ADVANCING AND BENCHMARKING PERSONALIZED TOOL INVOCATION FOR LLMS

Anonymous authorsPaper under double-blind review

000

001

002 003 004

010 011

012

013

014

016

017

018

019

021

023

025 026 027

028 029

031

032

033

034

037

040

041

042

043

044

045

046

047

048

051

052

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**, the first 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 (Oin et al., 2024; Yan et al., 2024; Lin et al., 2024), such as ensuring adherence to the required tool invocation 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 and behavioral history before invoking tools. Two common scenarios illustrate this challenge on personalized tool invocation: (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 Personalization. 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.

056

060

061

062

063

064

065 066

067

068

069

071 072

073

074

075

076

077

079

081

082

083

084

085

087

880

089

091

092

094

095

096

098 099 100

101 102

103 104

105

106

107

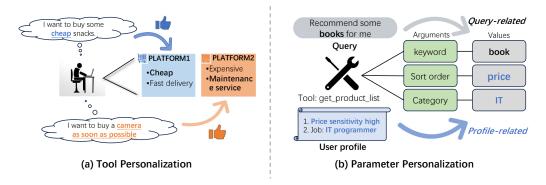


Figure 1: Examples of Personalized Tool Invocation. (a) Tool Personalization: Users may prefer different tools for similar functionalities depending on the query context. (b) Parameter Personalization: Certain tool parameters may be missing from the user's query and must be inferred from the user's profile.

In this work, we propose the novel task of personalized tool invocation, aiming to address the aforementioned critical challenges. 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**), the first benchmark designed to evaluate large models' ability in personalized tool invocation, consisting of 1,199 high-quality annotated data samples. Key contributions are summarized as follows:

- We propose the first paradigm for personalized tool invocation, incorporating both user tool personalization and parameter personalization, two key challenges in real-world applications.
- We develop a systematic personalized data synthesis framework and construct PTBench, the first benchmark for personalized tool invocation, enabling a comprehensive evaluation of models' 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 with few-shot learning, involving encouraging LLM to reason by providing examples (Yao et al., 2022), rewriting tool documentation with LLMs to enhance the comprehension (Yuan et al., 2024),

summarizing tool description with more concise and precise sentence (Xu et al., 2024), leveraging multi-agent collaboration to decompose the tool-calling task (Shi et al., 2024). Tuning-based methods leverage tool-learning samples to train existing LLMs, where the research problems comprise data collection and training strategy. Toolformer (Schick et al., 2024) and ToolkenGPT (Hao et al., 2024) add a special tool-related token into the vocabulary, switching the decoding process into tool selection and calling. Some works leverage advanced LLM to synthesize tool-calling samples to improve the tool-invocation ability of lightweight models, demonstrating the efficiency of the distillation from advanced models (Qin et al., 2024; Yang et al., 2023b; Liu et al., 2025).

2.2 Personalized LLMs

Personalized LLMs represent LLMs that have been adapted to align with user preferences and characteristics (Zhang et al., 2024c). Existing works mainly focus on the generation of personalized texts or applications in information systems. LLMs are customized as personal conversational AI assistants for various domains, including education (Kasneci et al., 2023; Dan et al., 2023; Park et al., 2024), healthcare (Belyaeva et al., 2023; Abbasian et al., 2024; Jin et al., 2024), finance (Liu et al., 2023; Lakkaraju et al., 2023), legal (Nguyen, 2023), and etc. User profiles are provided via prompts or hidden representation, leading the model to generate personalized text in the dialog. Personalized LLMs have been extensively applied in information systems such as recommender systems (Wu et al., 2023; Chen et al., 2024). LLMs are leveraged as an augmentation module for traditional recommender systems, serving as the content interpreter (Bao et al., 2023; Li et al., 2023; Yang et al., 2023a), the knowledge base (Xi et al., 2024; Wei et al., 2024), or the explainer (Lei et al., 2024; ?). Also, many works directly deploy LLMs as the direct recommenders via prompt techniques (?Hou et al., 2024) or fine-tuning (Zhang et al.). However, there is no work considering personalization in tool learning. This work is the first to propose personalized tool invocation for LLMs.

3 Personalized Tool Invocation

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

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

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

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

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

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

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

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

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

4.3 QUERY AND SOLUTION GENERATION

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

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

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

5 EXPERIMENTS

5.1 EXPERIMENTAL SETTINGS

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

Evaluation. We first evaluate the format accuracy by checking if the model's output can give formatted 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-

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

Туре	Model	Format	Tool-P	Param	Value	Tool	Invocati	on		Overall	
				Query	Profile	T-name T	T-param T	T-value	Trained U	Intrained (Overall
	GPT-4-turbo	97.78	54.84	81.23	68.32	91.78	77.09	35.18	18.34	18.56	18.47
	GPT-4o	90.12	44.84	71.44	61.04	82.83	69.91	28.69	13.50	17.08	15.51
API	Deepseek-v3	90.95	52.80	73.09	64.16	84.60	75.30	30.85	17.08	17.57	17.36
API	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	<u>58.26</u>	78.24	65.04	71.10	64.45	23.26	13.29	13.95	13.67
	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
OSS	Mistral-7B-v0.3	85.87	39.03	55.98	37.23	66.12	35.72	14.50	6.74	5.59	6.09
OSS	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	73.74	79.33	73.41	92.42	82.90	34.17	27.01	26.60	26.78

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

Baselines. We compare the latest open-source models and API-based models, as well as fine-tuned tool-calling models. Open-source models include DeepSeek-R1-Distill-Llama-8B(DeepSeek-AI, 2025), DeepSeek-R1-Distill-Qwen-7B(DeepSeek-AI, 2025), Qwen2.5-7B-Instruct(Team, 2024a;b), Llama-3.1-8B-Instruct (AI@Meta, 2024) and Mistral-7B-Instruct-v0.3(Jiang et al., 2023). API-based models include GPT-4-turbo¹, GPT-4o¹, Deepseek-v3(DeepSeek-AI, 2024), Deepseek-r1(DeepSeek-AI, 2025), Qwen-max(Team, 2024b) and Claude-3.5-sonnet². Models fine-tuned for tool-calling include Hammer2.1-7b(Lin et al., 2024), ToolACE-8B(Liu et al., 2025), watt-tool-8B³ and xLAM-7b-r(Zhang et al., 2024b; Liu et al., 2024; Zhang et al., 2024a).

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

5.2 Main Results

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

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

https://chatgpt.com

²https://www.anthropic.com

³https://ollama.com

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

Data	Untrained	Trained	Overall
All	26.60	27.01	26.78
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

Finding 2: Most models fall short on the tool-preference task, including the state-of-the-art model—GPT-4-turbo, indicating the high complexity of selecting a suitable one from several similar tools according to the user profile. Our model outperforms nearly all models in all aspects by a considerable improvement, presenting the necessity of personalized tool-invocation enhancement.

Finding 3: Our model demonstrates a significant improvement in its performance across various tasks on PTBench. Notably, the enhancement in the Tool Preference task is particularly pronounced when compared to the pre-trained Qwen2.5-7B-Instruct model. This also indicates that, even without additional manual verification of the training data, the model achieves a high accuracy, demonstrating the effectiveness of the proposed synthesis framework. Additionally, our model shows a significant improvement on untrained users, presenting the generalization of the model.

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

5.3 ABLATION STUDY

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

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

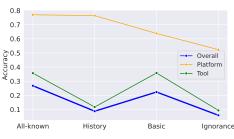
First, We use the four dataset variants to train and then evaluate the model with the consistent input. The results are reported in Table 2. From the result, we can observe that the existence of user history and basic features hold contributions to the overall performance of the model to an extent.

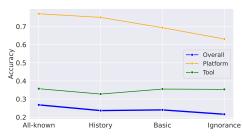
Additionally, we conduct experiments under two settings: (1) train the model with the All variant and evaluate the model with the four variants, illustrated in Figure 3a; (2) train the model with the four variants and evaluate the models with the All variant, illustrated in Figure 3b. The results exhibit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

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

Table 3: Evaluating Model Performance Without Personalized Information.

Model	Format	Platform	T-name	T-param	T-value	Overall
GPT-40	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





- (a) User profile ablation in evaluation.
- (b) User profile ablation in training.

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

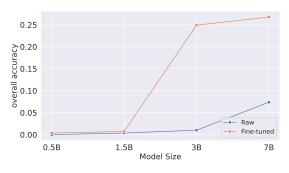
5.4 ERROR ANALYSIS

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

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

5.5 FURTHER ANALYSIS

Model Scaling. For the purpose of analyzing the influence of model size on the performance of our trained model, we utilize models with different sizes in the Qwen2.5 series, including 7B, 3B, 1.5B and 0.5B. The results are shown in Figure 5. We can observe that the 1.5B and 0.5B model only show slight improvement from the training, while 3B and 7B model gain substantial improvement from the training. This demonstrate that the personalized tool invocation is a high-level capability of LLMs, requiring a certain scale of parameters.



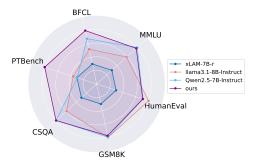


Figure 5: Study of model scaling. (Base model: Qwen2.5-series.)

Figure 6: General Capabilities Analysis. (Base model: Qwen2.5-7B-Instruct.)

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	95.69	41.38	87.93	64.65	20.69	5.17
Deepseek-v3	100.00	50.00	94.83	82.76	22.41	8.62
Qwen2.5-7B-Instruct	75.00	24.14	66.38	21.55	5.17	2.59
Llama-3.1-8B-Instruct	88.80	37.07	81.03	58.62	8.62	2.59
Hammer2.1-7b	92.24	30.17	85.34	37.07	7.76	4.31
xLAM-7b-r	78.45	18.97	76.73	35.34	31.02	2.59
Ours	100.00	67.24	94.83	81.89	25.00	15.52

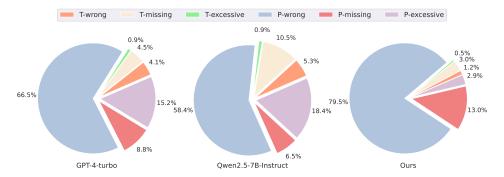


Figure 4: Error Analysis on PTBench. T-wrong, T-missing, and T-excessive represent wrong tools, missing tools and excessive tools, respectively. P-missing, P-excessive and P-error represent missing parameters, excessive parameters and wrong parameters, respectively.

Generalization to Unseen Scenarios. To further examine the generalizability of our model beyond the five common scenarios, we conducted an additional evaluation on unseen domains. Specifically, we synthesized 116 samples covering two new scenarios: finance and lifestyle. As shown in Table 4, our model consistently outperforms the baselines in these settings, demonstrating strong robustness and adaptability. These results provide evidence that the proposed methods and benchmark are not limited to the initial set of scenarios, but can extend to a broader range of real-world scenarios.

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(CommonSenceQA(Talmor et al., 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.

REFERENCES

- Mahyar Abbasian, Zhongqi Yang, Elahe Khatibi, Pengfei Zhang, Nitish Nagesh, Iman Azimi, Ramesh Jain, and Amir M Rahmani. Knowledge-infused llm-powered conversational health agent: A case study for diabetes patients. *arXiv preprint arXiv:2402.10153*, 2024.
- AI@Meta. Llama 3 model card. 2024. URL https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md.
 - Keqin Bao, Jizhi Zhang, Yang Zhang, Wenjie Wang, Fuli Feng, and Xiangnan He. Tallrec: An effective and efficient tuning framework to align large language model with recommendation. In *Proceedings of the 17th ACM Conference on Recommender Systems*, pp. 1007–1014, 2023.
 - Anastasiya Belyaeva, Justin Cosentino, Farhad Hormozdiari, Krish Eswaran, Shravya Shetty, Greg Corrado, Andrew Carroll, Cory Y McLean, and Nicholas A Furlotte. Multimodal llms for health grounded in individual-specific data. In *Workshop on Machine Learning for Multimodal Health-care Data*, pp. 86–102. Springer, 2023.
 - Jin Chen, Zheng Liu, Xu Huang, Chenwang Wu, Qi Liu, Gangwei Jiang, Yuanhao Pu, Yuxuan Lei, Xiaolong Chen, Xingmei Wang, et al. When large language models meet personalization: Perspectives of challenges and opportunities. *World Wide Web*, 27(4):42, 2024.
 - Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, and et al. Evaluating large language models trained on code. 2021.
 - Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
 - Yuhao Dan, Zhikai Lei, Yiyang Gu, Yong Li, Jianghao Yin, Jiaju Lin, Linhao Ye, Zhiyan Tie, Yougen Zhou, Yilei Wang, et al. Educhat: A large-scale language model-based chatbot system for intelligent education. *arXiv preprint arXiv:2308.02773*, 2023.
 - DeepSeek-AI. Deepseek-v3 technical report, 2024. URL https://arxiv.org/abs/2412.19437.
 - DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL https://arxiv.org/abs/2501.12948.
 - Shibo Hao, Tianyang Liu, Zhen Wang, and Zhiting Hu. Toolkengpt: Augmenting frozen language models with massive tools via tool embeddings. *Advances in neural information processing systems*, 36, 2024.
 - Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob Steinhardt. Aligning ai with shared human values. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021a.
 - Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021b.
- Yupeng Hou, Junjie Zhang, Zihan Lin, Hongyu Lu, Ruobing Xie, Julian McAuley, and Wayne Xin Zhao. Large language models are zero-shot rankers for recommender systems. In *European Conference on Information Retrieval*, pp. 364–381. Springer, 2024.
- Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=nZeVKeeFYf9.

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

- Mingyu Jin, Qinkai Yu, Dong Shu, Chong Zhang, Lizhou Fan, Wenyue Hua, Suiyuan Zhu, Yanda Meng, Zhenting Wang, Mengnan Du, et al. Health-llm: Personalized retrieval-augmented disease prediction system. *arXiv preprint arXiv:2402.00746*, 2024.
- Enkelejda Kasneci, Kathrin Seßler, Stefan Küchemann, Maria Bannert, Daryna Dementieva, Frank Fischer, Urs Gasser, Georg Groh, Stephan Günnemann, Eyke Hüllermeier, et al. Chatgpt for good? on opportunities and challenges of large language models for education. *Learning and individual differences*, 103:102274, 2023.
- Kausik Lakkaraju, Sai Krishna Revanth Vuruma, Vishal Pallagani, Bharath Muppasani, and Biplav Srivastava. Can llms be good financial advisors?: An initial study in personal decision making for optimized outcomes. *arXiv preprint arXiv:2307.07422*, 2023.
- Yuxuan Lei, Jianxun Lian, Jing Yao, Xu Huang, Defu Lian, and Xing Xie. Recexplainer: Aligning 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.
- Ruyu Li, Wenhao Deng, Yu Cheng, Zheng Yuan, Jiaqi Zhang, and Fajie Yuan. Exploring the upper limits of text-based collaborative filtering using large language models: Discoveries and insights. *arXiv preprint arXiv:2305.11700*, 2023.
- Qiqiang Lin, Muning Wen, Qiuying Peng, Guanyu Nie, Junwei Liao, Jun Wang, Xiaoyun Mo, Jiamu Zhou, Cheng Cheng, Yin Zhao, Jun Wang, and Weinan Zhang. Hammer: Robust function-calling for on-device language models via function masking, 2024. URL https://arxiv.org/abs/2410.04587.
- Weiwen Liu, Xingshan Zeng, Xu Huang, xinlong hao, Shuai Yu, Dexun Li, Shuai Wang, Weinan Gan, Zhengying Liu, Yuanqing Yu, Zezhong WANG, Yuxian Wang, Wu Ning, Yutai Hou, Bin Wang, Chuhan Wu, Wang Xinzhi, Yong Liu, Yasheng Wang, Duyu Tang, Dandan Tu, Lifeng Shang, Xin Jiang, Ruiming Tang, Defu Lian, Qun Liu, and Enhong Chen. ToolACE: Enhancing function calling with accuracy, complexity, and diversity. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=8EB8k6DdCU.
- Xiao-Yang Liu, Guoxuan Wang, Hongyang Yang, and Daochen Zha. Fingpt: Democratizing internet-scale data for financial large language models. *arXiv preprint arXiv:2307.10485*, 2023.
- Z Liu, Z Lai, Z Gao, E Cui, Z Li, X Zhu, L Lu, Q Chen, Y Qiao, J Dai, et al. Controlllm: augment language models with tools by searching on graphs (2023). *arXiv preprint arXiv:2310.17796*.
- Zuxin Liu, Thai Hoang, Jianguo Zhang, Ming Zhu, Tian Lan, Shirley Kokane, Juntao Tan, Weiran Yao, Zhiwei Liu, Yihao Feng, et al. Apigen: Automated pipeline for generating verifiable and diverse function-calling datasets. *arXiv* preprint arXiv:2406.18518, 2024.
- Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- Ha-Thanh Nguyen. A brief report on lawgpt 1.0: A virtual legal assistant based on gpt-3. *arXiv* preprint arXiv:2302.05729, 2023.
- Minju Park, Sojung Kim, Seunghyun Lee, Soonwoo Kwon, and Kyuseok Kim. Empowering personalized learning through a conversation-based tutoring system with student modeling. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*, pp. 1–10, 2024.

- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, dahai li, Zhiyuan Liu, and Maosong Sun. ToolLLM: Facilitating large language models to master 16000+ real-world APIs. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=dHng200Jjr.
 - Changle Qu, Sunhao Dai, Xiaochi Wei, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, Jun Xu, and Ji-Rong Wen. Tool learning with large language models: A survey. *Frontiers of Computer Science*, 19(8):198343, 2025.
 - Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36, 2024.
- Zhengliang Shi, Shen Gao, Xiuyi Chen, Yue Feng, Lingyong Yan, Haibo Shi, Dawei Yin, Pengjie Ren, Suzan Verberne, and Zhaochun Ren. Learning to use tools via cooperative and interactive agents. pp. 10642–10657, Miami, Florida, USA, November 2024. doi: 10.18653/v1/2024. findings-emnlp.624. URL 2024.findings-emnlp.624/.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4149–4158, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1421. URL https://aclanthology.org/N19-1421/.
- Qwen Team. Qwen2 technical report. arXiv preprint arXiv:2407.10671, 2024a.
- Qwen Team. Qwen2.5: A party of foundation models, September 2024b. URL https://qwenlm.github.io/blog/qwen2.5/.
- Wei Wei, Xubin Ren, Jiabin Tang, Qinyong Wang, Lixin Su, Suqi Cheng, Junfeng Wang, Dawei Yin, and Chao Huang. Llmrec: Large language models with graph augmentation for recommendation. In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*, pp. 806–815, 2024.
- Likang Wu, Zhi Zheng, Zhaopeng Qiu, Hao Wang, Hongchao Gu, Tingjia Shen, Chuan Qin, Chen Zhu, Hengshu Zhu, Qi Liu, Hui Xiong, and Enhong Chen. A survey on large language models for recommendation. *CoRR*, abs/2305.19860, 2023.
- Yunjia Xi, Weiwen Liu, Jianghao Lin, Xiaoling Cai, Hong Zhu, Jieming Zhu, Bo Chen, Ruiming Tang, Weinan Zhang, and Yong Yu. Towards open-world recommendation with knowledge augmentation from large language models. In *Proceedings of the 18th ACM Conference on Recommender Systems*, pp. 12–22, 2024.
- Yang Xu, Yunlong Feng, Honglin Mu, Yutai Hou, Yitong Li, Xinghao Wang, Wanjun Zhong, Zhongyang Li, Dandan Tu, Qingfu Zhu, Min Zhang, and Wanxiang Che. Concise and precise context compression for tool-using language models. pp. 16430–16441, Bangkok, Thailand, August 2024. doi: 10.18653/v1/2024.findings-acl.974. URL 2024.findings-acl.974/.
- Fanjia Yan, Huanzhi Mao, Charlie Cheng-Jie Ji, Tianjun Zhang, Shishir G. Patil, Ion Stoica, and Joseph E. Gonzalez. Berkeley function calling leaderboard. https://gorilla.cs.berkeley.edu/blogs/8_berkeley_function_calling_leaderboard.html, 2024.
- Fan Yang, Zheng Chen, Ziyan Jiang, Eunah Cho, Xiaojiang Huang, and Yanbin Lu. Palr: Personalization aware Ilms for recommendation. *arXiv preprint arXiv:2305.07622*, 2023a.
- Rui Yang, Lin Song, Yanwei Li, Sijie Zhao, Yixiao Ge, Xiu Li, and Ying Shan. GPT4tools: Teaching large language model to use tools via self-instruction. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023b. URL https://openreview.net/forum?id=cwjh8lqmOL.

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

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

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

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

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

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

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

A EXPERIMENTS

A.1 EVALUATION METRICS

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

• Format Accuracy indicates the instruction-following ability.

$$format_acc = \frac{\#parsable \, samples}{\#total} \tag{2}$$

Platform Accuracy indicates the tool preference recognition ability.

$$platform_acc = \frac{\#correct \, platform \, samples}{\#total} \tag{3}$$

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

$$query_param_acc = \frac{\#correct\ query\ params}{\#total\ query\ params} \tag{4}$$

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

$$profile_param_acc = \frac{\#correct\ profile\ params}{\#total\ profile\ params} \tag{5}$$

• Tool Name Accuracy indicates the tool selection ability.

$$tool_name_acc = \frac{\#correct \, name \, samples}{\#total} \tag{6}$$

• Tool Parameter Accuracy indicates the tool comprehension ability.

$$tool_param_acc = \frac{\#correct \, param \, samples}{\#total} \tag{7}$$

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

Dataset	#Scenario	#Platform	#API	#User	#Query
Train	5	15	360	74	7,096
Test(PTBench)	7	21	504	85	1,199
-Trained	5	15	360	74	474
-Untrained	5	15	360	6	609
-Unseen Scenarios	2	6	144	5	116
Total	5	15	360	80	8,197

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

$$tool_value_acc = \frac{\#correct\ value\ samples}{\#total} \tag{8}$$

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

$$trained_overall_acc = \frac{\#correct\ trained\ samples}{\#trained\ total} \tag{9}$$

Overall Accuracy on Untrained Users indicate the personalized tool selection ability on trained users.

$$untrained_overall_acc = \frac{\#correct\ untrained\ samples}{\#untrained\ total} \tag{10}$$

• Overall Accuracy indicate the overall personalized tool selection ability.

$$overall_acc = \frac{\#correct \, samples}{\#total}$$
 (11)

A.2 DETAILED RESULTS

the detailed component of the dataset are illustrated in Table 5

The detailed results of the trained and untrained subset on PTBench are illustrated in Table 6 and Table 7, respectively.

B HUMAN-IN-THE-LOOP VERIFICATION

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

B.1 Designing Evaluation Criteria

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

Table 6: Comparison with baseline models on trained users in PTBench. **Bold** and <u>underline</u> represent the best and the 2nd best results.

Type	Model	Format	Preference	erence Param Value		Too	Overall		
. 1			Platform	Query	Profile	T-name	T-param	T-value	
	GPT-4-turbo	0.9831	0.5569	0.7927	0.7080	0.9325	0.7869	0.3502	0.1834
	GPT-4o	0.8840	0.4157	0.6520	0.6164	0.8143	0.6941	0.2637	0.1350
API	Deepseek-v3	0.8903	0.5043	0.6868	0.6508	0.8376	0.7617	0.3059	0.1708
API	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
	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
OSS	Mistral-7B-Instruct-v0.3	0.8713	0.4198	0.5522	0.4113	0.6645	0.3734	0.1477	0.0674
USS	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	0.7826	0.7791	0.7653	0.9409	0.8628	0.3333	0.2701

B.3 TOOL PARAMETER CLASSIFICATION

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 model accuracy across different parameter types.

C EXAMPLES

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

Table 7: Comparison with baseline models on untrained users in PTBench. **Bold** and <u>underline</u> represent the best and the 2nd best results.

Type	vpe Model		Preference	Param	Value	Too	ol Invocat	ion	Overall
-J F -		Format	Platform	Query	Profile	T-name	T-param	T-value	
	GPT-4-turbo	0.9737	0.5419	0.8266	0.6637	0.9064	0.7586	0.3531	0.1856
	GPT-4o	0.9146	0.4746	0.7596	0.6057	0.8391	0.7028	0.3054	0.1708
API	Deepseek-v3	0.9245	0.5468	0.7629	0.6343	0.8522	0.7455	0.3104	0.1757
API	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
'	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
OSS	Mistral-7B-Instruct-v0.3	0.8489	0.3678	0.5653	0.3416	0.6584	0.3448	0.1429	0.0559
USS	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
	Ours	0.9507	0.7028	0.8035	0.7096	0.9112	0.8030	0.3481	0.2660

```
864
             [SYSTEM]
865
             You are given a user profile:
866
               "basic_features":{
    "username":"WineTraveler38",
867
868
869
               "user_history":{
870
                  "shopping":[{
                            "platform": "MegaMart",
871
                            "action": "Purchased a selection of premium imported wines"
872
                      }
873
                  ]
874
875
876
             Here is some platforms under the scenario:
877
             [ {
878
                       "name": "MegaMart",
                       "profile":{
879
                            "product range": "A wide-ranging selection, offering products from various categories."
880
                  }
882
             ]
883
             Here is some APIs under the platforms:
884
             [ {
885
                      'type': 'function',
886
                      'function': {
                           'name': 'registerUser',
887
                            'description': 'Registers a new user in the application.',
888
                            'parameters': {
                                 'type': 'object',
889
                                 'properties': {
                                      'username': {
    'type': 'string',
890
891
                                          'description': \"User's chosen username.\"
892
                                     },
893
894
                                 required': ['username', 'password', 'email']
895
                            'response': {
   'type': 'object',
                                 'properties': {
897
                                      'success': {
   'type': 'boolean',
898
                                          'description': 'Status of registration.'
899
                                     },
900
                                }
901
                          }
902
                      }
                  },
903
904
905
             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
             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
906
             ask the user for further information. You should respond in the format of {platform:[func1(param1_name = param1_value, param2...),
907
             func2...]} No other text MUST be included.
908
909
             Could you please register an account for me using my username, password and email address, and setting my home location to my
             place of residence? I prefer not to receive any marketing emails.
910
911
             [ASSISTANT]
912
               MegaMart:[
913
                  registerUser(
914
                    username='WineTraveler38', password='strongpassword123!',
                    email='jeanlucbordeaux@email.com', preferredLanguage='French',
915
                    marketingConsent=False, homeLocation='Paris, France'
916
917
```

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