6941	0.4430	0.5083	0.5162	0.6393	0.5358	0.2152	0.1456	
	Claude-3.5-sonnet	0.9662	0.5822	0.7519	0.6794	0.7152	0.6498	0.2236	0.1329	
932 OSS	DeepSeek-R1-Distill-Llama-8B	0.6203	0.2891	0.3495	0.3111	0.4958	0.3925	0.1013	0.0485	
	DeepSeek-R1-Distill-Qwen-7B	0.6013	0.1519	0.2148	0.0954	0.3503	0.1941	0.0147	0.0042	
	Qwen2.5-7B-Instruct	0.7827	0.3882	0.5900	0.4447	0.6856	0.5612	0.1772	0.0717	
	Llama-3.1-8B-Instruct	0.8819	0.3797	0.6384	0.5439	0.8039	0.6498	0.2236	0.0929	
	Mistral-7B-Instruct-v0.3	0.8713	0.4198	0.5522	0.4113	0.6645	0.3734	0.1477	0.0674	
	Hammer2.1-7b	0.9641	0.3650	0.7126	0.5468	0.8439	0.6582	0.2257	0.0739	
	ToolACE-8B	0.4114	0.1709	0.3147	0.2061	0.3987	0.2721	0.0865	0.0338	
	Watt-tool-8B	0.3966	0.2405	0.2708	0.2156	0.3586	0.2510	0.0992	0.0591	
	xLAM-7b-r	0.9641	0.3586	0.6732	0.5315	0.8881	0.6329	0.2194	0.0696	
	Ours	0.9662	<b>0.7826</b>	<u>0.7791</u>	<b>0.7653</b>	<b>0.9409</b>	<b>0.8628</b>	<u>0.3333</u>	<b>0.2701</b>	

933 Table 12: Comparison with baseline models on untrained users in PTBench. **Bold** and underline  
 934 represent the best and the 2nd best results.  
 935

936 Type	937 Model	938 Format	939 Preference		940 Param Value		941 Tool Invocation			942 Overall
			943 Platform	944 Query	945 Profile	946 T-name	947 T-param	948 T-value		
945 API	GPT-4-turbo	946 <b>0.9737</b>	0.5419	947 <b>0.8266</b>	0.6637	948 0.9064	0.7586	949 <b>0.3531</b>	0.1856	
	GPT-4o	0.9146	0.4746	0.7596	0.6057	0.8391	0.7028	0.3054	0.1708	
	Deepseek-v3	0.9245	0.5468	0.7629	0.6343	0.8522	0.7455	0.3104	0.1757	
	Deepseek-r1	0.8062	0.4712	0.6443	0.5403	0.7175	0.6059	0.2660	0.1494	
	Qwen-max	0.8276	0.5353	0.6828	0.5658	0.7635	0.6207	0.2496	0.1707	
	Claude-3.5-sonnet	0.9704	0.5829	0.8046	0.6275	0.7077	0.6404	0.2397	0.1395	
950 OSS	DeepSeek-R1-Distill-Llama-8B	0.6601	0.3120	0.4061	0.2935	0.5173	0.3695	0.0953	0.0394	
	DeepSeek-R1-Distill-Qwen-7B	0.6158	0.1429	0.2481	0.1106	0.3777	0.2250	0.0279	0.0066	
	Qwen2.5-7B-Instruct	0.7882	0.3727	0.6301	0.3943	0.6815	0.5287	0.1889	0.0755	
	Llama-3.1-8B-Instruct	0.8900	0.4253	0.6839	0.4906	0.7964	0.6059	0.2052	0.0985	
	Mistral-7B-Instruct-v0.3	0.8489	0.3678	0.5653	0.3416	0.6584	0.3448	0.1429	0.0559	
	Hammer2.1-7b	0.9655	0.3629	0.7420	0.5094	0.8374	0.6109	0.2266	0.0689	
	ToolACE-8B	0.3974	0.1659	0.3392	0.2039	0.3810	0.2562	0.0936	0.0378	
	Watt-tool-8B	0.3580	0.2184	0.2722	0.1859	0.3268	0.2003	0.0706	0.0411	
	xLAM-7b-r	0.9442	0.3054	0.6839	0.4695	0.8538	0.5632	0.2233	0.0771	
951 Ours		0.9507	<b>0.7028</b>	0.8035	<b>0.7096</b>	<b>0.9112</b>	<b>0.8030</b>	<u>0.3481</u>	<b>0.2660</b>	

952  
 953  
 954  
 955  
 956  
 957  
 958  
 959  
 960  
 961  
 962  
 963  
 964  
 965  
 966  
 967  
 968  
 969  
 970  
 971

```

972 [SYSTEM]
973 You are given a user profile:
974 {
975     "basic_features": {
976         "username": "WineTraveler38",
977         ...
978     }
979     "user_history": {
980         "shopping": [
981             {
982                 "platform": "MegaMart",
983                 "action": "Purchased a selection of premium imported wines"
984             }
985         ]
986     }
987 }
988 Here is some platforms under the scenario:
989 [
990     {
991         "name": "MegaMart",
992         "profile": {
993             "product range": "A wide-ranging product from various categories."
994             ...
995         }
996     }
997     ...
998 ]
999 Here is some APIs under the platforms:
1000 [
1001     {
1002         'type': 'function',
1003         'function': {
1004             'name': 'registerUser',
1005             'description': 'Registers a new user in the application.',
1006             'parameters': {
1007                 'type': 'object',
1008                 'properties': {
1009                     'username': {
1010                         'type': 'string',
1011                         'description': '\"User\'s chosen username.\"'
1012                     },
1013                     ...
1014                 },
1015                 'required': ['username', 'password', 'email']
1016             },
1017             'response': {
1018                 'type': 'object',
1019                 'properties': {
1020                     'success': {
1021                         'type': 'boolean',
1022                         'description': 'Status of registration.'
1023                     },
1024                     ...
1025                 }
1026             }
1027         },
1028         ...
1029     }
1030     ...
1031 ]
1032 The user will give you a query. Based on the profile, try to solve the query by using the platforms and APIs. The platform you choose
1033 should fit the user profile or the needs of the user's query. All the necessary information are provided in the user profile. DO NOT
1034 ask the user for further information. You should respond in the format of {platform:[func1(param1.name = param1.value, param2...),
1035 func2...]} No other text MUST be included.
1036
1037 [USER]
1038 Could you please register an account for me using my username, password and email address, and setting my home location to my
1039 place of residence? I prefer not to receive any marketing emails.
1040
1041 [ASSISTANT]
1042 {
1043     MegaMart:[
1044         registerUser(
1045             username='WineTraveler38', password='strongpassword123!',
1046             email='jeanlucbordeaux@email.com', preferredLanguage='French',
1047             marketingConsent=False, homeLocation='Paris, France'
1048         )
1049     ]
1050 }

```

Figure 7: Illustration of personalized tool invocation task. The blue part is the expected output.