

PSPA-BENCH: A Personalized Benchmark for Smartphone GUI Agent

Anonymous ACL submission

Abstract

Smartphone GUI agents operate directly on app interfaces, offering broad capability without deep system integration. However, real-world use is highly personalized: users follow diverse workflows and preferences, demanding customized rather than generic assistance. Existing benchmarks overlook this personalization dimension due to limited user-specific data and the lack of fine-grained evaluation. To address this gap, we introduce PSPA-BENCH, a benchmark dedicated to evaluating personalization in smartphone GUI agents. Specifically, PSPA-BENCH utilizes the *task decomposition graph*, TDG to represent the structure of personalized task at unit-instruction level. Using TDG, PSPA-BENCH generates task templates that instantiate personalized instructions without large-scale user logs, and supports fine-grained evaluation of GUI agents’ personalization abilities. PSPA-BENCH includes over 12,855 realistic personalized instructions spanning 10 scenarios and 22 commonly used mobile apps, along with four fine-grained evaluation metrics. We evaluate 11 state-of-the-art GUI agents and find that existing methods perform poorly under personalized settings, with even the best agent achieving limited success. Based on these results, we derive five insights to guide future research, positioning PSPA-BENCH as a foundation for advancing personalized GUI agents.

1 Introduction

Smartphone GUI agents (Zhang et al., 2025a; Liu et al., 2025a) have emerged as a promising paradigm for facilitating human-computer interaction in mobile environments. Different from conventional automation systems (e.g., Siri¹) that rely on system APIs, GUI agents (Zhang et al., 2025b; Wang et al., 2024a; Qin et al., 2025; Wen et al., 2024) operate directly on the graphical user interfaces of mobile applications. By emulating user

¹<https://www.apple.com/siri/>

actions such as tapping, swiping, and text input, these agents remain effective in ecosystems with limited or no backend integration.

However, unlike many other computing platforms, smartphone usage is inherently personalized (Li et al., 2015; Nguyen et al., 2025). As shown in Figure 1, users vary significantly in how they accomplish the same task through diverse action sequences (Tian et al., 2020). For instance, when ordering takeout, different users may prefer different search strategies and different food categories. Consequently, the demand for smartphone GUI agents is shifting away from *generic, one-size-fits-all* solutions toward *personalized, context-aware* assistants. Moreover, a GUI agent is not intended to act as a *one-time* helper; rather, it is expected to provide continuous support over *long-term use* (Lee et al., 2024b; Wang et al., 2025; Hu et al., 2025). This requires the agent to accumulate execution experience and adapt to evolving individual preferences, delivering a “more understanding of you” interaction experience progressively.

Existing benchmarks evaluate GUI agents from different perspectives. For example, Mobile-Bench (Deng et al., 2024), AndroidWorld (Rawles et al., 2025) and Spa-Bench (Chen et al., 2025a), focus on the general task-execution capabilities of GUI agents; TransBench (Lu et al., 2025) emphasize the transferability of GUI agents; AEIA-MN (Chen et al., 2025b) and Mobile Safety Bench (Lee et al., 2024a) focus on the safety of GUI agents. However, considering both the **sparsity of user data** and the **absence of suitable metrics for personalized evaluation**, the systematic assessment of GUI agents’ personalization capabilities remains largely unexplored.

To address this gap, we introduce PSPA-BENCH, a benchmark for systematic study of personalization in smartphone GUI agents. PSPA-BENCH uses the Task Decomposition Graph (TDG) to explicitly represent the structure of personalized

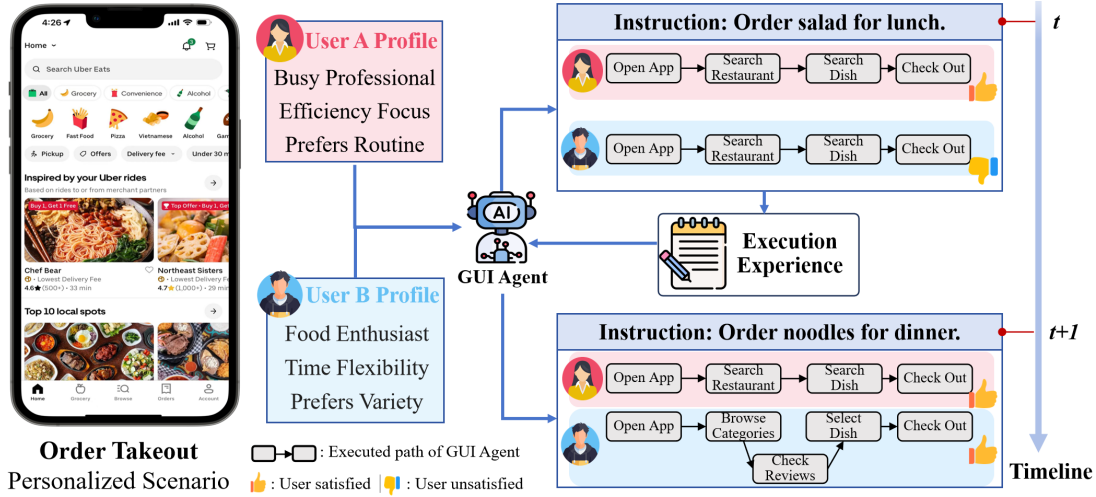


Figure 1: GUI agents are shifting from *general-purpose, one-time* task execution to *personalized, long-term* service. (i) Given the same instruction, an agent should adapt its execution path to user preferences; (ii) For the same user, the agent uses historical execution experience to provide long-term support; experience up to time t informs a refined execution at $t + 1$ for similar instructions. The formulation of personalized GUI agent task is provided in Appendix B.

GUI tasks. To compensate for sparse user data, PSPA-BENCH generates personalized instructions by building TDG-derived templates without relying on large-scale user logs. For fine-grained evaluation, PSPA-BENCH measures the ratio of completed unit instructions as the fine-grained metrics by through graph-trace alignment at the unit-instruction level. PSPA-BENCH covers 10 daily-use scenarios (e.g., shopping, travel, dining) across 22 mobile apps, defines 100 user personas, and yields 12,855 personalized instructions. A detailed comparison with existing smartphone GUI benchmarks is provided in Appendix A. With the proposed PSPA-BENCH, we conduct comprehensive evaluations with 11 GUI agents, and the evaluation results from PSPA-BENCH demonstrate that there is still a long way for the development of personalized GUI agents. Furthermore, we highlight 5 key findings that provide directions for the development of personalized GUI agents: (i) effective **perception** underpins personalized instruction execution; (ii) strong **reasoning** is required to handle ambiguity and complex behaviors; (iii) personalization amplifies the trade-off between **accuracy and efficiency**; (iv) **persistent memory** is essential for long-term adaptation; and (v) **self-evolution** further enhances this adaptive capability. To summarize, our contributions are as follows:

- The problem of personalization in smartphone GUI agents is systematically formalized, highlighting the scarcity of personalized user data and

the lack of evaluation metrics as core challenges.

- We propose PSPA-BENCH, a benchmark with a data generation method for sparse user data and process-oriented evaluation metrics based on task decomposition graphs. PSPA-BENCH covers 10 representative daily scenarios, 22 third-party apps, 100 user personas, and 12,855 personalized instructions
- 11 state-of-the-art GUI agents are benchmarked on PSPA-BENCH, and the results reveal that current methods perform poorly under personalized settings, and further distill 5 key directions for future research based on the evaluation results.

2 Benchmarking Setup

Despite rapid progress in GUI agents, the systematic assessment of their personalization capabilities remains largely unexplored. Current benchmarks face two main limitations. First, privacy regulations (e.g., GDPR (European Parliament and Council of the European Union, 2016)) and user concerns restrict the collection of user-specific data, hindering the creation of realistic yet privacy-compliant datasets. Second, evaluations rely on coarse metrics such as binary success signals (Chen et al., 2025a; Rawles et al., 2025; Yang et al., 2025) or task completion rates (Xing et al., 2024), which assume fixed action paths and fail to assess adaptability to diverse user preferences or learning over time.

To overcome these limitations, we propose PSPA-BENCH², a benchmark specifically designed for personalized GUI agents (Figure 2). At the core of PSPA-BENCH is the *task decomposition graph* (TDG), which explicitly captures the personalized structure of GUI agent tasks. Based on the TDG, we explain how PSPA-BENCH generate personalized task instructions without relying on large-scale user logs, and evaluate agent executions with unit-instruction level fine-grained metrics. Finally, we present the comparison GUI agents employed in our benchmark.

2.1 Task Decomposition Graph

To address two limitations identified above, i.e., the sparsity of real user data and the lack of metrics tailored to personalization, we introduce the task decomposition graph (TDG). The TDG explicitly captures the personalized structure of GUI tasks. This structure (i) enables controlled generation of personalized task instances even when user traces are scarce, and (ii) supports fine-grained, preference-aware evaluation beyond coarse end-state signals.

Specifically, GUI tasks such as "Buy a piece of clothing on a shopping app" implicitly combine universal steps with user-specific requirements (e.g., preferred product category or price range). To make these requirements explicit, we represent each task as a directed acyclic task decomposition graph $\mathcal{G} = (\mathcal{U}, \mathcal{E})$, where nodes are minimal GUI actions and edges denote execution dependencies. Any source-to-sink path defines a valid task completion. To model personalization, we distinguish fixed nodes \mathcal{U}_f for user-invariant steps (e.g., launching the app) and flexible nodes \mathcal{U}_p for preference-dependent actions (e.g., selecting preferred categories).

Therefore, compared with vague task descriptions, the TDG provides a more explicit and structured representation of the task with fine-grained unit instructions. This graph-based representation can serve as a unified foundation for generating personalized instructions and for conducting fine-grained, evaluation of agent behavior. The constructed TDGs are provided in Appendix C.

²The code and data of PSPA-BENCH are available at <https://anonymous.4open.science/status/PSPA-Bench>

2.2 Template-Driven Personalized Instruction Generation

To address user data sparsity, we use the TDG as a generative blueprint to create diverse, realistic personalized instructions without large-scale user logs. The task structure ensures universal steps are preserved while systematically varying preference-dependent details.

Template construction. The TDG explicitly separates fixed nodes \mathcal{U}_f from flexible nodes \mathcal{U}_p . We use the fixed nodes and dependencies \mathcal{E} as a procedural backbone, and convert flexible nodes into parameterized slots. This template preserves the shared task structure while clearly indicating where personalization occurs.

Specifically, each template $\text{TP}_{\mathcal{G}}$ encodes the fixed procedural backbone of the task (derived from \mathcal{U}_f) together with a slot schema $\mathcal{L} = \{l_{u_p}\}_{u_p \in \mathcal{U}_p}$ mapping each flexible node u_p to a slot l_{u_p} . Every slot $l \in \mathcal{L}$ is annotated with a semantic type $\tau(l)$ (e.g., `product_category`, `price_range`) and a domain of admissible values. The textual backbone preserves the execution ordering implied by \mathcal{E} , while slots mark precisely where personalization should be injected. We illustrate an example of the template for the shopping scenario as below:

Template Example for Shopping Scenario

```
Launch the shopping app, then search for
a <product_category> within <price_range>,
and add the top result to the cart.
```

Here, "`<product_category>`" and "`<price_range>`" are slots derived from flexible nodes in the TDG, while the rest reflects the fixed structure.

Slot instantiation. Once a template is selected, we instantiate it by filling each parameterized slot with values sampled from the user's preference distribution. For a user profile ρ , each slot l is assigned a value $v_l \sim D_{\tau(l)}^{\rho}$, yielding a concrete instruction $f(\text{TP}_{\mathcal{G}}, v_l) \mapsto I$ that respects the dependencies in \mathcal{G} . In this framework, the user profile ρ directly controls the *personalization strength*. A set of similar user profiles induces low-diversity slot values and thus weaker personalization, whereas a set of diverse profiles leads to more varied instructions, corresponding to stronger personalization signals. We provide the detailed instruction instantiation process in Appendix D.

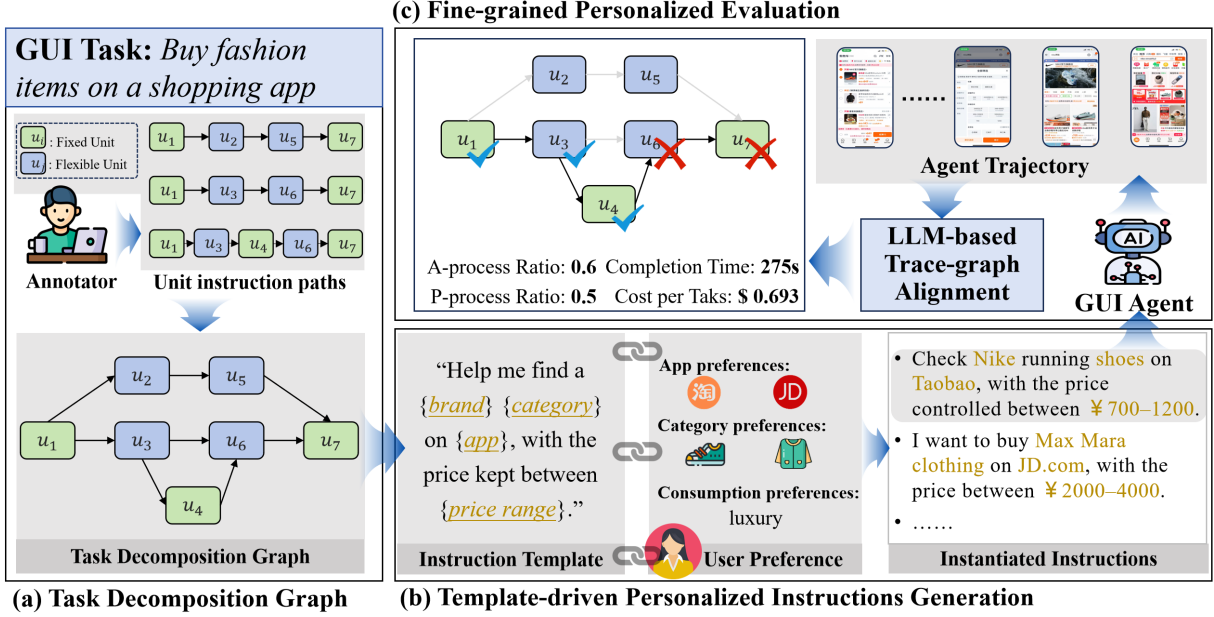


Figure 2: Benchmark framework of PSPA-BENCH. (a) **Task decomposition graph**: the task is decomposed into a directed acyclic graph of unit instructions, where fixed nodes denote universal steps and flexible nodes denote user-specific requirements. (b) **Template-driven personalized instruction generation**: the TDG is used to construct task templates, which are then instantiated with user preferences to generate personalized instructions. (c) **Fine-grained personalized evaluation**: The trace-graph alignment is used to compute APR (A-Progress Ratio) and PPR (P-Progress Ratio) for fine-grained evaluation of both immediate and long-term performance.

Task difficulty control. To control task difficulty, we vary instruction templates along two dimensions: (i) *Complexity level* (low, medium, high), based on the number of unit instructions in the corresponding TDG. (ii) *Clarity level* (high, medium, low), which governs how explicitly the instruction communicates its intent. We provide examples of templates in personalized scenarios in Appendix F

2.3 Fine-grained Personalized Evaluation

We evaluate agents along two complementary objectives in fine-grained level (Section 1): (i) an *immediate objective*: accurately and efficiently executing the current instruction, and (ii) a *long-term objective*: adapting to user-specific preferences over long-term use.

Let the trace be $\mathcal{T} = (a_1, \dots, a_L)$, where a_i denotes the i -th action. Since a TDG may contain multiple valid source-to-sink paths, we identify the path $P \subseteq \mathcal{G}$ that best matches the trace using a two-stage alignment procedure. First, we define an LLM-based alignment method $m : \mathcal{T} \rightarrow \mathcal{U} \cup \{\emptyset\}$ that assigns each action either to the unit instruction it fulfills or to \emptyset if it does not correspond to any instruction. Second, we select the optimal path P^* from all valid paths \mathcal{P} in the TDG by maximizing

the number of matched unit instructions:

$$P^* = \arg \max_{P \in \mathcal{P}} \sum_{u \in P} \mathbb{1}[\exists a_i \in \mathcal{T} : m(a_i) = u] \quad (1)$$

where $\mathbb{1}[\cdot]$ is the indicator function.

In case of ties (multiple paths achieving the same maximum match count), we apply the following tie-breaking rules in order: (i) prefer paths with more completed *flexible* nodes \mathcal{U}_p , as they better reflect personalization; (ii) prefer shorter paths (fewer total nodes), favoring efficiency; (iii) if ties persist, select the path whose matched nodes appear earliest in the trace, reflecting the agent’s primary execution intent. Our evaluation begins by aligning the agent’s execution trace with the corresponding TDG. The resulting *trace-graph* alignment indicates (i) which unit instructions were successfully completed and (ii) where failures occurred, with attribution to fixed nodes \mathcal{U}_f (universal steps) and flexible nodes \mathcal{U}_p (preference-sensitive steps).

Figure 3: Evaluation metrics taxonomy of PSPA-BENCH.

Performance	Immediate	Long-term
	APR	Δ APR
PPR	Δ PPR	
Efficiency	CT	Δ CT
	CPT	Δ CPT

Building on this, we assess these objectives along two metric dimensions (Figure 3):

- **Performance** measures execution quality at the unit-instruction level: **A-process Ratio (APR)**: completion ratio over all unit instructions on the executed path. **P-process Ratio (PPR)**: completion ratio restricted to flexible nodes \mathcal{U}_p , directly reflecting personalization.
- **Efficiency** quantifies the resources required for task execution: **Completion Time (CT)**: measured as the average elapsed time needed to finish a task. **Cost Per Task (CPT)**, which denotes the average monetary cost (USD) of completing a task, estimated from the number of model tokens consumed.

Each metric is paired with an incremental variant ‘ Δ ’ that measures improvements over time, thereby capturing the agent’s ability to adapt. For example, $\Delta\text{APR} = \text{APR}_{T_2} - \text{APR}_{T_1}$ denotes performance gains between task T_1 and T_2 , with analogous increments defined for PPR, CT, and CPT. Implementation details of evaluation are provided in Appendix E.

2.4 Comparison Methods

As shown in Table 1, PSPA-BENCH benchmarks a wide range of baselines and recent agent frameworks, summarized along five dimensions: **Perception, Reasoning, Memory, Evolution, and Collaboration**. **Perception** refers to how agents sense the GUI, either via accessibility (A11y) trees or visual inputs. **Reasoning** covers planning and reflection. **Memory** distinguishes temporary memory (within a task) from persistent memory (reused across tasks). **Evolution** considers whether behaviors or toolkits can be updated, and **Collaboration** characterizes the number and roles of agents.

For the **LLM category**, LLM-Base is a standard LLM (e.g., GPT-4o-mini). LLM+A11y augments it with A11y trees, while LLM+A11y+SoM adds set of marks (SoM) in screenshots for stronger visual grounding. LLM+A11y+Thinking further enhances reasoning with advanced models (e.g., GPT-O3-Mini). For **general agent frameworks**, we include ReAct (Yao et al., 2023), which interleaves reasoning and action; Self-Refine (Madaan et al., 2023), which iteratively improves via feedback; and Reflexion (Shinn et al., 2023), which uses reflection memory across tasks. For the **GUI agent**

frameworks, PSPA-BENCH evaluates four representatives: AppAgent (Zhang et al., 2025b), which uses multimodal grounding through GUI parsing and screenshots; M3A (Rawles et al., 2025), a GUI-specialized ReAct; Mobile-Agent V2 (Wang et al., 2024a), a multi-agent system with planning, decision, and reflection; and Mobile-Agent E (Wang et al., 2025), which extends V2 with hierarchical collaboration, persistent memory, and self-evolution.

3 Empirical Results

3.1 Experiments on Immediate Objective

We evaluate each method’s ability to execute a single personalized instruction effectively and efficiently. Under high-clarity settings and three task complexity levels (low, medium, high), we report unit-level metrics (APR, PPR) and efficiency metrics (CT, CPT) in Table 2. The results show clear performance differences across model classes, revealing the capabilities needed for personalized task execution.

(1) **Perception is the foundation for personalized instruction execution.** The drastic improvement from LLM-Base (APR = 0.002, PPR = 0) to +A11y (APR = 0.625, PPR = 0.586 at low complexity) highlights that access to GUI context is indispensable. Personalized tasks often rely on dynamic, user-specific content, such as preferred items, recent activity, or customized filters, that cannot be interpreted through language alone. Perception provides the necessary grounding, enabling the agent to map abstract instructions to the actual UI elements visible on the screen. Without this capability, the agent is effectively blind to the task’s personalized components.

(2) **Reasoning is essential for resolving ambiguity and executing complex personalized behaviors.** Perception allows the agent to observe, but reasoning allows it to decide. Many personalized instructions contain implicit preferences, conditional steps, or multi-hop objectives (e.g., “find a deal that matches your usual budget and dietary preferences”). Reasoning-enhanced models such as +A11y+Thinking and ReAct show significant gains over perception-only agents as complexity increases. For instance, ReAct achieves APR = 0.503 and PPR = 0.468 at medium complexity, outperforming simpler baselines. These models can handle flexible execution paths, maintain coherence across steps, and adapt to user-specific constraints

Table 1: Comparison of evaluated methods in PSPA-BENCH. The table presents baseline LLMs, general agent frameworks, and GUI agent frameworks, assessed across five key dimensions: Perception, Reasoning, Memory, Evolution, and Collaboration.

Method		Perception		Reasoning		Memory		Evolution		Collaboration	
Type	Name	Text	Vision	Plan-ning	Reflec-tion	Tempo-rary	Persis-tent	Rules	Toolk-its	Single	Multi-ple
LLM	Base		✓			✓				✓	
	+A11y	✓				✓				✓	
	+A11y+SoM	✓	✓			✓				✓	
	+A11y+Thinking	✓		✓		✓				✓	
General Agent Framework	ReAct	✓	✓	✓		✓				✓	
	Self-Refine	✓	✓		✓	✓					
	Reflexion	✓	✓	✓	✓	✓	✓				✓
GUI Agent Framework	AppAgent	✓	✓	✓		✓				✓	
	M3A	✓	✓	✓		✓				✓	
	Mobile-AgentV2		✓	✓	✓	✓	✓				✓
	Mobile-Agent E		✓	✓	✓	✓	✓	✓	✓		✓

Table 2: Mean evaluation results (APR, PPR, CT, CPT) are reported across five scenarios under high-clarity conditions with varying task complexity (low, medium, high). The results indicate each method’s effectiveness in achieving the immediate objective of personalized GUI tasks.

Method	Setting: Clarity fixed at middle, varying complexity											
	Complexity: low				Complexity: middle				Complexity: high			
	APR↑	PPR↑	CT↓	CPT↓	APR↑	PPR↑	CT↓	CPT↓	APR↑	PPR↑	CT↓	CPT↓
LLM-Base	0.002	0	236.7	0.612	0.001	0	256.6	0.501	0.002	0	293.4	0.816
+A11y	0.625	0.586	176.7	0.056	0.431	0.449	189.0	0.038	0.364	0.356	229.3	0.062
+A11y+SoM	0.643	0.596	243.6	0.671	0.463	0.461	254.3	0.525	0.373	0.372	293.3	0.930
+A11y+Thinking	0.674	0.624	197.7	0.673	0.466	0.464	233.0	0.726	0.382	0.378	274.8	0.884
ReAct	0.683	0.618	217.0	0.752	0.503	0.468	243.0	0.616	0.414	0.403	316.8	0.761
Self-Refine	0.596	0.542	249.1	0.758	0.438	0.441	731.4	0.796	0.308	0.339	845.3	0.931
Reflexion	0.665	0.611	250.7	0.657	0.465	0.466	332.5	0.694	0.438	0.406	464.5	0.776
AppAgent	0.654	0.592	218.7	0.329	0.469	0.495	227.8	0.337	0.373	0.402	292.6	0.502
M3A	0.685	0.623	317.6	0.951	0.521	0.485	340.7	1.010	0.460	0.403	450.6	1.507
Mobile-Agent V2	0.674	0.618	517.6	1.017	0.532	0.533	578.9	1.364	0.435	0.415	659.1	1.414
Mobile-Agent E	0.695	0.621	532.7	1.084	0.546	0.539	587.5	1.401	0.468	0.435	716.7	1.692

— all of which are critical in the context of personalization.

(3) **Personalization introduces a stronger trade-off between execution accuracy and efficiency.** While Mobile-Agent E achieves the highest APR (0.695) and PPR (0.621), it also exhibits the slowest execution (CT = 716.7 at high complexity), reflecting a strategy that prioritizes accuracy through exhaustive interaction. This behavior may stem from the need to explore interface options and disambiguate user preferences on the fly. In contrast, AppAgent adopts a more efficient but less exhaustive approach (CT = 292.6 at high complexity), with slightly lower accuracy. In personalized settings, where goals are less explicit and vary across users, striking the right balance between precision and efficiency becomes more challenging than in generic task settings.

3.2 Experiments on long-term Objective

We evaluate each method’s ability to support long-term adaptation to user preferences. To simulate sustained interaction, agents first perform tasks of medium complexity and clarity to accumulate experience, and are then evaluated against a baseline without prior experience. We report relative changes in Δ APR, Δ PPR, Δ CT, and Δ CPT, defined as the performance difference with and without experience accumulation. Evaluation is conducted under two settings: (i) increasing task complexity with fixed clarity, and (ii) increasing ambiguity with fixed complexity. Table 3 presents the results, demonstrating the importance of memory and experience processing for long-term personalization.

Formally, let M denote an evaluation metric (e.g., APR, PPR, CT, CPT). For a user u , we de-

Table 3: Mean values of evaluation metrics (Δ APR, Δ PPR, Δ CT, Δ CPT) are reported across five scenarios under two settings: (i) high clarity with varying task complexity (low, medium, high), and (ii) medium complexity with varying task clarity (low, medium, high). All methods first perform tasks with medium complexity and clarity to acquire initial experience. The results reflect each method’s ability to adapt to user preferences in long-term personalized GUI tasks. (Note: For readability, Δ APR, Δ PPR and Δ CPT values are scaled by 100.)

Method	Setting (i): Clarity fixed at high, varying complexity											
	Complexity: low				Complexity: middle				Complexity: high			
	Δ APR \uparrow	Δ PPR \uparrow	Δ CT \downarrow	Δ CPT \downarrow	Δ APR \uparrow	Δ PPR \uparrow	Δ CT \downarrow	Δ CPT \downarrow	Δ APR \uparrow	Δ PPR \uparrow	Δ CT \downarrow	Δ CPT \downarrow
LLM-Base	0	0	9.28	0.20	0	0	-17.68	0.73	0	0	-19.18	0.94
+A11y	0.66	-0.58	-12.73	-0.63	-0.39	0.05	-2.72	0.42	0.22	-0.72	-8.31	-0.27
+A11y+SoM	-0.09	0.57	12.01	-0.03	0.18	-0.91	4.30	0.66	-0.87	0.90	-18.63	0.62
+A11y+Thinking	-0.39	-0.80	7.37	-0.12	-0.76	-0.01	-18.62	0.82	-0.48	0.33	-7.53	0.04
ReAct	0.09	-0.63	-18.78	0.55	0.88	0.79	-3.92	0.84	-0.82	-0.61	-18.19	0.35
Self-Refine	-0.22	-0.46	13.15	-0.29	-0.44	0.09	-14.36	0.60	-0.85	0.97	10.89	-0.60
Reflexion	4.53	3.71	-38.35	-2.46	3.54	3.29	-40.62	-2.77	2.90	1.25	-26.76	-1.87
AppAgent	-0.38	-0.35	-9.18	0.28	0.77	-0.06	-15.22	0.43	0.52	0.12	10.84	-0.10
M3A	0.05	-0.14	18.98	-0.78	-0.94	0.27	-7.43	0.02	0.82	-0.50	3.58	0.51
Mobile-Agent V2	4.84	3.26	-36.48	-3.70	4.54	4.02	-41.45	-2.43	3.41	1.04	-26.02	-3.58
Mobile-Agent E	6.96	4.08	-62.08	-4.06	6.72	3.68	-57.02	-3.41	2.13	2.44	-35.51	-1.94

Method	Setting (ii): Complexity fixed at middle, varying clarity											
	Clarity: low				Clarity: middle				Clarity: high			
	Δ APR \uparrow	Δ PPR \uparrow	Δ CT \downarrow	Δ CPT \downarrow	Δ APR \uparrow	Δ PPR \uparrow	Δ CT \downarrow	Δ CPT \downarrow	Δ APR \uparrow	Δ PPR \uparrow	Δ CT \downarrow	Δ CPT \downarrow
LLM-Base	0	0	17.72	-0.35	0	0	-5.45	0.94	0	0	-0.11	0.40
+A11y	-0.43	-0.93	4.38	-0.01	-0.90	-0.44	16.33	-0.52	-0.71	-0.02	19.43	-0.52
+A11y+SoM	0.34	0.52	-10.49	0.46	-0.26	0.26	5.34	-0.07	-0.82	0.67	-7.17	0.63
+A11y+Thinking	-0.92	0.18	-7.10	0.97	0.02	-0.55	-5.81	0.65	0.38	-0.23	17.47	-0.72
ReAct	-0.32	-0.77	16.99	-0.75	-0.48	0.32	-12.69	0.11	0.06	-0.52	-16.28	0.79
Self-Refine	0.80	0.27	6.44	-0.30	0.45	0.79	-15.48	0.56	0.28	-0.83	-13.53	0.80
Reflexion	4.23	4.58	-37.50	-2.75	3.70	2.75	-39.54	-2.78	2.11	1.71	-31.34	-1.33
AppAgent	-0.35	0.49	5.99	-0.70	0.32	0.14	-16.25	0.26	-0.47	-0.51	18.92	-0.21
M3A	0.78	0.26	11.79	-0.65	0.15	-0.01	-12.19	0.44	-0.44	-0.95	-5.82	0.65
Mobile-Agent V2	4.27	3.15	-36.17	-2.95	3.72	3.62	-43.49	-2.1	2.26	1.61	-25.23	-1.31
Mobile-Agent E	6.54	7.15	-63.81	-3.08	5.13	4.94	-53.44	-2.28	5.55	4.08	-40.59	-2.75

note M_u^{with} as the metric value after experience accumulation and $M_u^{\text{w/o}}$ as the metric value without accumulation. The relative improvement is defined as: $\Delta M = \frac{1}{N} \sum_{u=1}^N (M_u^{\text{with}} - M_u^{\text{w/o}})$, where N is the number of evaluation tasks. For APR and PPR, a positive ΔM indicates better task completion, while for CT and CPT, a negative ΔM indicates improved efficiency. Table 3 reports the results, from which we draw two key findings:

(1) **Persistent memory is the foundation for long-term adaptation.** Models without persistent memory, such as ReAct and M3A, show little to no improvement, and in many cases even regress (e.g., negative Δ APR and Δ PPR across most settings). In contrast, methods with persistent memory, e.g., Reflexion and Mobile-Agent V2, achieve consistent gains. This confirms that retaining and reusing past experience is critical for adapting to personalized, evolving user behaviors.

(2) **Self-evolution enhances long-term adaptation.** Persistent memory provides the foundation, but Mobile-Agent E goes further by enabling self-evolution: it can update execution rules and create

shortcuts from accumulated experience in its persistent memory. This capability yields the strongest gains overall (e.g., Δ APR = +6.72, Δ PPR = +3.68 at medium complexity). Compared to Mobile-Agent V2, which mainly reuses stored memory, Mobile-Agent E actively improves its strategy over time, leading to better generalization to new but related personalized tasks.

3.3 Ranking of Methods on Personalized GUI Benchmark

To further assess method performance on PSPA-BENCH, we rank all methods with respect to both the immediate and long-term objectives. As shown in Fig. 4, for the immediate objective (left), there is a clear trade-off between execution accuracy (APR, PPR) and efficiency (CT, CPT): Mobile-Agent E achieves the highest accuracy but incurs substantial execution cost, whereas AppAgent and LLM+A11y are faster but less accurate. This suggests that GUI agents must balance precision and efficiency. For the long-term objective (right), methods with persistent memory, e.g., Reflexion,

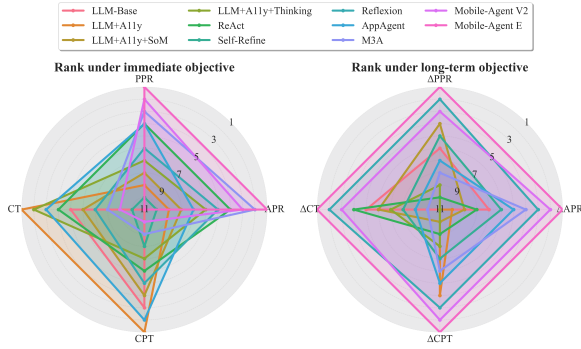


Figure 4: Radar plots showing the comparative ranks of different methods under two objectives. **Left:** Immediate objective, evaluated with APR, PPR, CT, and CPT. **Right:** Long-term objective, evaluated with Δ APR, Δ PPR, Δ CT, and Δ CPT. Each axis indicates the rank of a method on the corresponding metric, with lower values representing better performance.

Mobile-Agent V2, and especially Mobile-Agent E, achieve consistently higher improvements across Δ APR and Δ PPR. This confirms that retaining and reusing past experience is necessary for adapting to evolving user preferences, with stronger memory mechanisms yielding greater gains.

4 Related Works

4.1 Smartphone GUI Agent

Smartphone GUI agents (Zhang et al., 2025a; Liu et al., 2025a) automate tasks by interpreting visual and textual elements on mobile screens. Prior work (Nguyen et al., 2025) has explored several design dimensions. For perception, current agents rely on structured UI data such as A1ly trees (Rawles et al., 2025; Wen et al., 2024), or use vision-based methods for richer context (Zhang et al., 2025b; Gou et al., 2025). For reasoning and planning, frameworks like ReAct (Yao et al., 2023) and Reflexion (Shinn et al., 2023) have been integrated into M3A (Rawles et al., 2025) and Mobile-Agent V2 (Wang et al., 2024a) to handle complex tasks, while GUI-Explorer (Xie et al., 2025) uses prior knowledge of app operations. Recent studies also focus on learning from execution experience: Mobile-Agent E (Wang et al., 2025) and AppAgentX (Jiang et al., 2025) summarize repetitive low-level actions into high-level shortcuts to improve efficiency. Despite recent progress, most agents still target general automation, leaving personalization underexplored. Recently, FingerTip (Chen et al., 2025a) introduces a benchmark for evaluating personalization and shows that existing agents perform poorly on such

tasks. However, its coarse-grained metrics may lose important details when assessing adaptation to diverse preferences or long-term use.

4.2 Personalized LLM Agent

Personalized LLM agents (Dong et al., 2023; Tseng et al., 2024) aim to tailor assistance to individual preferences. Approaches fall into two paradigms: fine-tuning and training-free methods. Fine-tuning approaches, such as Personalized Soups (Jang et al., 2023) and (Liu et al., 2025b), adapt model parameters using explicit user feedback. While effective in capturing preferences across multiple users, these methods require large amounts of data and significant computational resources. Training-free methods instead keep the base model fixed, incorporating personalization via prompt engineering (Zhang et al., 2025c; Gao et al., 2024), retrieval-augmented generation (RAG) (Salemi et al., 2024b,a), or external memory (Wang et al., 2024b; Zhong et al., 2024). For example, CIPHER (Gao et al., 2024) uses user edits to infer preferences, while EMG-RAG (Wang et al., 2024b) organizes user-specific information with editable memory graphs. While effective in static tasks such as travel planning, applying these strategies to dynamic mobile GUI environments remains an open challenge.

5 Conclusion

We presented PSPA-BENCH, a benchmark for evaluating personalization in smartphone GUI agents. PSPA-BENCH introduces the TDG to explicitly separate universal and preference-sensitive steps, enabling both realistic personalized instruction generation and fine-grained evaluation. Covering 10 scenarios, 28 apps, and 100 user personas, PSPA-BENCH provides a large-scale, privacy-compliant testbed. Experiments with 11 representative models show that existing agents struggle on personalized tasks, revealing substantial room for improvement. At the same time, our analysis highlights that personalization requires agents to go beyond generic automation, emphasizing grounding in GUI context, structured reasoning, and mechanisms for accumulating and evolving user-specific experience. By providing a principled framework and comprehensive evaluation, PSPA-BENCH lays the foundation for advancing research on personalized GUI agents and offers a path toward more adaptive and user-centered mobile AI systems.

543 Limitations

544 While PSPA-BENCH provides a broad benchmark
545 for evaluating personalization in smartphone GUI
546 agents, several limitations should be noted. First,
547 although the benchmark covers 22 apps across 10
548 representative scenarios, this is still a limited sub-
549 set of the vast mobile application ecosystem, and
550 the findings may not fully generalize to other do-
551 mains such as healthcare, finance, or specialized
552 professional tools. Second, the user personas and
553 preferences in PSPA-BENCH are synthetically gener-
554 ated based on predefined templates rather than
555 collected from real users, which may not fully cap-
556 ture the complexity, inconsistency, and evolving
557 nature of authentic user behaviors. Third, the cur-
558 rent benchmark is primarily designed for Android
559 devices, limiting its direct applicability to other
560 platforms such as iOS, where UI conventions and
561 interaction patterns differ. Finally, the task decom-
562 position graphs are manually constructed by do-
563 main experts, which, while ensuring quality, may
564 not cover all possible task variations or edge cases
565 that arise in real-world usage. These limitations
566 also point to promising directions for future work,
567 including extending the benchmark to broader ap-
568 plication domains and platforms, incorporating real
569 user data while preserving privacy, and developing
570 automated methods for TDG construction.

571 References

572 Jingxuan Chen, Derek Yuen, Bin Xie, Yuhao Yang,
573 Gongwei Chen, Zhihao Wu, Li Yixing, Xurui Zhou,
574 Weiwen Liu, Shuai Wang, Kaiwen Zhou, Rui Shao,
575 Liqiang Nie, Yasheng Wang, Jianye Hao, Jun Wang,
576 and Kun Shao. 2025a. SPA-BENCH: a compre-
577 hensive benchmark for smartphone agent evaluation. In
578 *International Conference on Learning Representations*.
579

580 Yurun Chen, Xueyu Hu, Keting Yin, Juncheng Li,
581 and Shengyu Zhang. 2025b. AEIA-MN: evaluating
582 the robustness of multimodal LLM-powered mobile
583 agents against active environmental injection attacks.
584 *arXiv preprint arXiv:2502.13053*.

585 Shihan Deng, Weikai Xu, Hongda Sun, Wei Liu, Tao
586 Tan, Liu Jianfeng, Liu Jianfeng, Ang Li, Jian Luan, Bin
587 Wang, Rui Yan, and Shuo Shang. 2024. Mobile-
588 bench: An evaluation benchmark for LLM-based
589 mobile agents. In *Annual Meeting of the Association
590 for Computational Linguistics*, pages 8813–8831.

591 Xin Luna Dong, Seungwhan Moon, Yifan Ethan Xu,
592 Kshitiz Malik, and Zhou Yu. 2023. Towards next-
593 generation intelligent assistants leveraging LLM tech-

niques. In *ACM SIGKDD Conference on Knowledge
Discovery and Data Mining*, pages 5792–5793. 594
595

European Parliament and Council of the European 596
Union. 2016. Regulation (eu) 2016/679 (general 597
data protection regulation). Official Journal of the 598
European Union. 599

Ge Gao, Alexey Taymanov, Eduardo Salinas, Paul 600
Mineiro, and Dipendra Misra. 2024. Aligning LLM 601
agents by learning latent preference from user ed- 602
its. In *Advances in Neural Information Processing 603
Systems*. 604

Boyu Gou, Ruohan Wang, Boyuan Zheng, Yanan Xie, 605
Cheng Chang, Yiheng Shu, Huan Sun, and Yu Su. 606
2025. Navigating the digital world as humans do: 607
Universal visual grounding for GUI agents. In *Inter- 608
national Conference on Learning Representations*. 609

Xueyu Hu, Tao Xiong, Biao Yi, Zishu Wei, Ruixuan 610
Xiao, Yurun Chen, Jiasheng Ye, Meiling Tao, Xi- 611
angxin Zhou, Ziyu Zhao, Yuhuai Li, Shengze Xu, 612
Shenzhi Wang, Xinchun Xu, Shuofei Qiao, Zhaokai 613
Wang, Kun Kuang, Tiejong Zeng, Liang Wang, and 614
10 others. 2025. OS Agents: A survey on MLLM- 615
based agents for computer, phone and browser use. 616
In *Annual Meeting of the Association for Computa- 617
tional Linguistics*, pages 7436–7465. 618

Joel Jang, Seungone Kim, Bill Yuchen Lin, Yizhong 619
Wang, Jack Hessel, Luke Zettlemoyer, Hannaneh 620
Hajishirzi, Yejin Choi, and Prithviraj Ammanabrolu. 621
2023. PERSONALIZED SOUPS: Personalized large 622
language model alignment via post-hoc parameter 623
merging. In *Advances in Neural Information Process- 624
ing Systems Workshop Adaptive Foundation Models*. 625

Wenjia Jiang, Yangyang Zhuang, Chenxi Song, 626
Xu Yang, and Chi Zhang. 2025. AppAgentX: Evolv- 627
ing GUI agents as proficient smartphone users. *arXiv 628
preprint arXiv:2503.02268*. 629

J Richard Landis and Gary G Koch. 1977. The mea- 630
surement of observer agreement for categorical data. 631
biometrics, pages 159–174. 632

Juyong Lee, Dongyoon Hahm, June Suk Choi, 633
W. Bradley Knox, and Kimin Lee. 2024a. Mo- 634
bileSafetyBench: Evaluating safety of autonomous 635
agents in mobile device control. *arXiv preprint 636
arXiv:2410.17520*. 637

Sunjae Lee, Junyoung Choi, Jungjae Lee, Munim Hasan 638
Wasi, Hojun Choi, Steven Y. Ko, Sangeun Oh, and 639
Insik Shin. 2024b. MobileGPT: Augmenting LLM 640
with human-like app memory for mobile task automa- 641
tion. In *Annual International Conference on Mobile 642
Computing and Networking*, pages 1119–1133. 643

Huoran Li, Xuan Lu, Xuanzhe Liu, Tao Xie, Kaigui 644
Bian, Felix Xiaozhu Lin, Qiaozhu Mei, and Feng 645
Feng. 2015. Characterizing smartphone usage pat- 646
terns from millions of Android users. In *ACM Inter- 647
net Measurement Conference*, pages 459–472. 648

649	Guangyi Liu, Pengxiang Zhao, Liang Liu, Yaxuan Guo, Han Xiao, Weifeng Lin, Yuxiang Chai, Yue Han, Shuai Ren, Hao Wang, Xiaoyu Liang, Wenhao Wang, Tianze Wu, Linghao Li, Hao Wang, Guanqing Xiong, Yong Liu, and Hongsheng Li. 2025a. LLM-powered GUI agents in phone automation: Surveying progress and prospects. <i>arXiv preprint arXiv:2504.19838</i> .	Alireza Salemi, Sheshera Mysore, Michael Bendersky, and Hamed Zamani. 2024b. LaMP: When large language models meet personalization. In <i>Annual Meeting of the Association for Computational Linguistics</i> , pages 7370–7392.	706 707 708 709 710
656	Jiongnan Liu, Yutao Zhu, Shuting Wang, Xiaochi Wei, Erxue Min, Yu Lu, Shuaiqiang Wang, Dawei Yin, and Zhicheng Dou. 2025b. LLMs + Persona-Plug = personalized LLMs. In <i>Annual Meeting of the Association for Computational Linguistics</i> , pages 9373–9385.	Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. Reflexion: language agents with verbal reinforcement learning. In <i>Advances in Neural Information Processing Systems</i> .	711 712 713 714 715
662	Yuheng Lu, Qian Yu, Hongru Wang, Zeming Liu, Wei Su, Yanping Liu, Yuhang Guo, Maocheng Liang, Yunhong Wang, and Haifeng Wang. 2025. Trans-Bench: Breaking barriers for transferable graphical user interface agents in dynamic digital environments. In <i>Findings of the Association for Computational Linguistics: ACL</i> , pages 12464–12478, Vienna, Austria.	Yuan Tian, Ke Zhou, Mounia Lalmas, and Dan Pelleg. 2020. Identifying tasks from mobile app usage patterns. In <i>International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , pages 2357–2366.	716 717 718 719 720
669	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhunoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023. Self-Refine: Iterative refinement with self-feedback. In <i>Advances in Neural Information Processing Systems</i> .	Yu-Min Tseng, Yu-Chao Huang, Teng-Yun Hsiao, Wei-Lin Chen, Chao-Wei Huang, Yu Meng, and Yun-Nung Chen. 2024. Two tales of persona in LLMs: A survey of role-playing and personalization. In <i>Findings of the Association for Computational Linguistics: EMNLP</i> , pages 16612–16631.	721 722 723 724 725 726
677	Dang Nguyen, Jian Chen, Yu Wang, Gang Wu, Namyong Park, Zhengmian Hu, Hanjia Lyu, Junda Wu, Ryan Aponte, Yu Xia, Xintong Li, Jing Shi, Hongjie Chen, Viet Dac Lai, Zhouhang Xie, Sungchul Kim, Ruiyi Zhang, Tong Yu, Mehrab Tanjim, and 11 others. 2025. GUI agents: A survey. In <i>Findings of the Association for Computational Linguistics: ACL</i> , pages 22522–22538.	Junyang Wang, Haiyang Xu, Haitao Jia, Xi Zhang, Ming Yan, Weizhou Shen, Ji Zhang, Fei Huang, and Jitao Sang. 2024a. Mobile-Agent-v2: Mobile device operation assistant with effective navigation via multi-agent collaboration. In <i>Advances in Neural Information Processing Systems</i> .	727 728 729 730 731 732
685	Yujia Qin, Yining Ye, Junjie Fang, Haoming Wang, Shihao Liang, Shizuo Tian, Junda Zhang, Jiahao Li, Yunxin Li, Shijue Huang, Wanjuan Zhong, Kuanye Li, Jiale Yang, Yu Miao, Woyu Lin, Longxiang Liu, Xu Jiang, Qianli Ma, Jingyu Li, and 16 others. 2025. UI-TARS: Pioneering automated GUI interaction with native agents. <i>arXiv preprint arXiv:2501.12326</i> .	Zheng Wang, Zhongyang Li, Zeren Jiang, Dandan Tu, and Wei Shi. 2024b. Crafting personalized agents through retrieval-augmented generation on editable memory graphs. In <i>Conference on Empirical Methods in Natural Language Processing</i> , pages 4891–4906.	733 734 735 736 737 738
692	Christopher Rawles, Sarah Clinckemahillie, Yifan Chang, Jonathan Waltz, Gabrielle Lau, Marybeth Fair, Alice Li, William E. Bishop, Wei Li, Folawiyo Campbell-Ajala, Daniel Kenji Toyama, Robert James Berry, Divya Tyamagundlu, Timothy P. Lillicrap, and Oriana Riva. 2025. AndroidWorld: A dynamic benchmarking environment for autonomous agents. In <i>International Conference on Learning Representations</i> .	Zhenhailong Wang, Haiyang Xu, Junyang Wang, Xi Zhang, Ming Yan, Ji Zhang, Fei Huang, and Heng Ji. 2025. Mobile-Agent-E: Self-evolving mobile assistant for complex tasks. <i>arXiv preprint arXiv:2501.11733</i> .	739 740 741 742 743
700	Alireza Salemi, Surya Kallumadi, and Hamed Zamani. 2024a. Optimization methods for personalizing large language models through retrieval augmentation. In <i>International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , pages 752–762.	Hao Wen, Yuanchun Li, Guohong Liu, Shanhui Zhao, Tao Yu, Toby Jia-Jun Li, Shiqi Jiang, Yunhao Liu, Yaqin Zhang, and Yunxin Liu. 2024. Autodroid: Llm-powered task automation in Android. In <i>Annual International Conference on Mobile Computing and Networking</i> , pages 543–557.	744 745 746 747 748 749
701		Bin Xie, Rui Shao, Gongwei Chen, Kaiwen Zhou, Yinchuan Li, Jie Liu, Min Zhang, and Liqiang Nie. 2025. GUI-explorer: Autonomous exploration and mining of transition-aware knowledge for GUI agent. In <i>Annual Meeting of the Association for Computational Linguistics</i> , pages 5650–5667.	750 751 752 753 754 755
702		Mingzhe Xing, Rongkai Zhang, Hui Xue, Qi Chen, Fan Yang, and Zhen Xiao. 2024. Understanding the weakness of large language model agents within a complex Android environment. In <i>ACM SIGKDD Conference on Knowledge Discovery and Data Mining</i> , pages 6061–6072.	756 757 758 759 760 761

762 Qinglong Yang, Haoming Li, Haotian Zhao, Xiaokai
763 Yan, Jingtao Ding, Fengli Xu, and Yong Li. 2025.
764 FingerTip 20K: A benchmark for proactive and
765 personalized mobile LLM agents. *arXiv preprint*
766 *arXiv:2507.21071*.

767 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak
768 Shafran, Karthik R. Narasimhan, and Yuan Cao. 2023.
769 ReAct: Synergizing reasoning and acting in language
770 models. In *International Conference on Learning*
771 *Representations*.

772 Chaoyun Zhang, Shilin He, Jiaxu Qian, Bowen Li,
773 Liqun Li, Si Qin, Yu Kang, Minghua Ma, Guyue Liu,
774 Qingwei Lin, Saravan Rajmohan, Dongmei Zhang,
775 and Qi Zhang. 2025a. Large language model-brained
776 GUI agents: A survey. *Transactions on Machine*
777 *Learning Research*.

778 Chi Zhang, Zhao Yang, Jiaxuan Liu, Yanda Li, Yucheng
779 Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang
780 Yu. 2025b. AppAgent: Multimodal agents as smart-
781 phone users. In *CHI Conference on Human Factors*
782 *in Computing Systems*, pages 70:1–70:20.

783 Weizhi Zhang, Xinyang Zhang, Chenwei Zhang, Liang-
784 wei Yang, Jingbo Shang, Zhepei Wei, Henry Peng
785 Zou, Zijie Huang, Zhengyang Wang, Yifan Gao, Xi-
786 aoman Pan, Lian Xiong, Jingguo Liu, Philip S. Yu,
787 and Xian Li. 2025c. PersonaAgent: When large lan-
788 guage model agents meet personalization at test time.
789 *arXiv preprint arXiv:2506.06254*.

790 Wanjun Zhong, Lianghong Guo, Qiqi Gao, He Ye, and
791 Yanlin Wang. 2024. MemoryBank: Enhancing large
792 language models with long-term memory. In *AAAI*
793 *Conference on Artificial Intelligence*, pages 19724–
794 19731.

Appendix

A Comparison with Other Benchmarks

Table 4 summarizes the comparison between PSPA-BENCH and existing benchmarks for smartphone GUI agents. General-purpose GUI benchmarks such as AndroidArena, MobileAgentBench, and AndroidWorld are primarily designed to evaluate agents’ basic task-execution capabilities. While they offer large-scale task datasets, they lack support for personalization—there are no user profiles, preferences, or long-term evaluation protocols. SPA-Bench improves realism by leveraging physical devices and third-party apps, but it still does not consider personalization.

Personalized GUI benchmarks, such as LearnGUI and FingerTip, take steps toward user-aware evaluation. However, LearnGUI is limited to simulated environments without physical devices or third-party apps, reducing ecological validity. FingerTip collects real interactions on physical devices with user profiles, but the user preferences are fixed and static, limiting its ability to capture the dynamic and evolving nature of personalization in real-world use.

In contrast, PSPA-BENCH introduces several design choices that address these gaps:

1. Realistic Environment and App diversity.

Unlike many benchmarks that rely on simulated environments or limited app support, PSPA-BENCH is deployed on physical smartphones and spans 22 third-party apps, covering diverse real-world scenarios such as e-commerce, food delivery, and travel. This setting ensures that agents are evaluated under realistic constraints (e.g., UI latency, pop-ups, network variability), increasing ecological validity and practical relevance.

2. Dynamic Modeling of User Preferences.

Most existing personalized benchmarks (e.g., FingerTip) either fix user preferences or omit them entirely. In contrast, PSPA-BENCH models both short-term and long-term user preferences to generate realistic, user-aligned task instructions. Our user profiles and behaviors are not static; instead, they evolve over time, enabling the study of adaptive personalization strategies that reflect real usage patterns.

3. Multi-Dimensional Task Difficulty.

While previous benchmarks typically define difficulty via task execution length, PSPA-BENCH introduces a richer difficulty model. We consider multiple dimensions, including both task execution length and instruction ambiguity. This enables a more nuanced evaluation of an agent’s ability to reason about and adapt to varying personalized tasks.

4. Fine-Grained and Incremental Personalized Evaluation Metrics.

Existing evaluation metrics such as success signal or task completion ratio are often coarse-grained and overly reliant on fixed action sequences. PSPA-BENCH introduces process-oriented evaluation via task decomposition graphs, enabling step-by-step tracking of progress toward user-specific goals. Furthermore, we support incremental evaluation metrics to evaluate the performance changes of agents in long-term personalized use, which is missing in existing benchmarks.

To summarize, PSPA-BENCH combines realistic environments, dynamic user modeling, and fine-grained evaluation to address key gaps in existing benchmarks for personalized GUI agents.

B Definition of Personalized GUI Agent Task

Let Q be the space of GUI task instructions, $E = (\mathcal{S}, \mathcal{A})$ be a GUI environment that contains a set $\mathcal{S} = \{s_i\}$ of observable interface states and a set $\mathcal{A} = \{a_i\}$ of executable GUI actions for state transitions. The definition of Personalized GUI Agent Task is given as follows:

Definition B.1 Personalized GUI Agent Task. For a user u , the profile ρ_u induces a conditional distribution $\mathcal{D}_Q^{\rho_u}$ over the instruction space Q . At each time step t , a personalized instruction q_t^u is sampled from $\mathcal{D}_Q^{\rho_u}$. The agent operates in the GUI environment $E = (\mathcal{S}, \mathcal{A})$ under a user-specific policy π_u with two objectives: (i) **Immediate objective:** generate an action sequence that drives the environment E toward the target states satisfying the current instruction q_t^u ; (ii) **Long-term objective:** accumulate execution experience \mathcal{H}_u to update π_u , thereby improving expected performance on future personalized instruction sequences sampled from $\mathcal{D}_Q^{\rho_u}$.

Table 4: Comparison of General and Personalized GUI Benchmarks

Type	Dataset & Benchmark	Environment with physical phones	Third-part apps	User profiles	User preferences	Task instructions	Difficulty dimensions	Fine-grained progress evaluation	Incremental evaluation
General GUI benchmark	AndroidArena	✗	✗	✗	✗	221	✗	✓	✗
	LlamaTouch	✗	✓	✗	✗	495	Step length	✗	✗
	MobileAgentBench	✗	✗	✗	✗	100	Step length	✗	✗
	AndroidWorld SPA-bench	✗	✗	✗	✗	∞	Step length	✗	✗
Personalized GUI benchmark	LearnGUI	✗	✗	✗	✗	2353	✗	✗	✗
	FingerTip	✓	✓	91	Fixed	21437	Step length	✗	✗
	PSPA-BENCH	✓	✓	100	Dynamic	12855	Step length & instruction clarity	✓	✓

C Construction Task Decomposition Graph

We illustrate the constructed task decomposition graph with four personalized scenarios: shopping, dining, navigation, and travel. In these task decomposition graphs, the nodes with green color refer to fixed type, while nodes with blue color refers to personalized type.

D Instruction Instantiation

Once a template is selected, we instantiate it by filling each slot with concrete values drawn from the user’s preference distribution. Given a user profile ρ , which contains a set of preferences, each slot l is filled by sampling a value $v_l \sim D_{\tau(l)}^{\rho}$, where D_{τ}^{ρ} is a distribution over admissible values conditioned on ρ . The instantiation function $f(\text{TPG}, \{v_l\}_{l \in \mathcal{L}}) \mapsto I$ produces a concrete instruction I by replacing each slot with the sampled value, while ensuring that the chosen values respect the execution dependencies in \mathcal{G} . In this framework, We model user preferences as a combination of: (i) *Long-term preferences*, capturing stable patterns over time (e.g., favoring eco-friendly brands), with weight $w_{\text{long}}(t)$ modeled by a logarithmic growth function. (ii) *Short-term preferences*, capturing transient or contextual interests (e.g., recent searches or trends), with weight $w_{\text{short}}(t)$ modeled by a sinusoidal function. At each time step t , we compute the probability of sampling from each component as: $P_{\text{long}}(t) = w_{\text{long}}(t) / [w_{\text{long}}(t) + w_{\text{short}}(t)]$, and $P_{\text{short}}(t) = w_{\text{short}}(t) / [w_{\text{long}}(t) + w_{\text{short}}(t)]$, respectively.

These probabilities guide the instantiation process for each slot. For example, the `<product_category>` slot might be filled with “organic snacks” based on long-term preferences, or “Halloween costumes” based on short-term seasonal trends. This dynamic sampling ensures that

the personalized instructions reflect both enduring habits and recent behaviors, rather than being static or randomly constructed. The detailed formulation of preference weights, along with experiments evaluating the quality of the generated instructions, is provided in Appendix D.1 and D.2, respectively.

D.1 Details of Preference Weights Calculation

We now describe the preference weights used to model long-term and short-term user patterns. These formulations balance tractability with behavioral plausibility, so that generated instructions reflect both stable needs and dynamic fluctuations.

Long-Term Preference weights. Long-term preferences capture stable, gradually accumulating needs that reset upon satisfaction (e.g., recurring but not continuous interests). Let

$$t \in \mathbb{Z}_{\geq 0}, \quad L \in \mathbb{Z}_{\geq 0}, \quad C \in \mathbb{Z}_{\geq 1}, \quad (2)$$

denote the current day, the last day of satisfaction, and the cycle length in days, respectively. With a base weight $b \in (0, 1]$ and maximum allowed weight $W_{\text{max}} \in (0, 1]$, we define

$$\Delta(t) = \max\{0, t - L\}, \quad \tau(t) = \Delta(t) \bmod C. \quad (3)$$

The long-term preference weight is given by:

$$w_{\text{long}}(t) = \min \left\{ W_{\text{max}}, b + \frac{\ln(1 + \tau(t))(1 - b)}{\ln(1 + C)} \right\}. \quad (4)$$

The logarithmic form reflects the intuition that the *marginal increase* in preference weight diminishes over time: interest accumulates steadily but slows as it approaches saturation. Resetting $\tau(t)$ upon satisfaction models the release of long-term need after fulfillment. The normalization by $\ln(1 + C)$ guarantees that the weight grows smoothly to a normalized scale within each cycle, while clamping ensures that weights remain

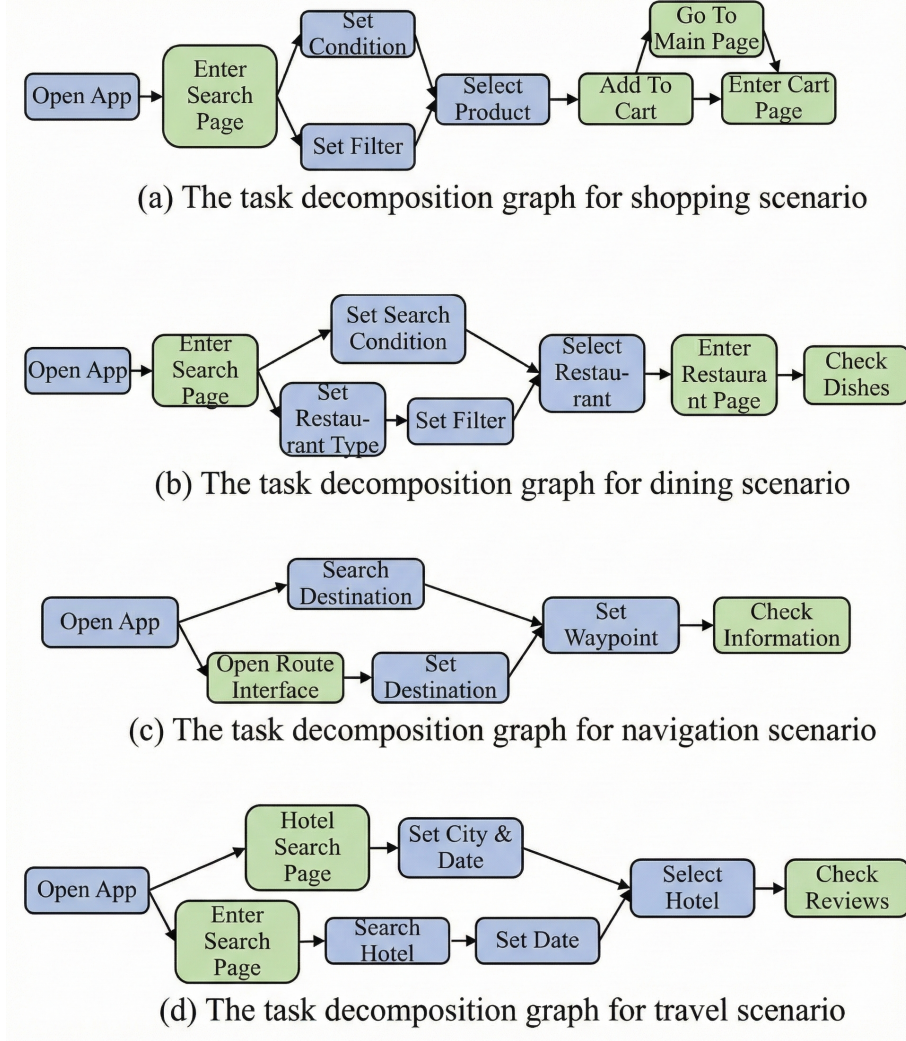


Figure 5: The task decomposition graphs of 4 personalized scenarios: shopping, dining, navigation, and travel. The nodes with green color refer to fixed type, while nodes with blue color refers to personalized type.

963 bounded by W_{\max} . This choice is particularly suitable for modeling durable interests such as dietary
 964 habits, weekly routines, or recurring educational
 965 needs, which build up gradually until satisfied.
 966

967 **Short-Term Preference weights.** Short-term
 968 preferences capture immediate, oscillatory needs
 969 that fluctuate with context. Let

$$970 \quad t \in \mathbb{Z}_{\geq 0}, \quad C \in \mathbb{Z}_{\geq 1}, \quad \phi \in \mathbb{R}, \quad (5)$$

971 denote the current day, the cycle length in days,
 972 and a phase offset, respectively. With bounds $0 <$
 973 $w_{\min} < w_{\max} \leq 1$, we define:

$$974 \quad w_{\text{short}}(t) = w_{\min} + (w_{\max} - w_{\min}) \cdot \frac{\sin\left(\frac{2\pi}{C}t + \phi\right) + 1}{2}. \quad (6)$$

975 The sinusoidal structure is chosen for its natural
 976 ability to model periodic fluctuations. It ensures
 977 smooth transitions between high and low

978 preference intensities, while the parameters w_{\min}
 979 and w_{\max} bound the oscillation range. This reflects
 980 short-term needs such as daily mood, time-of-day
 981 activity cycles, or context-specific demands
 982 (e.g., preferring entertainment in the evening versus
 983 productivity tools in the morning). Compared to
 984 stochastic fluctuations, sinusoidal oscillation
 985 avoids unnatural abrupt shifts and instead captures
 986 the rhythmic, recurrent nature of short-term
 987 preferences.

988 **Normalized Sampling probabilities.** Given
 989 both weights, we define the probabilities of drawing
 990 from long-term and short-term preferences as:

$$991 \quad P_{\text{long}}(t) = \frac{w_{\text{long}}(t)}{w_{\text{long}}(t) + w_{\text{short}}(t)}, \quad (7)$$

$$992 \quad P_{\text{short}}(t) = \frac{w_{\text{short}}(t)}{w_{\text{long}}(t) + w_{\text{short}}(t)}. \quad (8)$$

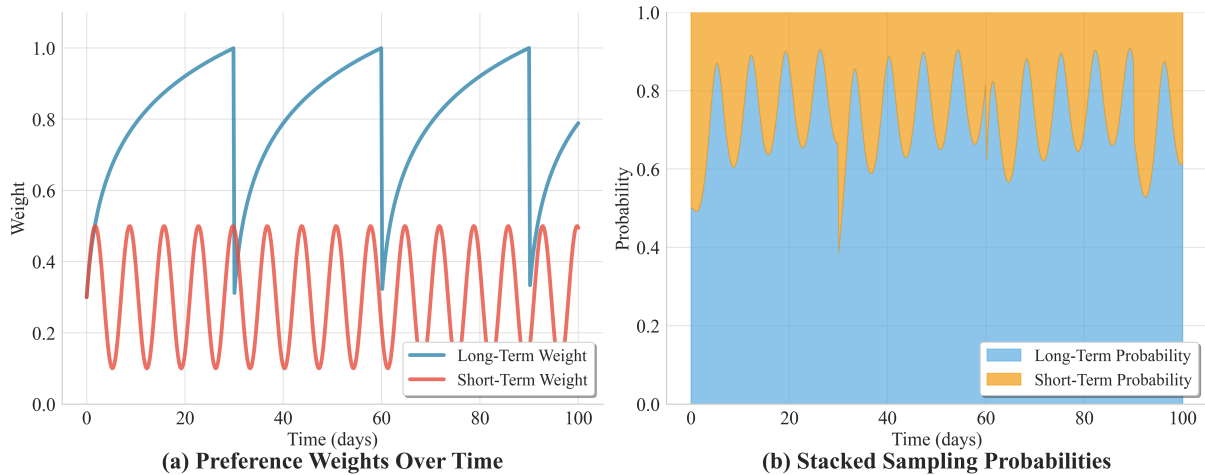


Figure 6: (a) **Preference Weights Over Time**: The blue curve represents the long-term preference weight, which gradually increases with time, reflecting the accumulation of stable needs. The red curve depicts the short-term preference weight, which oscillates periodically, reflecting dynamic, context-dependent needs. (b) **Stacked Sampling Probabilities**: The stacked area plot shows the normalized sampling probabilities over time. The blue region corresponds to the probability of sampling the long-term preference, while the orange region corresponds to the probability of sampling the short-term preference. The relative areas of these regions change over time, depending on the weights of each preference.

This normalization yields a valid probability distribution and balances stability with variability: when long-term needs are strong, they dominate; when short-term fluctuations peak, they take precedence. A visual illustration of these dynamics can be found in Figure 6.

D.2 Instruction Quality Evaluation

To rigorously evaluate the quality of generated task instructions, we conducted a human evaluation study across **10 distinct scenarios** involving a total of **100 users**. For each user, we selected the **first 10 generated instructions** from each method and used them as the evaluation material. This design ensures a fair and comparable assessment across methods, while also reflecting the diversity of tasks encountered in practical settings.

We compared three approaches for generating user instructions:

1. **Static Preference.** Each user is assigned a *fixed preference profile* that does not change over time. Instructions are generated by filling task templates strictly according to this static preference. While this captures general user tendencies, it cannot adapt to temporal shifts in user needs.
2. **Random Preference.** At each time step, a *random preference* is assigned to the user without considering historical data. Instructions

are then generated from this random preference. Although this produces variety, the lack of temporal or contextual grounding leads to incoherent or irrelevant instructions.

3. **Dynamic Preference Modeling (Ours).** We model user preferences as a distribution combining both *long-term patterns* (stable, accumulated needs) and *short-term patterns* (immediate, oscillating needs). At each time step t , the probability of sampling long-term versus short-term preferences is determined by their respective weights. Task instructions are generated based on the sampled preferences, enabling outputs that capture both the *stable characteristics* of the user and their *short-term dynamics*.

Annotators were presented with **pairs of instruction sets** (10 instructions per set, per user). Each pair was evaluated **blindly**, comparing two methods at a time. Annotators judged which set was better overall, based on two criteria:

- **Plausibility / Realism:** Are the instructions natural and reasonable, resembling tasks that a real user might genuinely perform?
- **Relevance to User Context:** Are the instructions aligned with the user’s scenario, needs, and historical behavior?

1048 Annotators were allowed to declare a **Tie** if the
1049 two sets were of equal quality. Each pair was eval-
1050 uated by multiple annotators, and final results were
1051 aggregated using majority voting.

1052 As shown in Table 5, the results clearly demon-
1053 strate the advantage of **dynamic preference mod-**
1054 **eling**:

- 1055 • **Ours vs. Static:** Our method achieved a sig-
1056 nificantly higher win-rate, confirming that ac-
1057 counting for temporal dynamics provides su-
1058 perior alignment with user preferences com-
1059 pared to fixed profiles.
- 1060 • **Ours vs. Random:** The large margin shows
1061 that random preferences fail to capture user
1062 needs, leading to inconsistent instruction qual-
1063 ity.
- 1064 • **Static vs. Random:** While Static outperforms
1065 Random to some extent, it still lags far behind
1066 our method due to its inability to adapt to
1067 evolving preferences.

1068 The **low tie percentage** across all comparisons
1069 indicates that evaluators could generally distin-
1070 guish between methods, supporting the reliability
1071 of the results.

1072 E Evaluation Details

1073 Table 6 organizes the evaluation metrics for
1074 the Personalized GUI Agent Task into two key
1075 dimensions—**performance** and **efficiency**—each
1076 assessed through both immediate and long-term
1077 objectives. Immediate metrics capture task-level
1078 outcomes, including the **A-process Ratio (APR)**,
1079 measuring the proportion of completed unit instruc-
1080 tions, and the **P-process Ratio (PPR)**, restricted to
1081 flexible preference-sensitive nodes. Efficiency is
1082 evaluated via **Completion Time (CT)**, the average
1083 execution time per task, and **Cost per Task (CPT)**,
1084 the average token-based monetary cost. Long-term
1085 metrics (ΔAPR , ΔPPR , ΔCT , ΔCPT) reflect im-
1086 provements across task sequences, defined as the
1087 difference between performance at a later task T_k
1088 and the initial task T_1 . This framework jointly
1089 measures execution quality, personalization, time,
1090 and cost, providing a comprehensive view of both
1091 short-term performance and long-term adaptability.

1092 E.1 Metric Comparison

1093 **Experiment Setup.** To evaluate the effectiveness
1094 of different metrics, we sampled 25 personalized

1095 GUI tasks from PSPA-BENCH and used **M3A** to
1096 execute each task and obtain the execution trajec-
1097 tory. For every task, we then measured the perfor-
1098 mance based on the trajectory using four metrics:
1099 **Success Signal**, **Task Completion Ratio (TCR)**,
1100 **A-process Ratio (APR)**, and **Human Ratings** (col-
1101 lected from third-party annotators). Success Signal
1102 provides a binary success/failure outcome, TCR
1103 measures the proportion of actions completed along
1104 a predefined reference path, APR reflects the aver-
1105 age completion rate across unit instructions along
1106 the execution path, and annotators rate the comple-
1107 tion ratio (human ratings) of each execution trajec-
1108 tory, which serves as the ground truth.

1109 **Results and Analysis.** Figure 7 compares the
1110 four metrics across tasks (sorted by human ratings).
1111 Success Signal is highly coarse-grained, producing
1112 only binary outcomes (0 or 1). This representation
1113 collapses diverse execution trajectories into two cat-
1114 egories, obscuring important differences between
1115 agents that achieve partial progress versus those
1116 that fail entirely. TCR provides more variation but
1117 is still tightly coupled to a fixed path, which leads to
1118 penalizing valid alternative behaviors. As a result,
1119 TCR diverges notably from human ratings, with
1120 weaker correlation ($\rho = 0.47$) and higher error
1121 ($\text{MSE}=0.062$).

1122 By contrast, APR demonstrates strong consis-
1123 tency with human judgments. As shown in Fig-
1124 ure 7, APR closely follows human ratings across
1125 all tasks, with almost all values lying within the
1126 human-acceptable tolerance band (± 0.05). Quan-
1127 titatively, APR achieves a high correlation with
1128 human ratings ($\rho = 0.95$) and near-zero error
1129 ($\text{MSE}=0.005$). These results show that APR cap-
1130 tures fine-grained progress, tolerates diverse execu-
1131 tion strategies, and aligns well with human evalua-
1132 tors. Compared to Success Signal and TCR, APR
1133 offers a more reliable assessment of agent perfor-
1134 mance in personalized, long-horizon GUI tasks.

1135 E.2 Metrics Calculation

1136 While the metrics above provide a fine-grained
1137 view of agent behavior, computing them reliably
1138 is non-trivial. (Xing et al., 2024) measure task
1139 completion by comparing executed actions against
1140 a fixed ground-truth sequence. Although useful,
1141 this method is tied to a single action path and can-
1142 not capture an agent’s ability to follow alternative
1143 strategies or adapt to user preferences. To address
1144 this limitation, we leverage the TDG as an inter-
1145 mediate representation and convert it into a struc-

Table 5: Pairwise win-rates of different methods in human evaluation.

Comparison (Pairwise)	Win-rate Ours \uparrow	Win-rate Baseline	Tie (%)
Ours vs. Static	72.5%	24.1%	3.4%
Ours vs. Random	81.3%	16.2%	2.5%
Static vs. Random	59.6%	38.7%	1.7%

Table 6: Evaluation metrics for the Personalized GUI Agent Task. Short-term metrics measure absolute performance on individual tasks, while long-term metrics capture improvements across instruction sequences. N denotes the number of tasks.

Dimension	Immediate Objective	Long-term Objective
Performance	APR: $\frac{1}{N} \sum \frac{\text{completed unit inst.}}{\text{total unit inst.}}$	$\Delta\text{APR: APR}_{T_k} - \text{APR}_{T_1}$
	PPR: $\frac{1}{N} \sum \frac{\text{completed flexible unit inst.}}{\text{total flexible unit inst.}}$	$\Delta\text{PPR: PPR}_{T_k} - \text{PPR}_{T_1}$
Efficiency	CT: $\frac{1}{N} \sum t_i$	$\Delta\text{CT: CT}_{T_k} - \text{CT}_{T_1}$
	CPT: $\frac{1}{N} \sum \text{Cost}_i$	$\Delta\text{CPT: CPT}_{T_k} - \text{CPT}_{T_1}$

1146 tured checklist. The checklist encodes execution
 1147 dependencies, alternative paths, and the distinction
 1148 between fixed and flexible unit instructions. We
 1149 then employ a LLM such as GPT-5 as an evaluator:
 1150 given the agent’s execution trace and the check-
 1151 list, the model determines which unit instructions
 1152 were successfully completed. This graph-driven
 1153 approach allows APR and PPR to be computed
 1154 consistently, even when tasks admit multiple valid
 1155 solutions or require personalization. To ensure the
 1156 validity and reproducibility of this LLM-based eval-
 1157 uation, we conduct a comprehensive human vali-
 1158 dation study, as detailed in Section E.3. Details of
 1159 checklist template are provided in Appendix E.4.

1160 **Path Selection Algorithm.** As described in
 1161 Equation 1, selecting the optimal path from the
 1162 TDG is crucial for accurate metric computation.
 1163 Algorithm 1 formalizes our two-stage alignment
 1164 procedure:

1165 The tie-breaking procedure (Lines 17–18) en-
 1166 sures deterministic path selection: we first prefer
 1167 paths with more matched flexible nodes (prioritiz-
 1168 ing personalization), then shorter paths (favoring
 1169 efficiency), and finally paths whose first matched
 1170 node appears earliest in the trace (reflecting pri-
 1171 mary intent). This hierarchical criterion guarantees
 1172 reproducible APR/PPR computation across differ-
 1173 ent evaluation runs.

Algorithm 1 Trace-Graph Alignment with Path Selection

Require: Trace $\mathcal{T} = (a_1, \dots, a_L)$, TDG \mathcal{G} with valid paths \mathcal{P}
Ensure: Optimal path P^* , alignment mapping M
 1: **Stage 1: Action-to-Instruction Alignment**
 2: **for** each action $a_i \in \mathcal{T}$ **do**
 3: $m(a_i) \leftarrow \text{LLM_Align}(a_i, \mathcal{U})$ {Map action to unit in-
 struction or \emptyset }
 4: **end for**
 5: **Stage 2: Path Selection**
 6: $\text{best_score} \leftarrow -1$; $\text{candidates} \leftarrow \emptyset$
 7: **for** each path $P \in \mathcal{P}$ **do**
 8: $\text{score}(P) \leftarrow |\{u \in P : \exists a_i, m(a_i) = u\}|$ {Count
 matched nodes}
 9: **if** $\text{score}(P) > \text{best_score}$ **then**
 10: $\text{best_score} \leftarrow \text{score}(P)$; $\text{candidates} \leftarrow \{P\}$
 11: **else if** $\text{score}(P) = \text{best_score}$ **then**
 12: $\text{candidates} \leftarrow \text{candidates} \cup \{P\}$
 13: **end if**
 14: **end for**
 15: **Stage 3: Tie-Breaking** (if $|\text{candidates}| > 1$)
 16: $P^* \leftarrow \text{TieBreak}(\text{candidates})$ {Apply rules: flexible \rightarrow
 length \rightarrow earliest}
 17: **return** $P^*, \{(a_i, m(a_i)) : m(a_i) \in P^*\}$

E.3 Human Validation of LLM-based Alignment

1174 To validate the reliability of our LLM-based trace-
 1175 graph alignment, we conduct a comprehensive hu-
 1176 man agreement study. 1177 1178

1179 **Annotation Protocol.** We randomly sampled
 1180 100 execution traces from our benchmark, covering
 1181 all four scenarios (Shopping, Dining, Navigation,
 1182 and Travel) with balanced complexity levels. Three
 1183 trained annotators independently labeled each trace
 1184 by marking which unit instructions in the corre-

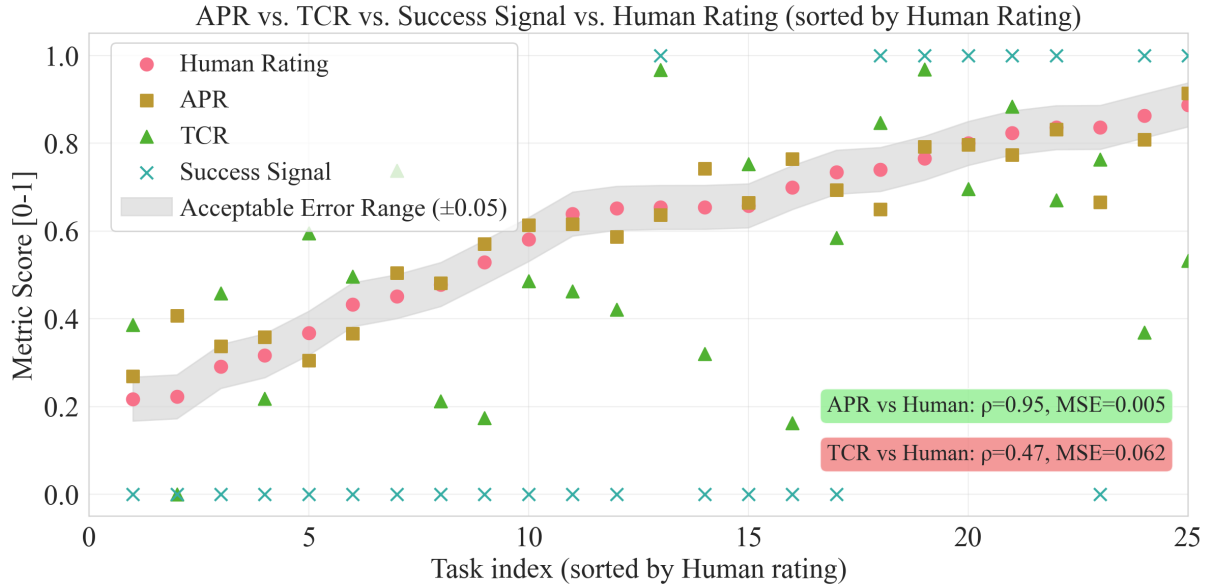


Figure 7: Comparison of evaluation metrics across tasks sorted by human ratings. Each point represents one task. Human ratings (pink dots) are shown with an acceptable tolerance band (± 0.05 , gray area). A-process Ratio (APR, brown squares) closely follows human ratings with minimal deviation, demonstrating strong alignment ($\rho=0.95$, $MSE=0.005$). In contrast, Task Completion Ratio (TCR, green triangles) shows weaker consistency with human ratings ($\rho=0.47$, $MSE=0.062$), while Success Signal (blue crosses) provides only binary outcomes, lacking granularity.

1185 sponding TDG were successfully completed. An-
 1186 notators were provided with the execution trace
 1187 (screenshots and action logs), the TDG checklist,
 1188 and detailed annotation guidelines. Prior to the
 1189 main annotation, annotators completed a calibra-
 1190 tion phase on 20 pilot traces to ensure consistent
 1191 interpretation of completion criteria.

1192 **Agreement Analysis.** We measure inter-
 1193 annotator agreement using Fleiss’ Kappa (κ),
 1194 which accounts for chance agreement among
 1195 multiple raters. The overall inter-annotator agree-
 1196 ment reached $\kappa = 0.87$, indicating *almost perfect*
 1197 *agreement* according to standard interpretation
 1198 guidelines (Landis and Koch, 1977). For cases
 1199 with disagreement, we resolved conflicts through
 1200 majority voting to establish the ground-truth labels.

1201 **LLM Alignment Accuracy.** We then compared
 1202 the LLM-based alignment results against the
 1203 human-annotated ground truth. Table 7 summa-
 1204 rizes the results across different instruction types:

1205 The LLM-based alignment achieves an overall
 1206 accuracy of 93% and Cohen’s $\kappa = 0.86$ with hu-
 1207 man annotations, demonstrating strong agreement.
 1208 Fixed instructions show slightly higher accuracy
 1209 (94%) than flexible instructions (91%), which is ex-
 1210 pected given that flexible instructions involve more

Table 7: Validation of LLM-based alignment against human annotations. We report Accuracy, Precision, Recall, F1-score, and Cohen’s κ for agreement with human ground truth.

Instruction Type	Acc.	Prec.	Rec.	F1	κ
Fixed Instructions	0.94	0.93	0.96	0.94	0.88
Flexible Instructions	0.91	0.90	0.93	0.91	0.83
Overall	0.93	0.92	0.95	0.93	0.86

1211 nuanced judgments about user preference satisfac-
 1212 tion.

1213 **Reproducibility Analysis.** To assess the repro-
 1214 ducibility of the LLM-based alignment, we ran
 1215 the evaluation pipeline five times on the same set
 1216 of 100 traces using identical prompts and model
 1217 settings (temperature = 0). The alignment results
 1218 showed 98.7% consistency across runs, with a stan-
 1219 dard deviation of only 0.8% for the computed APR
 1220 values. This high consistency demonstrates that
 1221 our LLM-based alignment method is reproducible
 1222 when using deterministic inference settings.

1223 **Error Analysis.** We analyzed the 7% of cases
 1224 where LLM alignment disagreed with human an-
 1225 notations. The primary sources of disagreement
 1226 were: (i) ambiguous partial completions where an

1227	instruction was partially but not fully satisfied (42%	• Prompt 2 (LLM-based Evaluation): Given an	1270
1228	of errors), (ii) alternative valid paths not explicitly	agent’s execution trace (as a sequence of screen-	1271
1229	enumerated in the checklist (31%), and (iii) edge	shots) and the corresponding checklist, the evalu-	1272
1230	cases involving UI state changes between screen-	ator LLM determines which unit instructions	1273
1231	screenshots (27%). These insights inform future improve-	were successfully completed. The prompt speci-	1274
1232	ments to both the checklist design and the evaluation	ifies evaluation rules, including handling of alter-	1275
1233	prompts.	native paths and dependency constraints.	1276
1234	E.4 Details of Checklists and Prompt	The complete checklist templates for all four	1277
1235	Templates	scenarios and the detailed prompt templates are	1278
1236	Checklist Structure. We provide four domain-	presented in Checklist 1–4 and Prompt 1–2.	1279
1237	specific checklist templates corresponding to the	F Templates of personalized GUI	1280
1238	four scenarios in our benchmark: Shopping, Din-	instruction	1281
1239	ing, Navigation, and Travel. Each checklist is de-	We provide examples of templates in personalized	1282
1240	derived from the Task Decomposition Graph (TDG)	scenarios. The templates are categorized by <i>clar-</i>	1283
1241	and contains the following structured elements for	<i>ity level</i> (high, medium, low) and <i>complexity level</i>	1284
1242	each unit instruction:	(low, medium, high). In each level, we provide	1285
1243	• Checkbox: A visual marker (□) indicating com-	3 templates with different expressions, which im-	1286
1244	pletion status.	prove the variety of personalized instructions.	1287
1245	• Step Name: The name of the action, annotated	G LLM Usage	1288
1246	with its operation type (<i>Fixed</i> or <i>Flexible</i>).	Large Language Models (LLMs) were used to aid	1289
1247	• Description: A clear explanation of what the	in the writing and polishing of the manuscript.	1290
1248	step accomplishes.	Specifically, we used an LLM to assist in refin-	1291
1249	• Dependency: The prerequisite step(s) that must	ing the language, improving readability, and ensur-	1292
1250	be completed before this step can be executed.	ing clarity in various sections of the paper. The	1293
1251	• Path Options: For flexible steps with multiple	model helped with tasks such as sentence rephras-	1294
1252	valid execution paths, we enumerate them as Path	ing, grammar checking, and enhancing the overall	1295
1253	A, Path B, etc., each with its own sub-steps.	flow of the text.	1296
1254	The distinction between <i>Fixed</i> and <i>Flexible</i> oper-	The authors take full responsibility for the con-	1297
1255	ations is important: Fixed steps are mandatory with	tent of the manuscript, including any text gener-	1298
1256	no alternatives, while Flexible steps allow multiple	ated or polished by the LLM. We have ensured that the	1299
1257	valid approaches or are influenced by user prefer-	LLM-generated text adheres to ethical guidelines	1300
1258	ences. This design lets the checklist accommodate	and does not contribute to plagiarism or scientific	1301
1259	personalized execution strategies while keeping	misconduct.	1302
1260	evaluation criteria consistent.		
1261	Prompt Templates. We design two prompts to		
1262	facilitate checklist-based evaluation:		
1263	• Prompt 1 (Checklist Generation): Converts a vi-		
1264	sual Task Decomposition Graph into a structured		
1265	textual checklist that can be processed by LLMs.		
1266	The prompt instructs the model to extract step		
1267	names, descriptions, dependencies, and operation		
1268	types while preserving the graph’s hierarchical		
1269	structure.		

Checklist 1: Shopping scenario

- Open Taobao Fixed**
Description: Open the Taobao homepage and enter the main platform page
- Enter Search Page Fixed**
Description: Enter the search interface to prepare for product search
Dependency: Must be executed after “Open Taobao”
- Set Conditions (choose one path)**
Description: Set search constraints (either method is acceptable) *Dependency:* Must first complete “Enter Search Page”
 - **Path A:**
 - Set Search Conditions Flexible**
Description: Enter keywords, categories, etc. for searching
 - **Path B:**
 - Set Filter Conditions Flexible**
Description: Set filters such as price range, sales sorting, etc.
- Select Product Flexible**
Description: Select a product from search results
Dependency: Must first complete “Set Conditions” in either path
- Add to Cart Fixed**
Description: Click the “Add to Cart” button
Dependency: Must first complete “Select Product”
- Navigate to Cart Page (choose one path) Flexible**
Description: Enter the shopping cart page (flexible path) *Dependency:* Must first complete “Add To Cart”
 - **Path A:**
 - Enter Cart Page (directly) Fixed**
Description: Directly enter the shopping cart page after adding to cart
 - **Path B:**
 - Return to Home Page Fixed**
Description: First return to the homepage
 - Enter Cart Page (from Home) Fixed**
Description: Enter the shopping cart page from the homepage
Dependency: Must first complete “Return to Home Page”

Checklist 2: Dining scenario

- Open App Flexible**
Description: Open the app and enter its main page
- Enter Search Page Fixed**
Description: Enter the search interface to prepare for product search
Dependency: Must be executed after “Open App”
- Set Conditions (choose one path) Flexible**
Description: Set Searching Conditions (flexible path)
Dependency: Must be executed after “Enter Search Page”
 - **Path A:**
 - Set Search Conditions Flexible**
Description: Directly search for the restaurant using keywords
 - **Path B:**
 - Set Restaurant Type Flexible**
Description: Search for the type or name of the restaurant
 - Set Filter Conditions Flexible**
Description: Set filter conditions such as price, distance, etc.
Dependency: Must first complete “Set Restaurant Type”
- Select Restaurant Flexible**
Description: Select a restaurant from the searching results
Dependency: Must be executed after “Set Searching Conditions”
- Enter Main Page of Restaurant Fixed**
Description: Click to enter the main page of the restaurant
Dependency: Must be executed after “Select Restaurant”
- Check Recommended Dishes Fixed**
Description: Check the restaurant’s recommended dishes
Dependency: Must be executed after “Enter Main Page of Restaurant”

Checklist 3: Navigation scenario

- Open App Flexible**
Description: Open the app and enter its main page
- Set Route To Destination Flexible**
Description: Set the navigation destination (flexible path)
Dependency: Must be executed after “Open App”
 - **Path A:**
 - Search For Destination Flexible**
Description: Directly search for the destination in the searching box
 - **Path B:**
 - Open Route Planning Interface Fixed**
Description: Enter the route planning and destination searching interface
 - Set Destination Flexible**
Description: Set destination in the route planning interface
Dependency: Must first complete “Open Route Planning Interface”
- Set Waypoint Flexible**
Description: Set the waypoint for the route
Dependency: Must be executed after “Set Route To Destination”
- Check Route Information Fixed**
Description: Check the information of the route
Dependency: Must be executed after “Set Waypoint”

Checklist 4: Travel scenario

- Open App Flexible**
Description: Open the App homepage and enter the main platform page
- Search For Hotel (choose one path)**
Description: Set search constraints (either method is acceptable) *Dependency:* Must first complete “Open App”
 - **Path A:**
 - Enter Hotel Search Page Fixed**
Description: Enter the search category for hotel
 - Set City And Date Flexible**
Description: Set the searching conditions like city, type, and time *Dependency:* Must first complete “Enter Hotel Search Page”
 - **Path B:**
 - Enter Search Page Fixed**
Description: Enter the search page of the APP
 - Search For Hotel In Target City Flexible**
Description: Set searching conditions like city and hotel type *Dependency:* Must first complete “Enter Search Page”
 - Set Date Fixed**
Description: Set dates of the stay *Dependency:* Must first complete “Set Date”
- Select Hotel Flexible**
Description: Select a hotel from the results *Dependency:* Must first complete “Search For Hotel”
- Check Reviews Fixed**
Description: Check reviews of the hotel on the APP *Dependency:* Must first complete “Select Hotel”

Prompt 1: Prompt for turning task decomposition graphs into checklist:

Your task is to turn a task decomposition graph in a picture into textual checklist that can be understood by another large language model easily. The task decomposition graph is provided below:

{TASK_GRAPH}

The checklist you generate should include:

- A checkbox
- Step Name: The name of the action from the graph.
- Description :A clear explanation of what the step does.
- Dependency :The prerequisite step(s) that must be completed before this step.
- Operation Type: Mark as Fixed if the step is mandatory with no alternatives, or Flexible if there are multiple paths/options, or the step is affected by user preference. (including grouped parallel steps).
- Path Options (if applicable): If the step has multiple paths (e.g., parallel steps leading to the same next step), list them as Path A, Path B, etc., with a description with the same format as the steps you generate.

An example of the checklist you generate is provided below, make sure you follow the correct format.

Example:

- Set Route To Destination Flexible**
Description: Set the navigation destination (flexible path)
Dependency: Must be executed after “Open App”
 - **Path A:**
 - Search For Destination Flexible**
Description: Directly search for the destination in the searching box
 - **Path B:**
 - Open Route Planning Interface Fixed**
Description: Enter the route planning and destination searching interface
 - Set Destination Flexible**
Description: Set destination in the route planning interface
Dependency: Must first complete “Open Route Planning Interface”

Prompt 2: Prompt for LLM-based evaluation:

You are an expert in evaluating the performance of a GUI Agent. Based on the input checklist and the agent’s execution flow images, you will determine one by one whether the agent has completed all sub-tasks in the checklist. The checklist will be provided to you in JSON format, while the images will be base64 encoded. Your output should also be a JSON, where you set the complete marker for each sub-task in the input JSON to 1 if the sub-task is completed and 0 if not completed.

{IMAGES}
{CHECKLIST}

Rules: 1. Check the subtasks and determine whether they are completed one by one. 2. For subtask with various paths, the completion of any path marks the completion of the subtask. 3. If the dependency of a subtask is not completed, the subtask itself is not completed as well.

Table 8: Instruction templates of shopping scenario.

Clarity level	Complexity level	Template Examples (English)
High clarity	Low	<ol style="list-style-type: none"> 1. I want to search for {brand} {category} on {APP}, with prices between {price range} yuan. 2. Help me find a {brand} {category} on {APP}, with prices in the range of {price range} yuan. 3. Check {APP} for a {brand} {category}, keeping the price within {price range} yuan.
	Medium	<ol style="list-style-type: none"> 1. I want to buy a {brand} {category} on {APP}, with prices between {price range} yuan, and stop at the product details page. 2. Help me find a {brand} {category} on {APP}, with prices in the range of {price range} yuan, and stay on the product details page. 3. Check {APP} for a {brand} {category}, with prices around {price range} yuan, and stay on the product details page.
	High	<ol style="list-style-type: none"> 1. I want to buy a {brand} {category} on {APP}, with prices between {price range} yuan, add the item to the cart, and finally stay on the cart page. 2. Help me find a {brand} {category} on {APP}, with prices around {price range} yuan, add it to the cart, and finally stay on the cart page. 3. Check {APP} for a {brand} {category}, within {price range} yuan, add it to the cart, and finally stay on the cart page.
Medium clarity	Low	<ol style="list-style-type: none"> 1. Search {brand} {category} on {APP}, price about {price range}. 2. Help me find {brand} {category} in {APP}, budget {price range}. 3. Check {APP} for {brand} {category}, price around {price range}.
	Medium	<ol style="list-style-type: none"> 1. Want to buy {brand} {category} on {APP}, price {price range}, just open the details page. 2. {APP}, find {brand} {category}, around {price range}, go to the details page. 3. Help me search {brand} {category} on {APP}, price {price range}, open the details page.
	High	<ol style="list-style-type: none"> 1. Buy {brand} {category} on {APP}, budget {price range}, add to cart and stay there. 2. {APP}, find {brand} {category}, price about {price range}, add to cart. 3. Help me add {brand} {category} ({price range}) in {APP} to the cart, and stay on the cart page.
Low clarity	Low	<ol style="list-style-type: none"> 1. {APP} {brand}{category} {price range} 2. Find {brand}{category} {price range} @ {APP} 3. {brand}{category} {price range} / {APP}
	Medium	<ol style="list-style-type: none"> 1. {APP} {brand}{category} {price range} details 2. {brand}{category} {price range} → details @ {APP} 3. {APP}: {brand}{category} {price range} details page
	High	<ol style="list-style-type: none"> 1. {APP} {brand}{category} {price range} cart 2. Add to cart {brand}{category} {price range} @ {APP} 3. {brand}{category} {price range} → cart / {APP}

Table 9: Instruction templates for music scenario.

Clarity level	Complexity level	Template Examples (English)
High clarity	Low	<ol style="list-style-type: none"> 1. I want to search for {singer} on {APP} and go to the singer's homepage. 2. Help me search for {singer} on {APP} and enter the singer's homepage. 3. Check {APP} for {singer} and click to open the singer's homepage.
	Medium	<ol style="list-style-type: none"> 1. I want to search for {singer} on {APP}, go to the singer's homepage, and view the albums. 2. Help me find {singer} on {APP}, enter the singer's homepage, and check the albums. 3. On {APP}, search for {singer}, open the singer's homepage, then view the albums.
	High	<ol style="list-style-type: none"> 1. I want to search for {singer} on {APP}, open the singer's homepage, check the albums, and select a song to play. 2. Help me find {singer} on {APP}, enter the singer's homepage, view the albums, and then play a song. 3. On {APP}, search for {singer}, open the homepage, view the albums, and finally play a song.
Medium clarity	Low	<ol style="list-style-type: none"> 1. Search {singer} on {APP}, go to homepage. 2. Help me find {singer} in {APP}, check homepage. 3. {APP} search {singer}, open homepage.
	Medium	<ol style="list-style-type: none"> 1. Want to find {singer} on {APP}, check homepage and albums. 2. {APP} search {singer}, enter homepage, view albums. 3. Help me search for {singer} on {APP}, view albums in homepage.
	High	<ol style="list-style-type: none"> 1. Find {singer} on {APP}, view albums, select a song to play. 2. {APP} search {singer}, choose a song from album to listen. 3. Help me find {singer} on {APP}, select a song from album and play.
Low clarity	Low	<ol style="list-style-type: none"> 1. {APP} {singer} homepage 2. Find {singer} @ {APP} 3. {singer} homepage / {APP}
	Medium	<ol style="list-style-type: none"> 1. {APP} {singer} album 2. {singer} album @ {APP} 3. {APP}: {singer} album page
	High	<ol style="list-style-type: none"> 1. {APP} {singer} play album song 2. Play {singer} album song @ {APP} 3. {singer} → play album song / {APP}

Table 10: Instruction templates for dining scenario.

Clarity level	Complexity level	Template Examples (English)
High clarity	Low	<ol style="list-style-type: none"> 1. I want to search for {cuisine type} restaurants on {APP}, with an average spending between {price range} yuan. 2. Help me find a {cuisine type} restaurant on {APP}, with average spending in the range of {price range} yuan. 3. Check {APP} for {cuisine type} restaurants, priced between {price range} yuan.
	Medium	<ol style="list-style-type: none"> 1. I want to search for {cuisine type} restaurants on {APP}, with an average spending between {price range} yuan, and stay on the restaurant details page. 2. Help me find a {cuisine type} restaurant on {APP}, with average spending in the range of {price range} yuan, and remain on the details page. 3. Check {APP} for {cuisine type} restaurants, priced between {price range} yuan, and stay on the details page.
	High	<ol style="list-style-type: none"> 1. I want to search for {cuisine type} restaurants on {APP}, with an average spending between {price range} yuan, enter the details page, and view all recommended dishes. 2. Help me find a {cuisine type} restaurant on {APP}, with average spending in the range of {price range} yuan, go to the details page, and view all recommended dishes. 3. Check {APP} for {cuisine type} restaurants, priced between {price range} yuan, enter the details page, and view all recommended dishes.
Medium clarity	Low	<ol style="list-style-type: none"> 1. Search {cuisine type} on {APP}, around {price range}. 2. Help me find {cuisine type} in {APP}, budget {price range}. 3. Check {APP} for {cuisine type}, price about {price range}.
	Medium	<ol style="list-style-type: none"> 1. Want to find {cuisine type} on {APP}, price between {price range}, open details page. 2. {APP} search {cuisine type}, in range {price range}, open details. 3. Help me search {cuisine type} on {APP}, price {price range}, view details.
	High	<ol style="list-style-type: none"> 1. Find {cuisine type} on {APP}, budget {price range}, check recommended dishes. 2. {APP} search {cuisine type}, price {price range}, view recommended dishes. 3. Help me find {cuisine type} ({price range}) on {APP}, view recommended dishes.
Low clarity	Low	<ol style="list-style-type: none"> 1. {APP} {cuisine type} {price range} 2. Find {cuisine type} {price range} @ {APP} 3. {cuisine type} {price range} / {APP}
	Medium	<ol style="list-style-type: none"> 1. {APP} {cuisine type} {price range} details 2. {cuisine type} {price range} → details @ {APP} 3. {APP}: {cuisine type} {price range} details page
	High	<ol style="list-style-type: none"> 1. {APP} {cuisine type} {price range} view recommended dishes 2. Recommended dishes {cuisine type} {price range} @ {APP} 3. {cuisine type} {price range} → view recommended dishes / {APP}