

# How Does Personalized Memory Shape LLM Behavior? Benchmarking Rational Preference Utilization in Personalized Assistants

Anonymous ACL submission

## Abstract

Large language model (LLM)-powered assistants have recently integrated memory mechanisms that record user preferences, leading to more personalized and user-aligned responses. However, irrelevant personalized memories are often introduced into the context, interfering with the LLM’s intent understanding. To comprehensively investigate the dual effects of personalization, we develop RPEVAL, a benchmark comprising a personalized intent reasoning dataset and a multi-granularity evaluation protocol. RPEVAL reveals the widespread phenomenon of irrational personalization in existing LLMs and, through error pattern analysis, illustrates its negative impact on user experience. Finally, we introduce RP-REASONER, which treats memory utilization as a pragmatic reasoning process, enabling the selective integration of personalized information. Experimental results demonstrate that our method significantly outperforms carefully designed baselines on RPEVAL, and resolves 80% of the bad cases observed in a large-scale commercial personalized assistant, highlighting the potential of pragmatic reasoning to mitigate irrational personalization. Our benchmark is publicly available at <https://anonymous.4open.science/r/RPEval-E4B0>.

## 1 Introduction

In human communication, people tend to express themselves as economically as possible—using minimal language to convey maximal intent. This often results in utterances that are heavily underspecified, relying on the listener to fill in the gaps through shared knowledge and mutual understanding (Grice, 1975). This fundamental cognitive mechanism underlies the increasing appeal of personalized assistants (PAs) (Zhang et al., 2024a; Zhao et al., 2025) powered by large language models (LLMs) (OpenAI, 2023). By continuously interacting with users, PAs incrementally build personalized memory stores that capture user-specific

information. When handling new queries, they read relevant memories and incorporate them into the context, enabling the generation of personalized response (Zhang et al., 2024c), thereby enhancing user experience and satisfaction.

This work centers on the duality of personalization, particularly its adverse effects. As shown in Figure 1, when a user requests sleep-aid audio, the memory module may activate a preference for strong rhythm music due to semantic similarity. This can mislead the LLM into recommending fast tracks, which clearly conflicts with the user’s current needs. This issue is particularly pronounced in real-world scenarios: (1) **memories are often sparse and fragmented**, leading to oversimplified or overly labeled modeling of users (He et al., 2017); (2) **user queries are typically open-ended**, frequently involving requests unrelated to previously stored information (Chaney et al., 2018). As a result, when generating personalized responses, PAs may inevitably introduce irrelevant memories into the context, causing LLMs to misuse memories and deviate from user intent. Inspired by Rational Speech Acts theory (Frank and Goodman, 2012), we propose **Rational Personalization** for LLMs: a problem where an LLM must evaluate the recalled user preferences with varying degrees of applicability, and decide whether and how to integrate them to accurately infer and respond to user intent.

To assess the rational personalization capabilities of PAs, we introduce **RPEval**, which consists of two components: (1) *Personalized Intent Reasoning Dataset*: simulating how users with different preferences express their intentions through natural, underspecified queries. The dataset is constructed under the design principles of *diversity, naturalness, and consistency* via a tailored *bootstrapping–inversion–validation–expansion* pipeline, covering over 8,000 preferences across twelve task categories. (2) *Multi-granularity Evaluation Protocol*: In the discriminative setting,

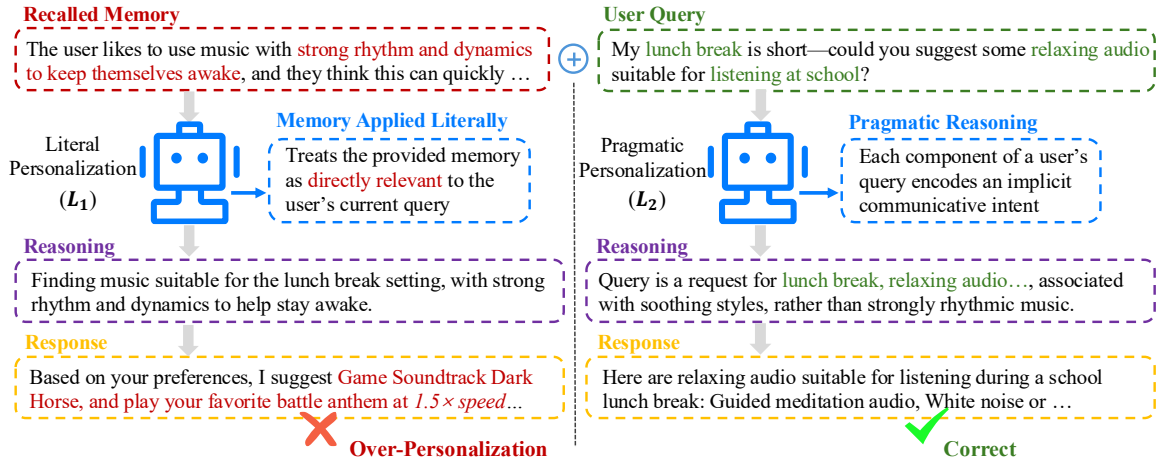


Figure 1: Different levels of PAs. In  $L_1$ , memory is directly concatenated with the query, whereas in  $L_2$ , the PA infers implicit cues from the user’s query to determine memory utilization strategy.

the protocol measures intent classification accuracy; in the generative setting, it operationalizes user-perceivable errors with a systematic taxonomy of irrational personalization phenomena (e.g., *Filter Bubble*) for assessing response quality. Our experiments reveal a 40%–90% accuracy gap between mainstream LLMs and humans on rational personalization, and we characterize the major failure modes of irrational personalization that underlie this gap. Notably, these irrational personalization behaviors become more pronounced as model capability increases, exhibiting an inverse scaling effect. Attribution analysis further indicates that these failures primarily stem from LLMs’ inherent attraction bias (Niu et al., 2025).

To tackle this challenge, we introduce **RP-Reasoner**, which reformulates personalized memory utilization as a reasoning mechanism grounded in pragmatics. Instead of mechanically concatenating memories, RP-REASONER approximates the user’s hidden process of query formulation, extracts cues from surface-level language, infers the underlying intent, and selectively integrates personalized information. Experimental results show that RP-REASONER not only improves intent prediction accuracy by about 35%, reduces error severity by 26% across mainstream LLMs on RPEVAL, but also resolves nearly 80% of the bad cases observed in large-scale commercial PAs, highlighting the substantial potential of pragmatic reasoning to enhance the rationality of PAs. Our key contributions are as follows:

- We propose the problem of Rational Personalization in LLMs. To the best of our knowledge, this is the first study to explore the dual effects of memory in LLM personalization.

- We develop RPEval, a benchmark with a personalized intent reasoning dataset and multi-granularity evaluation, which reveals widespread over-personalization in LLMs and demonstrates, via novel error analysis, how irrational personalization undermines user experience.
- We introduce RP-REASONER, validated on RPEVAL and a large-scale commercial PAs, highlighting the significant potential of pragmatic reasoning to mitigate irrational personalization.

## 2 Rational Personalization

**Problem Formulation.** In personalized response generation, the LLM takes  $(m, q)$  as input, where  $q$  is the user’s current query and  $m$  is the dialogue-context memory comprising multiple user preferences  $\{p_0, p_1, \dots, p_K\}$ . The LLM’s goal is to predict the user’s intent  $i$  and generate a response  $r$ :  $(i, r) = \text{LLM}(m, q)$ . Here,  $i_{\text{query}}$  denotes the true intent underlying  $q$ . While  $m$  may enrich  $q$  and support more accurate inference of  $i_{\text{query}}$ , it can also be irrelevant or even conflict with it. Unlike existing benchmarks (Zhao et al., 2025; Li et al., 2025b; Tan et al., 2025b), we emphasize that the evaluation of  $r$  should not be based solely on its consistency with  $m$ , but rather on its consistency with the user’s true intent  $i_{\text{query}}$ .

**Rational Personalization.** Inspired by Rational Speech Acts theory (Frank and Goodman, 2012), we categorize PAs into three levels based on their memory utilization strategies:

- **Non-personalized ( $L_0$ ):** Ignores  $m$  and predicts  $i$  via semantic matching between  $q$ ,  $i$ :  $P_{L_0}(i, r | q) \propto \text{Semantic}(q, i) \cdot P(i)$ , yielding generic responses, causing **under-personalization**.

- **Literal Personalized ( $L_1$ ):** Directly appends  $m$  with high semantic similarity to the context via an external memory retriever<sup>1</sup>:  $P_{L_1}(i, r | q, m) \propto \text{Semantic}(q, m, i) \cdot P(i | m)$ . While it tailors responses to  $m$ , it risks **over-personalization** when  $m$  is irrelevant to the current  $i_{\text{query}}$ .
- **Pragmatic Personalized ( $L_2$ ):** Viewing personalization as a pragmatic intent understanding task grounded in posterior Bayesian inference:  $P_{L_2}(i, r | m, q) \propto P_{\text{user}}(q | i, m) \cdot P(i | m)$ . By modeling the generative process of how a user formulates and expresses their intent,  $L_2$  can reverse-infer the true  $i$  to adaptively decide whether to integrate or bypass the memory  $m$ , thereby achieving **rational personalization**.

### 3 RPEval

In this section, we introduce RPEVAL, which consists of two components: a Personalized Intent Reasoning Dataset (§ 3.1) and a Multi-Granularity Evaluation Protocol (§ 3.2), providing a systematic evaluation of state-of-the-art LLMs (§ 3.3).

#### 3.1 Dataset Generation

To ensure both reliable and comprehensive intent annotations, we adopt a hierarchical data generation strategy: first narrow the intent space by annotating atomic preference–query pairs to improve consistency, and then expand the data to cover richer, more complex and realistic configurations.

**Atomic Data Format.** We first annotate the intent relation between a single atomic preference and a query, representing it as a quadruple  $(p, q, \text{rationale}, i_{\text{query}})$ . Here, *rationale* denotes the annotation rationale, and  $i_{\text{query}}$  is selected from a finite set of candidate intents: {Ignore, Support, Dominate}.

**Atomic Data Generation.** To systematically generate high-quality data points covering different intent labels, we propose an automated data construction pipeline, whose design is guided by three core objectives: *Diversity*, *Naturalness*, and *Consistency* (Figure 2(a), detailed in Appendix B).

- **Diversity (BOOTSTRAPPING):** To ensure comprehensive coverage of daily scenarios, we adopt a bootstrapping strategy for constructing meta-scenarios. Specifically, we manually define 20 base scenarios as few-shot examples, and then

use GPT-4.1 (OpenAI, 2023) to expand them into new scenarios, which are stored in a data repository. These newly generated scenarios are randomly sampled as few-shot exemplars for subsequent rounds of generation, thereby continuously enriching output diversity. In total, this process yields 100 meta-scenarios (§ B.1).

- **Naturalness (PREFERENCE INVERSION):** In real personalized dialogues, user queries are often brief and do not explicitly mention specific memories. In such cases, LLMs must autonomously decide whether and how to leverage memory. To capture this natural characteristic, we first generate everyday queries from meta-scenarios; then, for each query, we assign different intent labels and then generate the corresponding preferences in an inverted manner (§ B.2).

- **Consistency (ITERATIVE UPDATE):** Drawing on established principles from personalized systems, we developed a three-dimensional quality verification framework: (*Rationality*, *Relevance*, *Alignment*) for intent labeling. Within this framework, we performed iterative updates and evaluations, which not only imposed stringent quality constraints on the generated data but also decoupled the annotation results from the generation logic. This effectively eliminates the circular reasoning bias inherent in synthetic datasets, ensuring that intent labels are grounded in objective logic rather than model heuristics (§ B.3).

Following this pipeline, we curate a large-scale data pool of 8,255 samples. To ensure evaluation rigor, we sub-sample over 1,000 instances for manual verification, reaching a substantial initial inter-annotator agreement of 91.86%. For cases where disagreements persisted, we employ a systematic disambiguation protocol (Algorithm 2), resulting in a finalized gold-standard test set of 953 samples.

**Dataset Expansion** In the previous part, we generate and annotate the applicability relations between individual queries and explicit preferences (*single-explicit*), which already supports basic evaluation. Furthermore, to cover more complex and realistic scenarios, we design two expansion strategies, whose cross-combinations extend the original setting into 6 configurations (§ B.4, § B.5):

- **Explicit2Implicit.** In this strategy, we recast the explicit preference  $p$  as multi-turn dialogues, simulating realistic scenarios where user preferences are implicitly embedded in dialogue history. Un-

<sup>1</sup>In Appendix D.1, we show that semantic similarity alone fails to filter irrelevant memories.

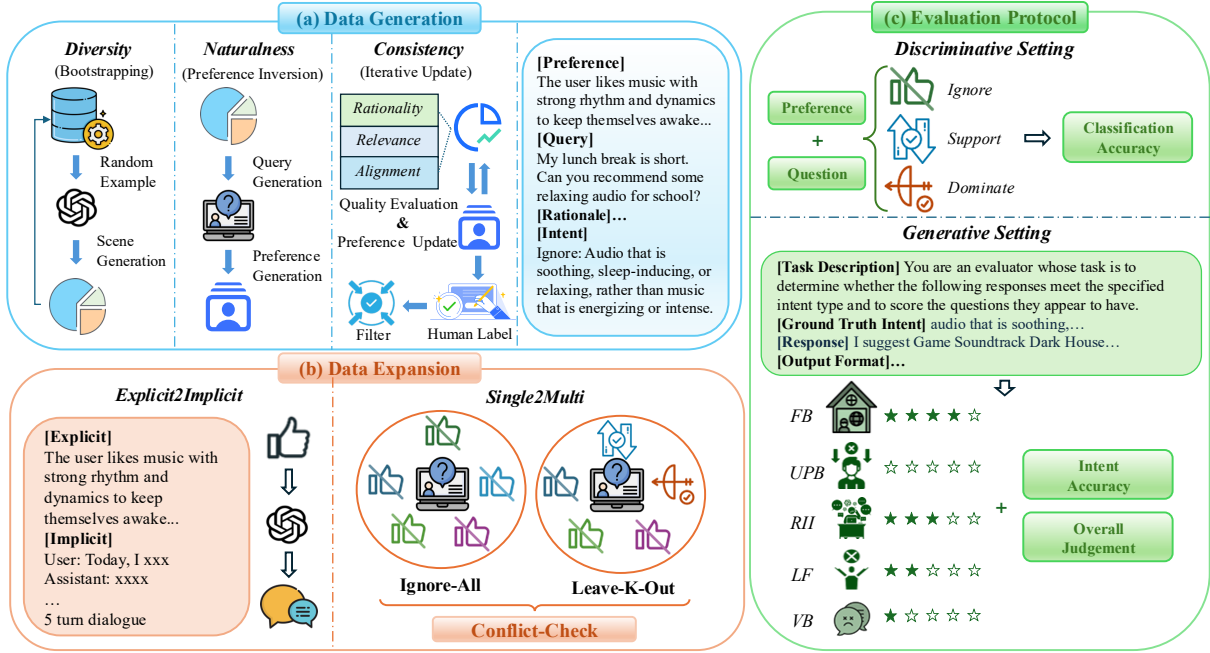


Figure 2: An illustration of our proposed RPEVAL: (a) Personalized Intent Reasoning Data Generation; (b) Dataset Expansion Strategies; (c) Multi-Granularity Evaluation Protocol.

der this setting, the PA is required to infer preferences rather than rely on explicit descriptions.

- **Single2Multi.** To better reflect real-world complexity, we construct *multi-preference* settings. Specifically, we combine different preferences associated with the same query in the initial dataset to derive two configurations: (1) *Ignore-All (IA)*:  $n \in [3, 8]$  irrelevant preferences are provided in the context to test whether the LLM can effectively suppress irrelevant memories under strong distractions. (2) *Leave-K-Out (LKO)*: the memory  $m$  is constructed by combining  $k \in [1, 3]$  relevant preferences with  $n - k$  irrelevant ones, followed by filtering contradictory entries to ensure consistency of the user profile. Under this configuration, the PA must assess the applicability of each preference and accurately identify useful preference amid multiple distractions.

### 3.2 Multi-Granularity Evaluation Protocol

In RPEVAL, we propose a multi-granularity evaluation protocol: a standard intent-matching evaluation in the discriminative setting and an *error-pattern-driven* evaluation in the generative setting.

**Discriminative Setting** In the discriminative setting, our primary objective is to systematically evaluate whether LLMs can correctly determine the applicability of preferences to a user query and the appropriate way to utilize them. We thus specifically design two types of multiple-choice

tasks: (1) *Single-Preference*: Given one preference and a query, the LLM is required to classify how the preference is used as Ignore, Support, or Dominate. (2) *Multi-Preference*: Given a query and several preferences, the LLM needs to independently decide whether and how each preference applies. Our evaluation metric is the classification accuracy across different configurations (§ B.6).

**Generative Setting** In the generative setting, we move beyond relying solely on intent match rates and instead ground evaluation in user-perceivable errors for a more reliable assessment. To this end, we draw on error pattern analysis methods from software engineering (Cemri et al., 2025) to systematically summarize failure modes caused by irrational personalization. Specifically, we first collect a set of representative error types from existing research in personalized and dialogue systems. Then, we manually annotate 200 interaction trajectories, including 100 failed cases extracted from RPEVAL and 100 bad cases from a commercial personalized assistant. The goal of this annotation was to align these representative error types with actual failure cases produced by current LLM-based PAs and validate them. Ultimately, we develop a *two-level* error taxonomy: *strategy-level* ( $\Delta$ ), defined by the memory applicability error taxonomy matrix shown in Table 2, and *response-level* ( $\circ$ ) errors:

- *Filter Bubble* (FB,  $\Delta$ ) (He et al., 2017): Occurs when the PA restricts its response to preference-

Table 1: Comparison with existing benchmarks on personalized memory. *Level* denotes the evaluation tier of personalized assistants within the RPA framework. *Task* specifies the concrete evaluation (**MemQA**: direct question answering over memory content; **Personalization**: generating personalized responses using memory). *Memory Usage Behaviors* represent the strategies of memory utilization, while *Error Phenomena* indicate the types of errors.

Benchmark	Level	Task	Memory Usage Behaviors			Error Phenomena				
			Dominate	Ignore	Support	UPB	FB	RII	LF	VB
MemBench	$L_1$	MemQA	✓	✗	✗	✓	✗	✗	✗	✗
LongMemEval	$L_1$	MemQA	✓	✗	✗	✓	✗	✗	✗	✗
PrefEval	$L_1$	Personalization	✓	✗	✗	✓	✗	✗	✗	✗
ImplexCONV	$L_1$	Personalization	✓	✗	✗	✓	✗	✗	✗	✗
RPEval (ours)	$L_2$	Personalization	✓	✓	✓	✓	✓	✓	✓	✓

Table 2: Strategy-level error taxonomy matrix based on the mismatch between the intended and actual memory utilization strategies. The horizontal axis represents the ground-truth strategy, while the vertical axis indicates the strategy reflected in the LLM’s response.

Predict ↓ / GT →	Ignore	Support	Dominate
Ignore	Correct	UPB	UPB
Support	RII	Correct	RII
Dominate	FB	FB	Correct

specific content while general suggestions would also be appropriate.

- *Redundant Information* (RII,  $\Delta$ ) (Eppler and Mengis, 2004): Occurs when the PA provides both preference-specific and general suggestions, even though the user’s intent only requires one.
- *Under-Personalization* (UPB,  $\Delta$ ) (Zhao et al., 2025; Zhang et al., 2024b): The PA ignores relevant preferences even when the user’s intent requires personalization.
- *Low Feasibility* (LF,  $\circ$ ) (Ji et al., 2023): The response includes impractical or ill-posed suggestions (e.g., recommending music with strong rhythm and dynamics for sleep).
- *Verbose Generation* (VG,  $\circ$ ) (Clark et al., 2021): The PA produces repetitive content, such as superfluous preference restatements.

During evaluation, we develop a *LLM-as-a-Judge* system based on GPT-4.1 (OpenAI, 2023). The system first assesses the alignment between the model response and the ground-truth intent, then assigns a *severity score* (0–5) for each error type, and finally produces an *overall error severity score* (0–5) to capture the degree of user-perceived experience degradation. We instruct 2 human annotators to follow exactly the same evaluation guidelines as the LLM judge, i.e., to assign ordinal scores from 0 to 5 for each of the five personalization error types. Figure 18(c) reports the agreement between the LLM judge and human annotators on these 0–5 or-

dinal ratings, measured by the quadratic-weighted Cohen’s kappa, a standard statistic for ordinal labels. The overall agreement is QWK = 0.87.

**Connect to Related Work.** In Table 1, we compare RPEVAL with existing benchmarks (Zhao et al., 2025; Tan et al., 2025a; Li et al., 2025b; Wu et al., 2025). These benchmarks are typically designed to test whether LLMs can accurately locate and utilize personalized information within long contexts. However, they often assume that user memories play a *Dominate* role, without considering the dual effects of personalization. In contrast, RPEVAL offers an **orthogonal perspective**: it focuses on more realistic scenarios, where PAs must handle preferences with varying applicability and decide on appropriate memory utilization strategies. Moreover, RPEVAL introduces a systematic *error-pattern-driven* analysis, providing a comprehensive account of how irrational personalization affects user experience. A more detailed survey of related work is provided in Appendix A.

### 3.3 Experimental Results and Key Findings

In this section, we conduct a systematic evaluation of mainstream LLMs using RPEVAL, covering small-scale open-source model (Qwen2.5-7B (Qwen et al., 2025)), large-scale open-source model (DeepSeek-V3 (DeepSeek-AI et al., 2025)), and closed-source models including GPT-4.1 (OpenAI, 2023) and the state-of-the-art hybrid reasoning model GPT-5 (OpenAI, 2025)). In all experiments, we explicitly prompt the LLMs to actively judge the applicability of contextual memories (see Appendix D.7; Reminder). In the *discriminative setting*, we compare the intent matching accuracy of different LLMs under the *single-explicit* and *multi-explicit* configurations. For reference, we provide the average accuracy of blind human annotations. In the *generative setting*, we conduct a fine-grained analysis of model responses from the

Table 3: Performance of major LLMs on the discriminative intent matching accuracy in RPEVAL.

	Single.				Multi-MACRO.			Multi-MICRO.		
	Ign.	Sup.	Dom.	ALL	IA	LKO	ALL	IA	LKO	ALL
Human	0.86	1.00	0.98	0.95	0.75	0.71	0.73	0.94	0.93	0.93
Qwen2.5-7B	0.06	0.84	0.24	0.38	0.12	0.02	0.06	0.45	0.36	0.39
Deepseek-v3	0.38	0.78	0.82	0.66	0.05	0.07	0.06	0.57	0.56	0.56
GPT-4.1	0.28	0.34	0.96	0.53	0.08	0.04	0.06	0.52	0.48	0.49
GPT-5	0.12	0.58	0.82	0.51	0.00	0.03	0.02	0.26	0.46	0.39
Hum. Gap ↓	55.8%	16.0%	2.0%	30.5%	84.0%	90.1%	91.8%	39.4%	39.8%	39.8%

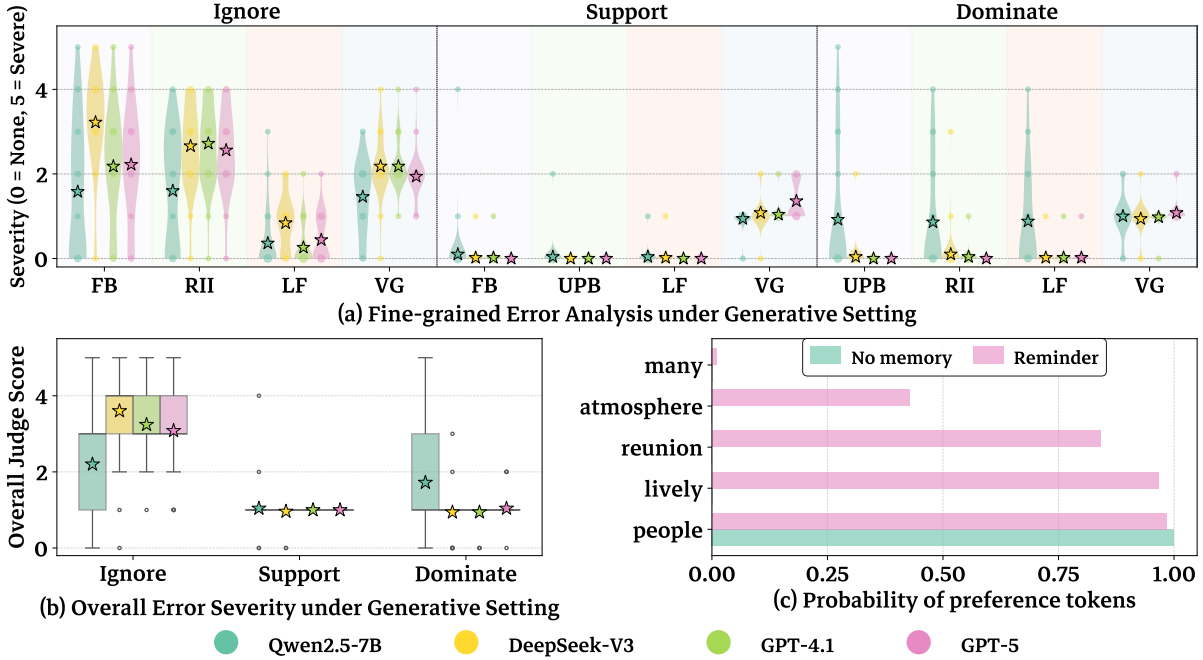


Figure 3: (a) Fine-grained Error Analysis; (b) Overall error severity in the generative setting with the *single-preference* configuration, (c) the reliability of the *LLM-as-a-judge* evaluation.

perspective of user-perceivable errors, and assign an overall severity score. Based on the experimental results, we summarize the following findings:

**Finding I:** LLMs struggle to suppress irrelevant memories and favor a “more-is-better” generation strategy. When the ground-truth intent is Ignore (Table 3), humans can reliably filter out irrelevant memory (e.g., 86% accuracy in the *single-explicit* configuration), whereas LLMs perform considerably worse (6%–38%). This indicates that current LLMs have significant weaknesses in suppressing irrelevant memories. Fine-grained error analysis (Figure 3(a)) further shows that UPB rarely occurs, whereas FB and RII are highly prevalent.

**Finding II:** Multi-preference represents a significant challenge. In the multi-preference setting, LLMs show about a 40% gap from humans in judging the applicability of each preference (Multi-Micro); whereas for overall all-correct accuracy (Multi-Macro), this gap expands to nearly 90%. This indicates that shifting from single

to multiple preferences greatly increases the task difficulty and amplifies the LLMs’ deficiencies in selection and filtering personalized information.

**Finding III:** Inverse scaling effect in rational personalization. In the discriminative setting, more capable base LLMs actually perform worse at ignoring irrelevant preferences. We hypothesize that this counterintuitive behavior stems from their stronger contextual attention, which makes them more inclined to over-utilize preference information rather than suppress it.

**Finding IV:** Attraction bias during generation underlies systematic failures. To better understand why LLMs fail systematically, we conduct a mechanistic analysis. As shown in Figure 3(c)<sup>2</sup>, during decoding, LLMs tend to reuse and amplify tokens or stylistic patterns already present in the context. When the user asks for a gift for a cat, LLM indiscriminately increases the probability mass of irrelevant preference tokens such as “reunion” and

<sup>2</sup>Refer to Appendix D.6 for details.

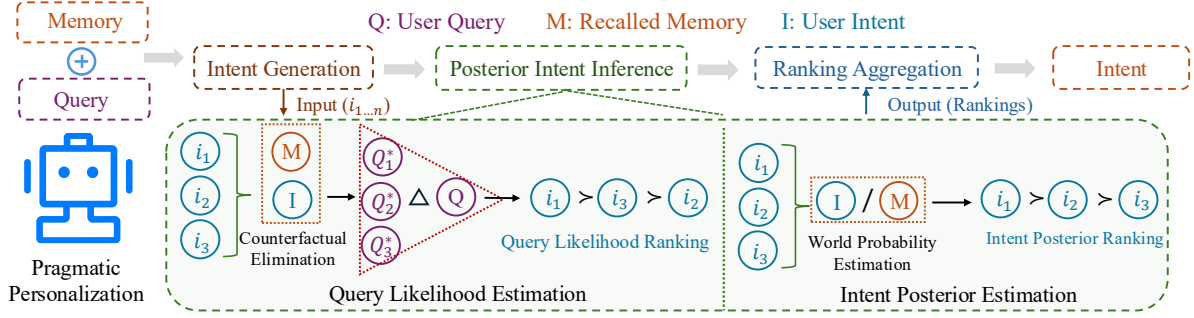


Figure 4: RP-Reasoner: An Implementation of Pragmatic Personalized Assistant.

“lively”, which steers the response toward a party-like style and away from the user’s actual intent.

In summary, our experimental results reveal a clear trend: While current LLMs can leverage memory to generate personalized responses, they lack rational mechanisms to determine whether and how to incorporate it. From a fine-grained user experience perspective, at the strategy level ( $\Delta$ ), LLMs are prone to FB and RII, which either over-constrain responses to preference-specific content and reduce diversity, or blend multiple types of suggestions and increase users’ cognitive burden. At the response level ( $\circ$ ), some LLMs introduce LF when attempting unnecessary personalization, or produce VG by redundantly restating preferences (see Appendix B.6 for detailed case studies).

## 4 Method

### 4.1 RP-Reasoner

Building on the preceding analysis, we observe that existing LLMs lack the capacity for rational memory utilization. To address this gap, we propose **RP-Reasoner**, representing a preliminary foray into constructing a Rational Personalized Assistant ( $L_2$ ) that leverages subtle cues within user queries for intent reasoning. Specifically, we first generate a set of candidate intents under various preference utilization modes:  $\mathcal{I} = \{i_1, \dots, i_n\}$ , and subsequently infer the intent according to the following Bayesian posterior (definition of  $L_2$  in § 2):

$$P(i | q, m) \propto \underbrace{P_{\text{user}}(q | i, m)}_{\text{query likelihood}} \cdot \underbrace{P(i | m)}_{\text{intent prior}},$$

**Query Likelihood Estimation (MLE).** This component simulates how a user would choose the query  $q$  to express a latent intent  $i_{\text{query}}$ , leveraging subtle cues in  $q$  to infer whether specific preferences are implicated. Formally, given memory  $m$ , the goal is to rank candidate intents according to  $P_{\text{user}}(q | i, m)$ . However, this corresponds to a likelihood estimation problem, which essentially

requires inverse modeling of the user’s query generation process. Such a distribution cannot be directly approximated by the world knowledge embedded in an LLM. To address this challenge, we draw inspiration from Approximate Bayesian Computation (Sunnåker et al., 2013) and propose an implicit estimation approach: for each candidate intent  $i$ , we prompt the LLM to estimate the semantic closeness between the observed query  $q$  and an idealized simulated query  $\hat{q}(i, m)$ , thereby approximating the likelihood function. Let  $d_i = \Delta(q, \hat{q}(i, m))$  denote the distance for intent  $i \in \mathcal{I}$ , we rank all candidate intents as follows:

$$\text{rank}_{mle}(i) = 1 + \sum_{j \in \mathcal{I} \setminus \{i\}} \mathbb{I}(\Delta_{q,j} > \Delta_{q,i})$$

This strategy can be understood through the lens of *counterfactual elimination* in pragmatics (Frank and Goodman, 2012): if the user truly intended  $i$ , then there should not exist a substantially better alternative query  $q'$  than the observed  $q$  to express it. Formally, this can be written as:

$$i_{mle}^* = \arg \max_{i \in \mathcal{I}} \mathbf{1} [\forall q' \in \hat{q}(i, m), \Delta(q, q') \geq d_i]$$

This means that if there exists a counterfactual expression  $q'$  that is clearly more suitable than the observed query  $q$  to express intent  $i$ , it indicates that the user’s actual intent is unlikely to be  $i$ .

**Intent Prior Estimation (IPE).** Given memory  $m$ , this component models personalization behaviors independent of  $q$ . It estimates the intent prior  $P(i | m)$  to capture historical preferences. We rank candidate intents as:

$$\text{rank}_{ipe}(i) = 1 + \sum_{j \in \mathcal{I} \setminus \{i\}} \mathbb{I}(P(j | m) > P(i | m))$$

**Aggregation.** RP-REASONER forms the posterior by fusing the query likelihood and the intent prior. Specifically, it minimizes their rank sum:

$$i^* = \arg \min_{i \in \mathcal{I}} (\text{rank}_{mle}(i) + \text{rank}_{ipe}(i))$$

Ties are resolved via  $i_{\text{predict}} \sim \text{Uniform}(\{i^*\})$ .

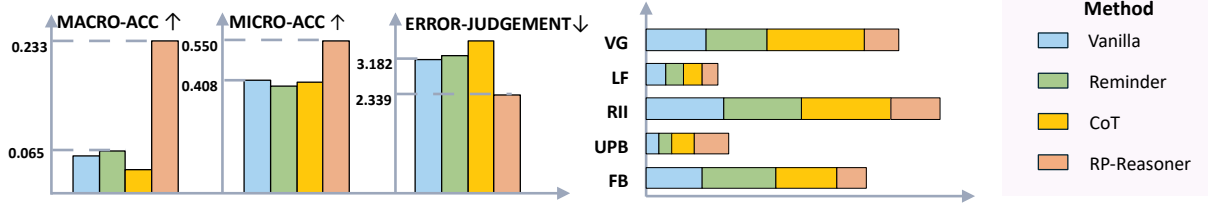


Figure 5: RP-REASONER achieves notable gains in *multi-preference* generative settings.

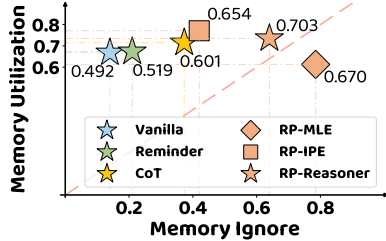


Figure 6: Ablation study.

	RPEval			Real World		
	Macro↑	Micro↑	Judge↓	Macro↑	Micro↑	Judge↓
Vanilla	0.013	0.400	3.533	0.482	0.548	1.899
Reminder	0.013	0.395	3.580	0.126	0.241	3.271
CoT	0.067	0.395	3.487	0.342	0.497	2.216
RP-Reasoner	<b>0.240</b>	<b>0.522</b>	<b>2.493</b>	<b>0.734</b>	<b>0.834</b>	<b>1.070</b>

Table 4: Experimental results of RP-REASONER on both the RPEVAL and real-world failure cases from large-scale commercial PA.

## 4.2 Experimental Results

In this section, we systematically evaluate the effectiveness of RP-REASONER. Specifically, we design three prompt-based baselines (Vanilla, Reminder, and CoT (Wei et al., 2023), CoT-SC (Wang et al., 2023), Self-Refine (Madaan et al., 2023)) and conduct comparative experiments across different backbone LLMs. The experimental analysis is as follows:

**Performance Comparison.** As an initial exploration of applying pragmatic reasoning to mitigate irrational personalization, RP-REASONER achieves significant improvements. As shown in Figure 5, we report the average performance of four models in the *multi-preference* generative setting<sup>3</sup>. The results demonstrate that, compared with the best baseline, RP-REASONER yields a relative improvement of 258% in Macro-acc, 35% in Micro-acc, and a 26% reduction in error severity. A fine-grained error analysis further reveals that RP-REASONER leverages memory in a more rational manner: compared with the baselines, it introduces only a small amount of UPB while effectively mitigating FB and RII issues, and simultaneously exerts partial control over LF and VG.

**Ablation Study.** As shown in Figure 6, we further analyze the roles of different components of RP-REASONER in memory utilization. The vertical axis represents memory utilization (Support & Dominate), while the horizontal axis represents memory ignoring (Ignore). The results show that the MLE module tends to infer from subtle cues in the query whether the user intends to incorpo-

rate preferences, making it more conservative in memory utilization. In contrast, the IPE module reasons about intent plausibility and only considers preferences that the user is likely to accept, thereby exhibiting a more permissive attitude toward memory utilization. When combined, the two modules strike a balance between over-reliance on and excessive neglect of historical memory.

**Real-World Validation.** In Table 4, we further evaluate personalized response generation on both the RPEVAL benchmark and a large-scale personalized dialogue assistant that has been deployed and is actively maintained by the business team. The results reveal that: (1) the data patterns observed in RPEVAL are highly consistent with those in real-world scenarios, where all baseline methods exhibit significant irrational personalization issues, thereby validating the reliability of our benchmark; (2) RP-REASONER achieves substantial improvements in both settings, successfully resolving about 80% of the error cases in real business deployment. These findings demonstrate that RP-REASONER not only excels in benchmark evaluations but also delivers tangible value in practical applications.

## 5 Conclusion

In this study, we propose the challenge of Rational Personalization for LLM-based assistants and accordingly develop RPEVAL, an evaluation framework comprising a personalized intent reasoning dataset and a multi-granularity evaluation protocol. This framework enables a systematic analysis of the dual effects of personalization. Furthermore, we propose RP-REASONER to explore and validate the potential of pragmatic reasoning in building more rational and user-aligned intelligent assistants.

<sup>3</sup>Detailed comparisons with other baselines, along with full results and error analysis, are provided in Appendix D.7.

## 562 Limitations

563 (1) Our benchmark focuses on evaluating person-  
564 alized assistants’ ability to rationally utilize user  
565 preferences of varying applicability within context,  
566 rather than testing LLMs’ long-context memory re-  
567 call capabilities or the factual accuracy of memory-  
568 based QA. While these are also important direc-  
569 tions, they represent orthogonal research dimen-  
570 sions and are therefore beyond the scope of this  
571 work. (2) Although we made substantial efforts, in-  
572 cluding establishing unified annotation guidelines  
573 and conducting double-blind human annotation to  
574 ensure the consistency and verifiability of intent  
575 labeling, the inherently subjective nature of person-  
576 alized assistants means that no perfectly objective  
577 standard exists. Nevertheless, we believe that our  
578 proposed memory utilization criteria and rigorous  
579 annotation protocols provide a solid starting point.  
580 As personalized assistants are deployed at larger  
581 scales in real-world applications, future work can  
582 leverage more authentic user data to better capture  
583 such cases.

## 584 Ethical Considerations

585 In this work, we introduce RPEVAL, a benchmark  
586 for evaluating LLMs’ ability to rationally utilize  
587 user preferences. Our study places a strong em-  
588 phasis on responsible and ethical practices, with  
589 particular attention to data privacy, ethical consid-  
590 erations of data quality, bias mitigation, and research  
591 integrity. Since all memory contents and queries  
592 were newly created, we conducted a rigorous man-  
593 ual screening process to ensure that the dataset  
594 contains no personally identifiable information or  
595 inappropriate content. This work involves human  
596 annotation in two places: (1) dataset construction  
597 and filtering (§ 3.1); (2) evaluation with LLM-as-a-  
598 Judge (§ 3.2). The process was mainly conducted  
599 by four expert annotators, who are in-house NLP  
600 researchers with more than three years of experi-  
601 ence in dialogue assistant research. In addition,  
602 developers from the commercial dialogue assis-  
603 tant product team participated in summarizing and  
604 screening real-world bad cases (§ 3.2), defining the  
605 rationality of quality annotation protocols (§ 3.1),  
606 and providing feedback on the validity of error cat-  
607 egorization standards (§ 3.2). All annotators were  
608 thoroughly briefed with the annotation objectives,  
609 and any uncertain cases were resolved through dis-  
610 cussion among the annotators. In total, approxi-  
611 mately **200 human hours** were spent on annota-

tion. Annotators were compensated on a monthly  
basis, and their salaries included the working hours  
dedicated to annotation.

RPEVAL may entail dual societal impacts. On  
the one hand, introducing rational preference uti-  
lization mechanisms enables assistants to better  
capture user intent, reduce inappropriate memory  
calls, improve interaction efficiency and satisfac-  
tion, and promote more transparent, controllable,  
and responsible personalization. On the other hand,  
more precise memory usage may increase risks of  
privacy leakage and misuse, and, in the absence  
of effective regulation, could be exploited for over-  
profiling and manipulation. We therefore call for  
research and applications in this direction to be ac-  
companied by strict privacy protection and safety  
safeguards, ensuring that technological progress  
truly serves users and societal well-being.

## References

- Gediminas Adomavicius, Bamshad Mobasher,  
Francesco Ricci, and Alexander Tuzhilin. 2011.  
[Context-aware recommender systems](#). *AI Magazine*,  
32(3):67–80.
- Mert Cemri, Melissa Z. Pan, Shuyi Yang, Lakshya A.  
Agrawal, Bhavya Chopra, Rishabh Tiwari, Kurt  
Keutzer, Aditya Parameswaran, Dan Klein, Kannan  
Ramchandran, Matei Zaharia, Joseph E. Gonzalez,  
and Ion Stoica. 2025. [Why do multi-agent llm sys-  
tems fail?](#) *Preprint*, arXiv:2503.13657.
- Allison J. B. Chaney, Brandon M. Stewart, and Bar-  
bara E. Engelhardt. 2018. [How algorithmic con-  
founding in recommendation systems increases ho-  
mogeneity and decreases utility](#). In *Proceedings of  
the 12th ACM Conference on Recommender Systems*,  
RecSys ’18, page 224–232, New York, NY, USA.  
Association for Computing Machinery.
- Elizabeth Clark, Tal August, Sofia Serrano, Nikita  
Haduong, Suchin Gururangan, and Noah A. Smith.  
2021. [All that’s ‘human’ is not gold: Evaluat-  
ing human evaluation of generated text](#). *Preprint*,  
arXiv:2107.00061.
- DeepSeek-AI, Aixin Liu, Bei Feng, and et al.  
2025. [Deepseek-v3 technical report](#). *Preprint*,  
arXiv:2412.19437.
- Martin J Eppler and Jeanne Mengis. 2004. The con-  
cept of information overload: A review of literature  
from organization science, accounting, marketing,  
mis, and related disciplines. *The Information Society*,  
20(5):325–344.
- Michael C Frank and Noah D Goodman. 2012. Predict-  
ing pragmatic reasoning in language games. *Science*,  
336(6084):998–998.

664	H. Paul Grice. 1975. Logic and conversation. In Peter Cole and Jerry L. Morgan, editors, <i>Syntax and Semantics, Volume 3: Speech Acts</i> , pages 41–58. Academic Press, New York.	
665		
666		
667		
668	Jerry Zhi-Yang He, Sashrika Pandey, Mariah L. Schrum, and Anca Dragan. 2025. <a href="#">Context steering: Controllable personalization at inference time</a> . <i>Preprint</i> , arXiv:2405.01768.	
669		
670		
671		
672	Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. 2017. <a href="#">Neural collaborative filtering</a> . <i>Preprint</i> , arXiv:1708.05031.	
673		
674		
675	Ziwei Ji, Nayeon Lee, Jason Fries, Tao Yu, Danqi Zhang, and Chelsea Finn. 2023. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. <i>arXiv preprint arXiv:2309.05922</i> .	
676		
677		
678		
679		
680	Jia-Nan Li, Jian Guan, Songhao Wu, Wei Wu, and Rui Yan. 2025a. <a href="#">From 1,000,000 users to every user: Scaling up personalized preference for user-level alignment</a> . <i>Preprint</i> , arXiv:2503.15463.	
681		
682		
683		
684	Xintong Li, Jalend Bantupalli, Ria Dharmani, Yuwei Zhang, and Jingbo Shang. 2025b. <a href="#">Toward multi-session personalized conversation: A large-scale dataset and hierarchical tree framework for implicit reasoning</a> . <i>Preprint</i> , arXiv:2503.07018.	
685		
686		
687		
688		
689	Junru Lu, Siyu An, Mingbao Lin, Gabriele Pergola, Yulan He, Di Yin, Xing Sun, and Yunsheng Wu. 2023. <a href="#">Memochat: Tuning llms to use memos for consistent long-range open-domain conversation</a> . <i>Preprint</i> , arXiv:2308.08239.	
690		
691		
692		
693		
694	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023. <a href="#">Self-refine: Iterative refinement with self-feedback</a> . <i>Preprint</i> , arXiv:2303.17651.	
695		
696		
697		
698		
699		
700		
701		
702	Adyasha Maharana, Dong-Ho Lee, Sergey Tulyakov, Mohit Bansal, Francesco Barbieri, and Yuwei Fang. 2024. <a href="#">Evaluating very long-term conversational memory of llm agents</a> . <i>Preprint</i> , arXiv:2402.17753.	
703		
704		
705		
706	Jakob Nielsen. 1994. <i>Usability Engineering</i> . Academic Press.	
707		
708	Jingcheng Niu, Xingdi Yuan, Tong Wang, Hamidreza Saghir, and Amir H. Abdi. 2025. <a href="#">Llama see, llama do: A mechanistic perspective on contextual entrainment and distraction in LLMs</a> . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 16218–16239, Vienna, Austria. Association for Computational Linguistics.	
709		
710		
711		
712		
713		
714		
715		
716	OpenAI. 2023. Gpt-4 technical report. <a href="https://openai.com/research/gpt-4">https://openai.com/research/gpt-4</a> . Accessed: August 7, 2025.	
717		
718		
	OpenAI. 2025. Gpt-5 technical report. <a href="https://openai.com/zh-Hans-CN/index/introducing-gpt-5/">https://openai.com/zh-Hans-CN/index/introducing-gpt-5/</a> . Accessed: August 7, 2025.	719 720 721
	Eli Pariser. 2011. <i>The Filter Bubble: What the Internet Is Hiding from You</i> . Penguin Press, New York.	722 723
	Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jingren Zhou, and 25 others. 2025. <a href="#">Qwen2.5 technical report</a> . <i>Preprint</i> , arXiv:2412.15115.	724 725 726 727 728 729 730
	Nils Reimers and Iryna Gurevych. 2019. <a href="#">Sentence-bert: Sentence embeddings using siamese bert-networks</a> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing</i> . Association for Computational Linguistics.	731 732 733 734 735
	Jingbo Shang, Zai Zheng, Jiale Wei, Xiang Ying, Felix Tao, and Mindverse Team. 2024. <a href="#">Ai-native memory: A pathway from llms towards agi</a> . <i>Preprint</i> , arXiv:2406.18312.	736 737 738 739
	Ben Shneiderman. 1987. <i>Designing the User Interface: Strategies for Effective Human-Computer Interaction</i> . Addison-Wesley.	740 741 742
	Mikael Sunnåker, Alberto Giovanni Busetto, Elina Numminen, Jukka Corander, Matthieu Foll, and Christophe Dessimoz. 2013. Approximate bayesian computation. <i>PLoS computational biology</i> , 9(1):e1002803.	743 744 745 746 747
	John Sweller. 1988. Cognitive load during problem solving: Effects on learning. <i>Cognitive Science</i> , 12(2):257–285.	748 749 750
	Haoran Tan, Zeyu Zhang, Chen Ma, Xu Chen, Quanyu Dai, and Zhenhua Dong. 2025a. <a href="#">Membench: Towards more comprehensive evaluation on the memory of llm-based agents</a> . <i>Preprint</i> , arXiv:2506.21605.	751 752 753 754
	Juntao Tan, Liangwei Yang, Zuxin Liu, Zhiwei Liu, Rithesh Murthy, Tulika Manoj Awalganekar, Jianguo Zhang, Weiran Yao, Ming Zhu, Shirley Kokane, Silvio Savarese, Huan Wang, Caiming Xiong, and Shelby Heinecke. 2025b. <a href="#">Personabench: Evaluating ai models on understanding personal information through accessing (synthetic) private user data</a> . <i>Preprint</i> , arXiv:2502.20616.	755 756 757 758 759 760 761 762
	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. <a href="#">Self-consistency improves chain of thought reasoning in language models</a> . <i>Preprint</i> , arXiv:2203.11171.	763 764 765 766 767
	Zheng Wang, Zhongyang Li, Zeren Jiang, Dandan Tu, and Wei Shi. 2024. <a href="#">Crafting personalized agents through retrieval-augmented generation on editable memory graphs</a> . <i>Preprint</i> , arXiv:2409.19401.	768 769 770 771

772 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten  
773 Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and  
774 Denny Zhou. 2023. *Chain-of-thought prompting elicits reasoning in large language models*. *Preprint*,  
775 arXiv:2201.11903.

777 Di Wu, Hongwei Wang, Wenhao Yu, Yuwei Zhang,  
778 Kai-Wei Chang, and Dong Yu. 2025. *Longmemeval: Benchmarking chat assistants on long-term interactive memory*. *Preprint*, arXiv:2410.10813.

781 Hongli Yu, Tinghong Chen, Jiangtao Feng, Jiangjie  
782 Chen, Weinan Dai, Qiyang Yu, Ya-Qin Zhang, Wei-  
783 Ying Ma, Jingjing Liu, Mingxuan Wang, and Hao  
784 Zhou. 2025. *Memagent: Reshaping long-context llm with multi-conv rl-based memory agent*. *Preprint*,  
785 arXiv:2507.02259.

787 Kai Zhang, Yejin Kim, and Xiaozhong Liu. 2025. *Personalized llm response generation with parameterized memory injection*. *Preprint*, arXiv:2404.03565.

790 Zeyu Zhang, Quanyu Dai, Xiaohe Bo, Chen Ma, Rui Li,  
791 Xu Chen, Jieming Zhu, Zhenhua Dong, and Ji-Rong  
792 Wen. 2024a. A survey on the memory mechanism of  
793 large language model based agents. *ACM Transactions on Information Systems*.

795 Zeyu Zhang, Quanyu Dai, Luyu Chen, Zeren Jiang, Rui  
796 Li, Jieming Zhu, Xu Chen, Yi Xie, Zhenhua Dong,  
797 and Ji-Rong Wen. 2024b. *Memsim: A bayesian simulator for evaluating memory of llm-based personal assistants*. *arXiv preprint arXiv:2409.20163*.

800 Zhehao Zhang, Ryan A Rossi, Branislav Kveton, Yi-  
801 jia Shao, Diyi Yang, Hamed Zamani, Franck Der-  
802 noncourt, Joe Barrow, Tong Yu, Sungchul Kim, and  
803 1 others. 2024c. Personalization of large language  
804 models: A survey. *arXiv preprint arXiv:2411.00027*.

805 Siyan Zhao, Mingyi Hong, Yang Liu, Devamanyu Haz-  
806 arika, and Kaixiang Lin. 2025. Do llms recognize  
807 your preferences? evaluating personalized preference  
808 following in llms. *arXiv preprint arXiv:2502.09597*.

809 Wanjun Zhong, Lianghong Guo, Qiqi Gao, He Ye, and  
810 Yanlin Wang. 2023. *Memorybank: Enhancing large  
811 language models with long-term memory*. *Preprint*,  
812 arXiv:2305.10250.

<b>Contents</b>		813
<b>1 Introduction</b>	<b>1</b>	814
<b>2 Rational Personalization</b>	<b>2</b>	815
<b>3 RPEval</b>	<b>3</b>	816
3.1 Dataset Generation . . . . .	3	817
3.2 Multi-Granularity Evaluation Protocol . . . . .	4	819
3.3 Experimental Results and Key Findings . . . . .	5	820
<b>4 Method</b>	<b>7</b>	822
4.1 RP-Reasoner . . . . .	7	823
4.2 Experimental Results . . . . .	8	824
<b>5 Conclusion</b>	<b>8</b>	825
<b>A Related Work</b>	<b>12</b>	826
<b>B Supplemental Details for RPEval</b>	<b>13</b>	827
B.1 Diversity . . . . .	13	828
B.2 Naturalness . . . . .	13	829
B.3 Consistency . . . . .	13	830
B.4 Dataset Expansion . . . . .	16	831
B.5 Basic Statistics . . . . .	16	832
B.6 Evaluation Metric Building . . . . .	16	833
<b>C Supplemental Details for RP-Reasoner</b>	<b>28</b>	834
C.1 Implementation Details . . . . .	28	835
C.2 Inference Cost Optimization . . . . .	28	836
<b>D Extended Analysis</b>	<b>33</b>	837
D.1 Empirical Analysis of Semantic Similarity . . . . .	33	838
D.2 Model Background . . . . .	33	840
D.3 Baseline Details . . . . .	33	841
D.4 Comparison with Additional Baselines . . . . .	33	842
D.5 Reliability of LLM-as-judge . . . . .	34	844
D.6 Analysis of the Mechanism . . . . .	34	845
D.7 Discriminative Setting . . . . .	34	846
D.8 Generative Setting . . . . .	35	847
<b>E Future Works</b>	<b>43</b>	848

## 849 A Related Work

850 **Personalized Memory Benchmarks.** LLM-  
851 based agents have been widely applied across do-  
852 mains, marking the advent of a new era of per-  
853 sonal assistants. A key research direction is how  
854 to endow agents with *memory*, enabling them to  
855 retain past dialogues and tasks, and update their  
856 understanding of users for more personalized and  
857 consistent services. Early evaluation benchmarks  
858 mainly focused on accuracy, often formulated as  
859 *memQA* tasks (Tan et al., 2025b; Maharana et al.,  
860 2024). For instance, MemoryBank (Zhong et al.,  
861 2023) contains multi-day chat histories from 15  
862 users with 194 human-written probing questions.  
863 MemSim/MemBench (Zhang et al., 2024b; Tan  
864 et al., 2025a) further introduce a Bayesian relation  
865 network to generate reliable QA pairs automati-  
866 cally, while LongMemEval (Wu et al., 2025) cov-  
867 ers diverse core long-term memory abilities such  
868 as information extraction, multi-session reasoning,  
869 and knowledge updates.

870 However, in realistic personalized assistant sce-  
871 narios, memQA is fundamentally different from  
872 user queries: users rarely ask direct memory ques-  
873 tions, but instead pose open-ended, life-oriented  
874 queries. To address this gap, recent benchmarks  
875 shift to *memory-supported downstream personal-*  
876 *ization tasks*. For example, PrefEval (Zhao et al.,  
877 2025) evaluates whether LLMs can proactively  
878 leverage user preferences in long texts, while Im-  
879 plicitConv (Li et al., 2025b) tests assistants’ capa-  
880 bility for implicit personalized reasoning. Never-  
881 theless, these benchmarks still focus on memory-  
882 retrieval-centered objectives, with error analysis  
883 limited to UPB, essentially testing models only at  
884 the  $A_0$  level—whether preferences can be directly  
885 mapped to responses. In contrast, our proposed  
886 RPEVAL offers an orthogonal perspective: it em-  
887 phasizes realistic scenarios where personal assis-  
888 tants must handle preferences with varying appli-  
889 cability, balance them with general content, and  
890 infer user intent. Moreover, RPEVAL introduces  
891 a systematic error-pattern analysis, revealing how  
892 irrational personalization negatively impacts user  
893 experience.

894 **Personalized Assistants.** Approaches to build-  
895 ing personalized assistants broadly fall into three  
896 lines: (1) *long-context methods* (Yu et al., 2025),  
897 which extend LLM input windows to directly con-  
898 sume lengthy interaction histories but inevitably  
899 introduce irrelevant content and suffer from the

“lost-in-the-middle” effect (Zhao et al., 2025); (2)  
*parametric memory* (Li et al., 2025a; Zhang et al.,  
2025; He et al., 2025), which encodes user pref-  
erences into model parameters via fine-tuning or  
prompt tuning, often leading to overfitted person-  
alization and limited flexibility; and (3) *retrieval-*  
*augmented generation* (Lu et al., 2023; Shang et al.,  
2024; Li et al., 2025b; Wang et al., 2024), which  
retrieves external memories to support personal-  
ization but relies on similarity- or rule-based re-  
trieval, making it prone to noisy or off-target re-  
calls. Despite different mechanisms, all three face a  
common challenge: achieving rational personaliza-  
tion—leveraging memory when appropriate while  
avoiding excessive or improper personalization so  
that memory serves intent inference rather than  
mere restatement of history. To address this, we  
introduce RPEVAL, a benchmark that evaluates not  
only whether models use memory but also whether  
they decide *when and how* to use it under realistic  
conditions. Building on this perspective, we further  
propose RP-REASONER, a simple reasoning mod-  
ule grounded in pragmatic Bayesian inference that  
moves beyond “how to remember” toward “how  
to use memory well,” enabling more robust and  
contextually appropriate personalization.

**The Dual Effects of Personalization** Personal-  
ization has long been recognized as a double-edged  
sword in traditional information systems (Pariser,  
2011; Adomavicius et al., 2011; He et al., 2017).  
However, in the nascent field of LLM-based per-  
sonalized assistants, this problem remains largely  
unmodeled, and corresponding benchmarks are  
lacking. Furthermore, traditional context-aware  
recommendation datasets and evaluation protocols  
predicated on large-scale exposure or click logs are  
not directly applicable to LLMs due to fundamental  
differences in data and decision structures. LLM-  
based personalization relies on sparse, free-form  
textual memories rather than structured interaction  
logs, and its action space consists of open-ended  
natural language rather than ranking over fixed item  
sets. To bridge this gap, we are the first to sys-  
tematically characterize and evaluate the problem  
of **Rational Personalization** in LLM-powered as-  
sistants, providing both a novel benchmark and a  
formal reasoning framework.

947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994

## B Supplemental Details for RPEval

In this section, we discuss the details of our benchmark construction process.

### B.1 Diversity

In Table 5, we present the scenarios generated by RPEVAL through bootstrapping. RPEVAL covers a total of 100 everyday life scenarios, each defined by two elements: *What* and *Why*. Here, *What* specifies the concrete situation or activity (e.g., family trip planning, friends’ gathering, healthy diet plan), while *Why* captures the underlying motivation or need that drives the situation (e.g., strengthening family bonds, relaxing with friends, maintaining health). In subsequent steps, a large number of personalized reasoning data points will be derived from each meta-scenario.

### B.2 Naturalness

Based on each scenario, we prompt GPT-4.1 (See Figure 8) to generate 5–10 user daily queries. These queries are intentionally brief and typically do not explicitly reference any memory, such as “*Our family is going on a trip, do you have any recommendations?*”. In total, we obtain around 800 daily queries. For each query, we then assign one of three memory utilization intent labels (Ignore, Support, Dominate). For every  $(q, i)$  pair, we prompt GPT-4.1 (See Figure 9) to generate approximately 5–10 corresponding user memories. For example, in a *family travel* scenario, when the intent is Ignore, the model may generate a memory such as a *personal preference for horror movies*, which is unrelated to travel and unsuitable for the family context. Altogether, this process yields roughly **15,000**  $(p, q, i)$  intent reasoning data points.

### B.3 Consistency

**Quality Verification Standard.** Constructing a high-quality and verifiable benchmark for personalized memory intent reasoning is non-trivial. It requires systematic efforts and rigorous design principles. In practice, we identify two major challenges: (1) When both the user query and candidate persona information are simultaneously exposed to the model or human annotators, they often assume that the persona *must* be used. This leads to overly positive judgments and distorts the assessment of whether personalization is rational. (2) Annotators may disagree on whether invoking a particular memory is reasonable for a given query, resulting

in labels that lack verifiability and stability.

To address these issues, we draw on quality standards from Human-Computer Interaction (HCI), Usability Engineering, and classic principles from recommender systems/personalized AI. We summarize three-dimensional *Quality Verification Standard* for memory utilization intent labeling and implement quality assurance through iterative updates and manual cross-validation. Each candidate sample  $(p, q, i)$  must satisfy these criteria before being admitted into the benchmark.

In implementation, we ensure data quality through a combination of LLM-based verification and data updating, followed by human fine-grained annotation for further assurance.

**Data Consistency Guarantee (Automatic)** In the automated pipeline, we build a GPT-4.1-based data quality verifier (See Figure 11), which scores each sample across three dimensions on a 0–5 scale. We then update the personas according to the proposed standards (See Figure 10) to ensure better compliance with the quality criteria. Finally, only samples that achieve full scores (i.e., 5 in all dimensions) are retained, resulting in about **8,000** samples that are regarded as relatively high-quality data.

We then randomly select a subset of the full-score samples and invite human annotators to re-label them. The results show an accuracy of about 80%. These high-quality samples are included in the complete dataset we release. In addition, we will conduct human fine annotation on a separate test set in the next step to further ensure reliability and robustness.

**Data Consistency Guarantee (Human)** Although strict automatic quality control provides strong guarantees for data generation, we further construct a higher-quality test set by randomly sampling a subset and employing annotators for rigorous double-blind labeling. In total, we annotate about 1,000 data points. The annotation procedure follows the standard in table 6, which is provided to annotators during the process.

**Inter-annotator Agreement Statistics** For the full 8K dataset, we randomly sampled approximately 1,000 instances for human annotation. Each instance was independently annotated by two annotators in a blind setting. We then compared both annotators’ labels with the LLM-generated label. Whenever at least one annotator disagreed with

995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044

Table 5: The meta-scenarios constructed via bootstrapping are designed to simulate the domains of everyday life relevant to PAs, and serve as the basis for building personalized intent reasoning data.

Attribute 1: What
<p><b>1.1: Family &amp; Parenting:</b> Family outing &amp; trip planning, Weekend parent–child activities, Family short-trip planning, Family weekend outdoor activity plan, Family weekend exercise plan, Parent–child craft activities, Parent–child outdoor sports suggestions, Parent–child painting activities, Parent–child reading list for holidays, Parent–child holiday decor DIY, Family game night activities, Weekend family watchlist, Family healthy breakfast pairing, Weekend picnic prep, Family weekend farm experience, Family diary / photo album organizing, Elder’s birthday dinner venue, Family weekend movie list (family-friendly), Family-friendly city day trip, Family short self-drive planning, Family travel packing list, Family carry-on essentials, Family city food exploration</p> <p><b>1.2: Friends / Colleagues &amp; Social:</b> Friends’ gathering – activity ideas, Friends’ gathering – restaurant choice, Friends’ gathering mini-games, Friends’ birthday surprise planning, Birthday party planning, Wedding gift for friends, Colleague farewell gift selection, Friends’ gathering theme design.</p> <p><b>1.3: Couples &amp; Dating:</b> Date activity plan, Couples’ date restaurant choice, Partner anniversary surprise, Anniversary trip planning, Niche gift ideas for a partner, Couples’ holiday surprise plan, Handmade gift for a partner, Couple shared reading list</p> <p><b>1.4: Personal Growth &amp; Health:</b> Personal workout plan, Short-term fitness training plan, Healthy eating plan, Short-term fat-loss diet plan, After-work light-meal menu, Morning run route planning, Learning a new skill (e.g., instrument), Short-term language learning resources, Short-term upskilling course recommendations, Weekend kitchen “new dishes” try-outs, Healthy snack recommendations, Keep-fit equipment shopping for home, Family weekend sports plan</p> <p><b>1.5: Pets:</b> Pet travel/outdoor gear prep, Weekend outdoor activities with pets, Pet holiday gift recommendations</p> <p>...</p> <p><b>1.13: Work &amp; Low-Social Leisure:</b> Job-interview outfit advice, After-work leisure suggestions (low-social), After-work solo activities</p>
Attribute 2: Why
<p><b>2.1: Family &amp; Parenting Bonds:</b> Strengthen family relationships, Enhance parent–child bonds, Provide sense of companionship, Make children happy, Foster creativity and patience</p> <p><b>2.2: Friendship &amp; Social Bonds:</b> Relax and strengthen friendships, Everyone can easily reach, Create lively atmosphere, Avoid awkward silence</p> <p><b>2.3: Romantic &amp; Couple Bonds:</b> Express care and affection, Create romance, Leave memorable moments, Show attentiveness</p> <p><b>2.4: Health &amp; Growth:</b> Improve routines, Increase physical strength, Maintain health, Healthy weight management, Efficient fat loss, Self-improvement, Holiday learning and growth, Build reading interest, Personal development</p> <p><b>2.5: Pet Care:</b> Ensure comfort and safety, Let pets feel love and care</p> <p><b>2.6: Holidays &amp; Rituals:</b> Express feelings, Create festive atmosphere, Balance tradition and innovation, Enhance sense of ceremony, Make gatherings harmonious</p> <p>...</p> <p><b>2.11: Impression Management:</b> Leave a good impression, Appear professional but not rigid, First meeting confidence</p> <p><b>2.12: Self-Improvement &amp; Learning:</b> Refresh the mind, Prepare for overseas communication, Broaden horizons, Skill growth</p> <p><b>2.13: Diet &amp; Health Specifics:</b> Control diet for weight, Balanced nutrition, Quick energy recovery</p> <p><b>2.14: Leisure &amp; Rest Specifics:</b> Not staying at home all day, Low-social relaxation, Independent unwinding, Enjoy downtime without boredom</p>

Table 6: Three-dimensional Quality Verification Standard for intent reasoning samples. Each candidate  $(p, q, i)$  must satisfy these criteria before being admitted into the benchmark.

<p><b>Dimension 1: Rationality</b></p> <p><b>Definition:</b> The system-generated response is reasonable.</p> <p><b>Theoretical Support:</b> One of Nielsen’s 10 Usability Heuristics (Nielsen, 1994): User Control and Freedom, which means that the user’s current intent should take precedence, and the system should not override the user’s choices due to historical preferences.</p> <p><b>Why:</b> If the rationality constraint is violated, the system may lead to factual errors or cause strong user dissatisfaction, disrupting the natural flow of the interaction. Ensuring that rationality takes priority helps avoid misleading the user and facilitates the smooth achievement of task goals.</p> <hr/> <p>Ignore: does not improve, or may even mislead, the response. Example: <i>Persona: likes blue</i> → irrelevant when recommending a reading list.</p> <p>Support: can improve the quality of the response but is not mandatory. Example: <i>Persona: likes spicy food</i> → offering spicy options when recommending restaurants is better, but not compulsory.</p> <p>Dominate: must be strictly followed; otherwise it leads to factual errors or strong user rejection. Example: <i>Persona: vegetarian</i> → recommending steak is a severe conflict.</p>
<p><b>Dimension 2: Relevance</b></p> <p><b>Definition:</b> Personalization should avoid introducing irrelevant information into the dialogue response.</p> <p><b>Theoretical Support:</b> Grice’s Conversational Maxims (Quantity, Relevance) (Grice, 1975) states that information should be sufficient but not excessive, and it should be relevant to the dialogue goal. Sweller’s Cognitive Load Theory (Sweller, 1988) indicates that irrelevant information increases cognitive load, hindering the accomplishment of task objectives.</p> <p><b>Why:</b> Irrelevant personalized information increases cognitive load, creates noise, and reduces the conciseness of the dialogue, thereby affecting the user’s task execution efficiency.</p> <hr/> <p>Ignore: the user would almost never recall this memory. Example: <i>Query: checking weather</i> → irrelevant to <i>likes Sichuan cuisine</i>.</p> <p>Support: the user may recall it, but not necessarily. Example: <i>Query: ordering food</i> → might recall <i>likes Sichuan cuisine</i>.</p> <p>Dominate: the user would inevitably recall this memory. Example: <i>Query: ordering food</i> → always recall <i>vegetarianism</i>.</p>
<p><b>Dimension 3: Alignment</b></p> <p><b>Definition:</b> Personalization should align with the user’s query focus.</p> <p><b>Theoretical Support:</b> In human-to-human conversation, the listener typically responds directly to the question asked, rather than introducing personal habits or background information. One of Shneiderman’s Eight Principles of Interface Design is “Consistency” (Shneiderman, 1987). Human-computer dialogue design should follow this principle, providing responses that directly address the user’s query.</p> <p><b>Why:</b> If personalized information does not align with the task goal, it can make the system feel intrusive or unnatural, undermining trust and reliability. Personalization should be closely aligned with the task goal to ensure a smooth and effective dialogue experience.</p> <hr/> <p>Ignore: query focus is unrelated to the preference. Example: <i>Query: holiday dining</i> ↔ <i>Habit of eating fast food on weekdays</i>.</p> <p>Support: query focus and preference are compatible, but general advice is also acceptable. Example: <i>Query: holiday dining</i> ↔ <i>likes spicy food</i>.</p> <p>Dominate: query focus and preference are fully aligned. Example: <i>Query: holiday dining</i> ↔ <i>Vegetarian identity</i>.</p>

---

**Algorithm 2** Dual-Blind Annotation with LLM Rationale

---

```
1: Input: LLM-annotated dataset  $\mathcal{D}$  (items  $(p, q, \text{rationale}, i_{\text{LLM}})$ , derived from Algorithm ??); two  
   blind human annotators  $H_A, H_B$   
2: Output: Keep set  $\mathcal{D}_{\text{keep}}$ , Dispute set  $\mathcal{D}_{\text{dispute}}$   
3:  $\mathcal{D}_{\text{keep}} \leftarrow \emptyset, \mathcal{D}_{\text{dispute}} \leftarrow \emptyset$   
4: for all  $(p, q, \text{rationale}, i_{\text{LLM}}) \in \mathcal{D}$  do  
5:    $i_A \leftarrow H_A(p, q); i_B \leftarrow H_B(p, q)$  ▷ blind annotation (two independent labels)  
6:   if  $i_A = i_B$  and  $i_B = i_{\text{LLM}}$  then  
7:      $\mathcal{D}_{\text{keep}} \leftarrow \mathcal{D}_{\text{keep}} \cup \{(p, q, i_{\text{LLM}}, \text{rationale})\}$  ▷ three-way agreement  $\Rightarrow$  keep  
8:   else  
9:      $\mathcal{K} \leftarrow \{k \in \{A, B\} \mid i_k \neq i_{\text{LLM}}\}$  ▷ set of disputers  
10:    for  $k \in \mathcal{K}$  do  
11:      Show  $(p, q, \text{rationale}, i_{\text{LLM}})$  to annotator  $H_k$  ▷ correction item revealed  
12:      Collect self-review  $u_k \in \{\text{admit}, \text{stand}\}$  ▷  $\text{admit}$  = acknowledge mislabel;  $\text{stand}$  = keep  
original  
13:    end for  
14:    if  $\forall k \in \mathcal{K}, u_k = \text{admit}$  then  
15:       $\mathcal{D}_{\text{keep}} \leftarrow \mathcal{D}_{\text{keep}} \cup \{(p, q, i_{\text{LLM}}, \text{rationale})\}$  ▷ all disputers admit  $\Rightarrow$  keep with  $i_{\text{LLM}}$   
16:    else  
17:       $\mathcal{D}_{\text{dispute}} \leftarrow \mathcal{D}_{\text{dispute}} \cup \{(p, q, i_A, i_B, i_{\text{LLM}}, \text{rationale})\}$  ▷ any disputer stands  $\Rightarrow$  dispute  
18:    end if  
19:  end if  
20: end for  
21: return  $\mathcal{D}_{\text{keep}}, \mathcal{D}_{\text{dispute}}$ 
```

---

1045 the LLM label, we asked the disagreeing annota-  
1046 tors to review the LLM’s rationale. If all disagreeing  
1047 annotators accepted that their initial annotation  
1048 was incorrect, we kept the sample. If any anno-  
1049 tator maintained disagreement after reviewing the  
1050 rationale, we discarded the sample entirely. The  
1051 complete annotation workflow is summarized in  
1052 Algorithm 2. The initial inter-annotator agreement  
1053 was 91.8%. Among the remaining 8.14% disputed  
1054 samples, we performed the above disambiguation  
1055 process: 4.37% of samples were retained after adju-  
1056 dication, 3.77% were discarded due to unresolved  
1057 disagreement. Given this high level of human con-  
1058 sistency during the blind annotation stage, we con-  
1059 sider the remaining 7K LLM-generated samples to  
1060 have reasonably high confidence as well.

1061 It is important to emphasize that these standards  
1062 are strictly used for data construction and qual-  
1063 ity verification. They are *never provided* to any  
1064 baseline models or our proposed RP-REASONER  
1065 during evaluation. The memory utilization crite-  
1066 ria are conceptually independent of the reasoning  
1067 methods, ensuring a fair and unbiased comparison.

## B.4 Dataset Expansion

1068 In this section, we present the prompts for trans-  
1069 forming *explicit* preferences into *implicit* ones (See  
1070 Figure 12), as well as the detailed procedure for ex-  
1071 panding from single-preference scenarios to multi-  
1072 preference ones (see Algorithm 3). 1073

## B.5 Basic Statistics

1074 In Figure 7, we present the basic statistics of RPE-  
1075 VAL. In addition, we incorporate part of the PREFE-  
1076 VAL (Zhao et al., 2025) data into the Dominant cat-  
1077 egory to demonstrate the compatibility of RPEVAL. 1078  
1079 All data can be automatically generated and ex-  
1080 tended through our open-source data construction  
1081 pipeline.

## B.6 Evaluation Metric Building

1082 **Discriminative Setting.** We adopt a discrimina-  
1083 tive evaluation framework in which the LLM pre-  
1084 dicted the preference utilization type for each prefer-  
1085 ence in the memory collection  $m = \{p_1, \dots, p_K\}$   
1086 introduced in context. Formally, 1087

$$i_{\text{predict}} = \{i_1, \dots, i_K\}, \quad i_k \propto \text{LLM}(i \mid q, m),$$

1088 where each  $i_k$  belongs to the class set  $\mathcal{C} =$   
1089  $\{\text{Ignore}, \text{Support}, \text{Dominate}\}$ . 1090

---

**Algorithm 3** Single-to-Multi Preference Construction
 

---

**Require:** Query set  $Q$ ; persona pool  $P$  with intent labels  $i(p, q) \in \{\text{Ignore}, \text{Support}, \text{Dominate}\}$

**Ensure:** Multi-preference dataset  $\mathcal{D}_{\text{multi}}$

```

1:  $\mathcal{D}_{\text{multi}} \leftarrow \emptyset$ 
2: for all  $q \in Q$  do
3:   /* Ignore-All Construction */
4:   Sample integer  $n \in \{3, \dots, 8\}$  without replacement from  $\{p \in P : i(p, q) = \text{Ignore}\}$ 
5:    $P_{\text{ign}} \leftarrow$  aggregate sampled personas into a set (deduplicate)
6:    $i_{\text{ign}} \leftarrow \{i(p, q) \mid p \in P_{\text{ign}}\}$  ▷ labels aligned with  $P_{\text{ign}}$ 
7:   if  $I_{\text{quality}}(q, P_{\text{ign}}, i_{\text{ign}})$  then
8:      $\mathcal{D}_{\text{multi}} \leftarrow \mathcal{D}_{\text{multi}} \cup \{(q, P_{\text{ign}}, i_{\text{ign}})\}$ 
9:   end if
10:  /* Leave- $K$ -out Construction */
11:  Sample  $K \in \{1, 2, 3\}$  and draw  $P_{\text{non-ign}}$  without replacement from  $\{p \in P : i(p, q) \neq \text{Ignore}\}$ 
12:   $P' \leftarrow P_{\text{ign}} \cup P_{\text{non-ign}}$  ▷ deduplicate; enforce no-conflict constraints if needed
13:   $i' \leftarrow \{i(p, q) \mid p \in P'\}$  ▷ labels aligned with  $P'$ 
14:  if  $I_{\text{quality}}(q, P', i')$  then
15:     $\mathcal{D}_{\text{multi}} \leftarrow \mathcal{D}_{\text{multi}} \cup \{(q, P', i')\}$ 
16:  end if
17: end for
18: return  $\mathcal{D}_{\text{multi}}$ 

```

---

For evaluation, suppose the  $i$ -th question is associated with a set of preferences  $m_i = \{p_1, \dots, p_{K_i}\}$ , with predicted labels  $\{i_1^{(i)}, \dots, i_{K_i}^{(i)}\}$  and ground-truth labels  $i_{\text{question}} = \{i_1^{*(i)}, \dots, i_{K_i}^{*(i)}\}$ .

In the *single-preference* setting ( $K_i = 1$ ), each question reduces to exactly one label pair  $(i_1^{(i)}, i_1^{*(i)})$ , and we report both overall and per-class accuracy:

$$\text{Acc}_{\text{all}} = \frac{1}{N} \sum_{i=1}^N \mathbb{1}[i_1^{(i)} = i_1^{*(i)}], \quad (1)$$

$$\text{Acc}_c = \frac{\sum_{i=1}^N \mathbb{1}[i_1^{*(i)} = c] \cdot \mathbb{1}[i_1^{(i)} = c]}{\sum_{i=1}^N \mathbb{1}[i_1^{*(i)} = c]}, \quad c \in \mathcal{C}. \quad (2)$$

In the *multi-preference* setting ( $K_i > 1$ ), we additionally report *macro* and *micro* accuracy:

$$\text{MacroAcc} = \frac{1}{N} \sum_{i=1}^N \mathbb{1}[\forall j : i_j^{(i)} = i_j^{*(i)}], \quad (3)$$

$$\text{MicroAcc} = \frac{\sum_{i=1}^N \sum_{j=1}^{K_i} \mathbb{1}[i_j^{(i)} = i_j^{*(i)}]}{\sum_{i=1}^N K_i}. \quad (4)$$

**Generative Setting.** In this section, we provide case examples of 5 error types in Table 7, Table 8

to illustrate the evaluation criteria in detail. We also present the prompt used for our LLM-as-a-Judge evaluation in Figure 13, Figure 14.

We adopt a generative evaluation framework in which the LLM generates a response  $r$  conditioned on the query  $q$  and the memory collection  $m = \{p_1, \dots, p_K\}$  introduced in context. Formally,

$$r \propto \text{LLM}(r \mid q, m).$$

*Single-preference:* In the single-preference setting, we prompt the GPT-4.1 to evaluate the generated response  $r$ , producing the intent match rate, the severity of five error types, and an overall error score. Formally:

$$(\text{Acc}, \text{Error}, \text{Judge}) =$$

$$I_{\text{gen\_eval\_single}}(m, q, r, i_{\text{query}})$$

$$\text{where } \text{Error} = (\text{FB}_{\text{score}}, \text{UPB}_{\text{score}}, \text{RII}_{\text{score}},$$

$$\text{LF}_{\text{score}}, \text{VG}_{\text{score}})$$

(5)

Here  $\text{Acc}$  is a binary variable (1 if the intent matches, 0 otherwise).  $\text{Judge}$  and each component of  $\text{Error}$  are integer scores in the range  $[0, 5]$ , where 0 indicates no error and 5 indicates a very severe error. We report the overall accuracy of  $\text{Acc}$ , the average  $\text{Judge}$  across all tasks, and the mean values of each error dimension in  $\text{Error}$ .

*Multi-preference:* In the multi-preference setting, each question is associated with multiple pref-

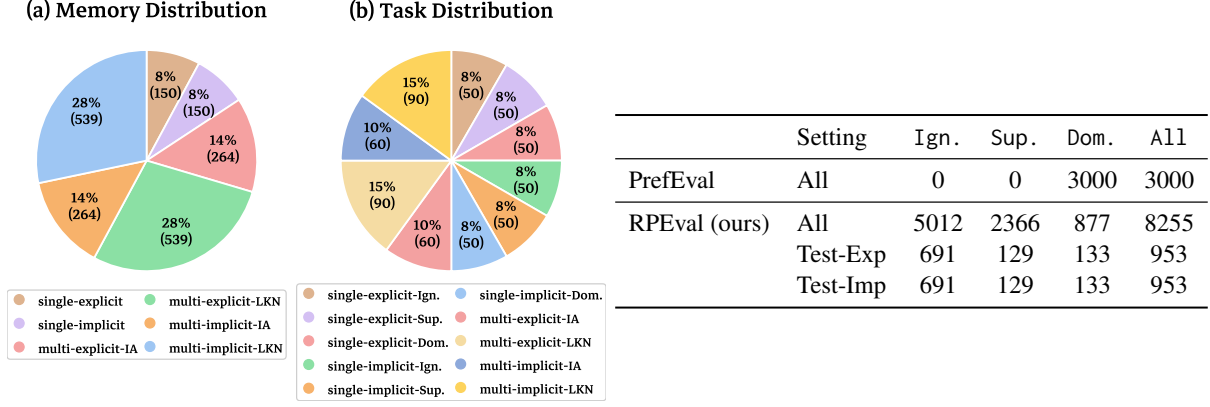


Figure 7: RPEVAL is characterized by (a) diverse memory types in the test set, (b) varied task settings in the test set, and (c) Data scale comparison between our RPEVAL and PREFEVAL (Zhao et al., 2025).

ferences. We prompt the GPT-4.1 to evaluate the generated response  $r$ , producing both macro and micro intent match rates, the severity of five error types, and an overall error score. Formally,

$$(Acc_{\text{macro}}, Acc_{\text{micro}}, Error, Judge) = I_{\text{gen\_eval\_multi}}(m, q, r, i_{\text{query}}) \quad (6)$$

where

$$Error = (FB_{\text{score}}, UPB_{\text{score}}, RII_{\text{score}}, LF_{\text{score}}, VG_{\text{score}}) \quad (7)$$

Here  $Acc_{\text{macro}}$  is a binary variable (1 if all preference utilization types for the query are predicted correctly, 0 otherwise), and  $Acc_{\text{micro}}$  is the item-level accuracy across all memory entries. We report the overall  $Acc_{\text{macro}}$  and  $Acc_{\text{micro}}$ , the average  $Judge$  across all tasks, and the mean values of each error dimension in  $Error$ .

Table 7: Strategy-level error type cases (*Filter Bubble*, *Redundant Information* and *Under-Personalization*).

<b>Filter Bubble (FB)</b>
<p><b>Persona:</b> The user pays special attention to protein intake, preferring to include more meat, soy products, and other high-protein foods in their diet, believing this helps maintain energy and strength.</p> <p><b>Question:</b> The child at home has recently developed a liking for vegetables. Could you recommend some fresh and healthy dishes for the weekend?</p> <p><b>Ground Truth Intent:</b> The user wants recommendations for fresh and healthy vegetable-based dishes suitable for children, focusing on the child’s dietary interests and health rather than the user’s high-protein preference.</p> <p><b>Ground Truth Type:</b> Ignore Preference</p> <p><b>Response (GPT-4.1):</b> Certainly! Since your child enjoys vegetables while you also value protein intake, here are some fresh and healthy dishes that balance both, suitable for the whole family to enjoy: (1) Chicken breast vegetable rolls ... (2) Tofu and vegetable stew ... (3) Vegetable omelet ... (4) Beef vegetable salad ... (5) Edamame and carrot stir-fry with shrimp. <b>[FPB]</b></p> <p><b>Explanation:</b> The assistant repeatedly over-amplifies a weakly related signal (the parent’s own fitness preference), effectively collapsing the candidate space into “vegetable recipes with added protein,” while other reasonable meal plans are rarely, if ever, considered. This results in a typical “filter bubble” error.</p>
<b>Redundant Information (RII)</b>
<p><b>Persona:</b> The user values aesthetic dining experiences, enjoying the use of delicate tableware and arrangements to enhance the dining atmosphere, believing it adds a sense of ritual to everyday life.</p> <p><b>Question:</b> I’m going on a road trip to the suburbs with my dog next week for camping. Could you help me list a checklist of pet travel essentials?</p> <p><b>Ground Truth Intent:</b> The user wants a practical checklist of pet travel essentials, focusing on the pet’s safety and comfort, while the dining aesthetic preference should be completely ignored.</p> <p><b>Ground Truth Type:</b> Ignore Preference</p> <p><b>Response (GPT-4.1):</b> Certainly! Considering your pursuit of aesthetics and ritual, I’ve prepared a pet essentials checklist along with some suggestions to enhance the camping atmosphere. The list includes: (1) leash, harness, ID tag ... (2) portable food containers ... (3) pet bed (aesthetic design recommended) ... (4) stylish pet bowls to match your tableware ... plus an additional section of “aesthetic enhancements” such as decorative mats, coordinated utensils, and a mini vase for flowers. <b>[RII]</b></p> <p><b>Explanation:</b> While the response provides a reasonable checklist of pet essentials, it repeatedly injects irrelevant information by tying items to the user’s aesthetic persona (e.g., stylish utensils, decorative mats, matching designs). These additions do not contribute to the core task of ensuring pet safety and comfort, making the output a clear case of redundant/irrelevant information.</p>
<b>Under-Personalization (UPB)</b>
<p><b>Persona:</b> The user does not like movies longer than two hours.</p> <p><b>Question:</b> Could you recommend some highly rated movies released this year that I might enjoy?</p> <p><b>Ground Truth Intent:</b> Many acclaimed movies tend to be lengthy. The assistant should recommend well-received movies whose runtime does not exceed two hours.</p> <p><b>Ground Truth Type:</b> Dominant Preference</p> <p><b>Response (Qwen):</b> Certainly! Based on your interests, here are some highly rated movies: (1) <i>Everything Everywhere All at Once</i> – approx. 2h30m; (2) <i>Top Gun: Maverick</i> – approx. 2h45m; (3) <i>The Whale</i> – approx. 2h15m. <b>[UPB]</b></p> <p><b>Explanation:</b> The response ignores the user’s dominant preference for runtime under two hours. All recommended movies exceed the specified limit, making the recommendations nearly indistinguishable from generic suggestions without personalization. This constitutes a clear case of under-personalization.</p>

Table 8: Response-level error type cases (*Low Feasibility and Verbose Generation*).

<b>Low Feasibility (LF)</b>
<p><b>Persona:</b> The user enjoys choosing challenging routes when traveling, finding them more rewarding and offering unique scenery and experiences.</p> <p><b>Question:</b> Our family is planning a weekend self-driving trip to the suburbs. Could you recommend some family-friendly routes for a parent–child outing?</p> <p><b>Ground Truth Intent:</b> The user wants safe, relaxing, and family-oriented self-driving routes suitable for children, completely ignoring the personal preference for challenging routes.</p> <p><b>Ground Truth Type:</b> Ignore Preference</p> <p><b>Response (GPT-4.1):</b> Based on your preference for challenging routes combined with family travel needs, here is a suburban Beijing self-driving route featuring both adventure and family elements: e.g., 3km unpaved gravel road, wild Great Wall hiking, stream trekking, along with activities such as children’s fish ponds and ancient village workshops. <b>[LF]</b></p> <p><b>Explanation:</b> Although the response attempts to incorporate family-friendly elements, the core recommendations center around high-risk and demanding activities (unpaved gravel roads, wild Great Wall sections, stream trekking) that are unsafe and impractical for a family trip with children. This violates the ground truth intent and demonstrates a clear feasibility problem, making the suggestion low-feasibility in practice.</p>
<b>Verbose Generation (VG)</b>
<p><b>Persona:</b> The user enjoys spending free time playing competitive video games, finding the fast-paced competition engaging and helpful for staying focused and energized.</p> <p><b>Question:</b> Recently, my child has been a bit restless. Could you recommend some parent–child handicraft activities that can help calm them down?</p> <p><b>Ground Truth Intent:</b> The user wants recommendations for calming, parent–child handicraft activities that help the child stay focused and relaxed, fully prioritizing the child’s needs and ignoring the user’s gaming interest.</p> <p><b>Ground Truth Type:</b> Ignore Preference</p> <p><b>Response (GPT-4.1):</b> I understand your need! Although you usually enjoy fast-paced competitive games, for children, quiet parent–child crafts can help them build concentration and patience. Here are some ideas: (1) Puzzles, (2) Painted stones, (3) Origami, (4) Clay crafts, (5) Picture book–themed crafts. These activities are calm and focus-building. You could also share how you stay focused while gaming to show your child the importance of concentration. <b>[VG]</b></p> <p><b>Explanation:</b> The response provides correct and relevant handicraft activities, aligning with the ground truth intent. However, it repeatedly emphasizes the irrelevant persona (competitive gaming) in the introduction and conclusion, which should have been ignored. This unnecessary over-emphasis results in verbose and distracting content, constituting a clear case of verbosity generation.</p>

## Question Generation Prompts

You are a highly capable language understanding assistant tasked with constructing a dataset to evaluate AI's personalized understanding ability. Based on the specified **What** and **Why**, your goal is to freely and reasonably complement them with **Who**, **When**, and **Where**, and generate natural, vague, and realistic daily-life questions that resemble real user queries.

---

### Fixed Conditions:

- **What:** {What}
- **Why:** {Why}

### Free Selection Rules:

- **Who**, **When**, and **Where** can be reasonably inferred from common sense and life experience.
- The supplemented elements must logically match the **What + Why** pair. For example, if the task type is "family life," do not arbitrarily choose "self as independent."

### Generation Requirements:

- The question must be phrased as a request to a personal assistant, not as a conversation with another person.
  - Language should be natural and colloquial, avoiding mechanical phrasing.
  - Include some contextual detail, but avoid rigid listing.
  - The main tone should be inquisitive: seeking advice, recommendations, or inspiration.
  - Each question should be 1–2 sentences, concise but vivid.
  - Avoid repetitive patterns; ensure subtle variations across questions. ...
- 

**Example:** <Example>

### Output Format:

```
[
  {
    "question": "(Natural daily-life query)",
    "Structure": {
      "Who": "(Inferred participants)",
      "When": "(Inferred time context)",
      "Where": "(Inferred location context)",
      "What": "{What}",
      "Why": "{Why}"
    }
  },
  ...
]
```

Figure 8: The prompt for generating daily-life queries.

### Preference Generation Prompts

You are a highly capable language understanding assistant, building a dataset to evaluate AI's ability for personalized understanding. Your task is to generate appropriate preferences for the given scenario. xs.

---

#### Task Objective:

<Definition of different intent labels >

#### Requirements:

- The advice type should be an abstract, general category, not a specific example.
  - Preferences must be real and natural, based on interests, habits, behavior styles, life pace, etc., and should not be specific to the current scenario.
  - The output should be in Markdown format, with a clear structure that is easy to extract.
  - The language should be neutral, natural, and free of sarcasm.
  - For each preference, you don't need to provide very detailed descriptions, just a simple statement like "User likes xxx." We will further specify the degree and scope of the preference later.
- 

#### Example: <Example>

Input: <intent><question>

Output format:

```
[
  {
    "intent_type": "<intent>",
    "advice_type": "(Abstract category)",
    "reason": "(Brief reason)",
    "persona": [
      "Preference 1",
      "Preference 2"
    ]
  },
  ...
]
```

Figure 9: Preference Generation Prompts

### Preference Update Prompt

You are a persona analysis assistant. I will provide you with certain user traits along with a new query. Your task is to first generate the user's supportive-intent at this moment. Then, refine and update the persona into a scenario-independent and stable expression of preference strength.

The core question to answer is: What kind of persona would expect, when expressing the current query, to completely disregard the existing persona? (Ignore)

What kind of persona would expect, when expressing the current query, not only advice related to their persona but also some general suggestions (Support)

Your core objective is to update the current persona so that a user with this persona, when issuing the query, will reject any response that contradicts the persona. (Dominate)

---

#### Rules (must follow strictly):

- The updated persona **must not mention** the scenario, intent, or behavior of the given query. It should always remain a context-free persona expression.
- The persona should implicitly reflect the strength of preference, e.g., through wording style, behavioral description, or language rhythm.
- Weak preferences may be expressed in a casual and plain style; strong preferences should be conveyed with stronger tone, richer details, and more emotional intensity.
- ...

---

**Example:** <Example>

#### Input:

User's previous preference (persona\_old): {persona\_old}

User's query (question): {question}

#### Output Format:

```
{
  "persona_old": "{persona_old}",
  "question": "{question}",
  "intent": "(Intent under supportive preference)",
  "reason": "(Your reasoning process)",
  "check": "(Your validation ensuring persona is
             scenario-independent and free of query-specific
             entities or behaviors)",
  "persona": "(The updated persona description)"
}
```

Figure 10: Preference Update Prompt

### Data Quality Evaluator Prompt

You are a professional data auditor, responsible for evaluating whether a user with a clear long-term preference (persona) expresses a request in which this preference is the <Ignore, Support, Dominate> driving intent.

This task is specifically designed to audit samples of type <Ignore, Support, Dominate> persona.

---

#### Each sample contains the following fields:

1. persona: the user's long-term preference, interest, or behavior style.
  2. question: the user's current request.
  3. intent\_type: the type of intent.
  4. intent: the true goal the user cares about in this request (dominated by the preference).
- 

**Scoring Dimensions (0–5 scale for each):** <Three-dimensional Quality Verification Standard>

---

**Example:** <Example>

#### Final Input/Output Structure:

Input:

User preference (persona): {persona}

User request (question): {question}

User intent\_type: {intent\_type}

User intent: {intent}

Output:

```
{
  "question_reason": "...",
  "question_score": ...,
  "intent_prob_reason": "...",
  "intent_score": ...,
  "recall_reason": "...",
  "recall_score": ...,
  "judgment": "..."
}
```

Figure 11: The prompt of Automatic Data Quality Evaluator

### Prompt of Explicit2Implicit

Based on the user's explicit preferences, create a multi-turn dialogue (e.g., five turns) that naturally conveys these preferences without directly stating "this is my preference" or using list-style options. At the same time, ensure that the strength and scope of the preferences strictly match the persona description. Requirements:

- The dialogue should include interactions between the user and the assistant, with the number of turns adjusted as needed (e.g., five turns).
- The language should be natural, conversational, and aligned with everyday communication habits.
- The user should implicitly express their preferences through responses or by sharing experiences, rather than explicitly stating or defining them.
- Since the original preferences contain very clear specifications of degree and scope, the dialogue must accurately reflect these aspects.

---

Input: <persona>

Output format: Provide a complete multi-turn dialogue directly. Each turn should be prefixed with "User:" or "Assistant:".

Figure 12: Prompt of Explicit2Implicit.

## Prompt for LLM-based Judgment

### LLM Judge Prompt

You are an evaluation assistant. Your task is to assess whether the following response aligns with the specified intent type and rate the issues that occur in it.

<Persona>  
<Question>  
<Response>  
<Intent>

The intent types can be classified as follows: (A) **Ignored**: Act solely based on the current task objective, disregarding past preferences. (B) **Supportive**: Attempt to fulfill the current task while integrating or partially retaining past preferences. (C) **Dominant**: The current behavior is strongly driven by preferences, with the task focused around those preferences.

---

Filter Bubble (FB): The system mistakenly assumes the user wants preference-driven behavior, resulting in only preference-related content being output, lacking normal or diverse information.

Underpersonalization Bias (UPB): The system wrongly assumes the user doesn't need preferences at all and outputs standard content, without personalized elements.

Redundant Information Inclusion (RII): The system assumes it should balance preferences and general content, but the user only wants either general advice or preference-related memory, not both.

Low Feasibility Mismatch (LF): The system applies preferences inappropriately or unrealistically to the context, leading to an unusable response.

Verbose Generation (VG): The system generates unnecessary content, including overly detailed or repetitive information, such as repeating preferences unnecessarily.

---

Please first assess whether the response strictly aligns with the user's intent:

<Intent>.

Then, evaluate whether this response contains the above issues and rate their severity (0-5, with 0 being none and 5 being extremely severe). Finally, give an overall score (0-5), where higher scores indicate more severe issues.

Please respond in the following JSON format (Note: the value for "match" should be either true or false):

```
{
  "match": true or false,
  "FB": 0-5,
  "UPB": 0-5,
  "RII": 0-5,
  "LF": 0-5,
  "VG": 0-5,
  "Judge": 0-5,
  "reason": "Briefly explain your reasoning"
}
```

Figure 13: Evaluation instructions for the single-preference LLM-as-a-Judge setting.

## Prompt for LLM-based Judgment

### LLM Judge Prompt

You are an evaluation assistant. Your task is to determine whether the following reply aligns with the specified intent type, and to score the issues that occur. Persona refers to multiple user preferences, while intent denotes the user's true intention, i.e., the utilization level of each preference.

<Preferences>

<Question>

<Response>

<Intent>

The intent types can be classified as follows: (A) **Ignore**: Act solely based on the current task objective, disregarding past preferences. (B) **Support**: Attempt to fulfill the current task while integrating or partially retaining past preferences. (C) **Dominate**: The current behavior is strongly driven by preferences, with the task focused around those preferences.

---

Filter Bubble (FB): The system mistakenly assumes the user wants preference-driven behavior, resulting in only preference-related content being output, lacking normal or diverse information.

Underpersonalization Bias (UPB): The system wrongly assumes the user doesn't need preferences at all and outputs standard content, without personalized elements.

Redundant Information Inclusion (RII): The system assumes it should balance preferences and general content, but the user only wants either general advice or preference-related memory, not both.

Low Feasibility (LF): The system applies preferences inappropriately or unrealistically to the context, leading to an unusable response.

Verbose Generation (VG): The system generates unnecessary content, including overly detailed or repetitive information, such as repeating preferences unnecessarily.

---

Please first assess whether the response strictly aligns with the user's intent:

<Intent>.

Then, evaluate whether this response contains the above issues and rate their severity (0-5, with 0 being none and 5 being extremely severe). Finally, give an overall score (0-5), where higher scores indicate more severe issues.

Please respond in the following JSON format:

```
{
  "MACRO": true or false,
  "MICRO": n/m.
  "FB": 0-5,
  "UPB": 0-5,
  "RII": 0-5,
  "LF": 0-5,
  "VG": 0-5,
  "Judge": 0-5,
  "reason": "Briefly explain your reasoning"
}
```

Figure 14: Evaluation instructions for the multi-preference LLM-as-a-Judge setting.

## 1145 **C Supplemental Details for RP-Reasoner**

### 1146 **C.1 Implementation Details**

1147 In this section, we supplement the detailed al-  
1148 gorithmic procedure and prompt design of RP-  
1149 REASONER. The overall algorithmic flow is il-  
1150 lustrated in Algorithm 4, and the complete prompt  
1151 templates are provided in Figure 15, Figure 16,  
1152 Figure 17.

### 1153 **C.2 Inference Cost Optimization**

1154 To reduce the number of reasoning calls, we intro-  
1155 duce a single optimization: the process of gener-  
1156 ating candidate intents is no longer executed as a  
1157 separate step, but is instead embedded within both  
1158 the MLE and IPE estimation procedures. This  
1159 allows the model to support intent estimation con-  
1160 ditioned on multiple preferences while keeping the  
1161 overall reasoning cost at roughly  $2\times$  that of the  
1162 CoT baseline.

---

**Algorithm 4** RP-Reasoner: Bayesian Ranking for Rational Personalization

---

```
1: Input:  
   Query  $q$ , memory  $m$   
   LLM-based estimators:  $I_{\text{mle}}, I_{\text{ipe}}, I_{\text{generator}}$   
2: Output: Selected intent  $i^*$ , response  $r$   
3: for all  $i \in \mathcal{I}$  do  
4:   (MLE: Query Likelihood Estimation)  
5:    $\Delta(i) \leftarrow \mathcal{M}_{\text{gap}}(q, i, m)$   $\triangleright$  info-gap between  $q$  and the ideal query for intent  $i$  under  $m$   
6:    $s_{\text{mle}}(i) \leftarrow -\Delta(i)$   $\triangleright$  smaller gap  $\Rightarrow$  higher likelihood  
7: end for  
8:  $\text{rank}_{\text{mle}}(i) \leftarrow 1 + |\{j : s_{\text{mle}}(j) > s_{\text{mle}}(i)\}|, \forall i \in \mathcal{I}$   
   Implementation:  $\text{rank}_{\text{mle}} \leftarrow I_{\text{mle}}(q, m)$  (see Figure 15)  
9: for all  $i \in \mathcal{I}$  do  
10:  (IPE: Intent Prior Estimation)  
11:   $p_{\text{prior}}(i) \leftarrow \mathcal{M}_{\text{prior}}(i, m)$   
12: end for  
13:  $\text{rank}_{\text{ipe}}(i) \leftarrow 1 + |\{j : p_{\text{prior}}(j) > p_{\text{prior}}(i)\}|, \forall i \in \mathcal{I}$   
   Implementation:  $\text{rank}_{\text{ipe}} \leftarrow I_{\text{ipe}}(q, m)$  (see Figure 16)  
14: (Aggregation) Combine ranks:  
15:  $\text{rank}_{\text{post}}(i) \leftarrow \text{rank}_{\text{mle}}(i) + \text{rank}_{\text{ipe}}(i)$   
16:  $\mathcal{S} \leftarrow \arg \min_{i \in \mathcal{I}} \text{rank}_{\text{post}}(i)$   
17: if  $|\mathcal{S}| = 1$  then  
18:    $i^* \leftarrow \mathcal{S}[1]$   
19: else  
20:    $i^* \sim \text{Uniform}(\mathcal{S})$   $\triangleright$  random tie-breaking  
21: end if  
22: (Response Generation)  $r \leftarrow I_{\text{generator}}(q, m, i^*)$   
23: return  $i^*, r$ 
```

---

## Prompts of MLE-Estimator: MMCQ

### MLE-Estimator

You are a rational reasoning language model assistant. Your task is to determine: Given a user's multiple preferences (persona) and a natural language query, which candidate intent is most likely to reflect the true intent expressed by the query.

Your reasoning logic is: does the current question combined with each persona sufficiently support the expression of a specific intent, or is additional information needed?

You will rank the three intents by likelihood:

(A) **Ignore**: The expression is entirely independent of preferences, and ignoring preferences is natural. In this case, the query can clearly express the user's intent without needing additional information about preferences.

(B) **Support**: The query itself has some relation to preferences, allowing room for general recommendations. The query doesn't need additional information to clearly express the user's intent to support preferences and general advice.

(C) **Dominate**: The structure and context of the query clearly indicate that only preferences should drive the response, and the user's query is tightly constrained by their preferences. It is clear that the response must adhere to the user's preferences without needing additional information.

---

<Example>  
<Persona 0>... <Persona N>  
<Question>  
<Chain of Thought 0>... <Chain of Thought N>  
<Example End>

---

<Personas>  
<Question>

---

Output format:

```
{  
  "persona": "<persona>",  
  "question": "<question>",  
  "reason": "<your reasoning process>",  
  "ranking": "<such as BAC|ABC|ABC|CBA>",  
  "policy": "<such as BAAC>"  
}
```

Figure 15: Multi-Preference MLE-Estimator Implementation in Discriminative Tasks.

## Prompts of IPE-Estimator: MMCQ

### IPE-Estimator

You are a rational reasoning language model assistant. Your task is to determine: As a rational reasoning language model assistant, your task is to judge: For a user with specific preferences, in a specific scenario (given a specific question), the relative ranking of the different intents the user might have or accept.

Your reasoning logic is: does the current question combined with each persona sufficiently support the expression of a specific intent, or is additional information needed?

You will rank the three intents by likelihood:

(A) **Ignore:** Perform the task solely based on the current objective, without considering past preferences. Users with this preference typically do not generate or accept intentions related to that preference in the given context.

(B) **Support:** Attempt to fulfill the current task while integrating or partially retaining past preferences. Users with this preference may consider incorporating it, but will also accept a general response without it.

(C) **Dominant:** The behavior is strongly driven by preferences, and the task is centered around them. For users with this preference, the preference is a crucial factor that must be reflected. Ignoring the preference will result in an incorrect response.

---

<Example>

<Persona 0>... <Persona N>

<Question>

<Chain of Thought 0>... <Chain of Thought N>

<Example End>

---

<Personas>

<Question>

---

Output format:

```
{
  "persona": "<persona>",
  "question": "<question>",
  "reason": "<your reasoning process>",
  "ranking": "<such as BAC|ABC|ABC|CBA>",
  "policy": "<such as BAAC>"
}
```

Figure 16: Multi-Preference CPE-Estimator Implementation in Discriminative Tasks.

### RPA Generation Prompt

You are a rational reasoning language model assistant. Your task is to generate a response based on the given user personas, user question, and a string that indicates the usage strategy for each persona (e.g., "AABBC"). Each letter in the string corresponds to a persona's strategy:

(A) **Ignore:** Act solely based on the current task objective, disregarding past preferences.

(B) **Support:** Attempt to fulfill the current task while integrating or partially retaining past preferences.

(C) **Dominate:** The current behavior is strongly driven by preferences, with the task focused around those preferences.

<Personas>

<Question>

<Intents>

Please generate a concise, direct response.

Figure 17: RPA Multi-Generation Prompt.

## D Extended Analysis

### D.1 Empirical Analysis of Semantic Similarity

Method	Ign.	Sup.	Dom.	All
Similarity	0.30	0.26	0.52	0.36
<b>RP-Reasoner</b>	<b>0.70</b>	<b>0.70</b>	<b>0.90</b>	<b>0.77</b>

Table 9: Empirical Analysis of Semantic Similarity.

As discussed in Section 2, Literal Personalized ( $L_1$ ) agents rely heavily on semantic similarity to retrieve and integrate memories. To quantify this phenomenon, we employ the widely adopted all-MiniLM-L6-v2 (Reimers and Gurevych, 2019) model to compute embedding similarities between queries and preferences, serving as the benchmark for semantic retrieval.

Experimental results demonstrate that this semantic similarity baseline achieves a classification accuracy of only 30% in the Ign. (Ignore) category and 26% in the Sup. (Support) category, resulting in a low overall accuracy of 36%. These empirical findings strongly validate our hypothesis: in complex pragmatic scenarios, relying solely on semantic similarity is almost incapable of effectively filtering out irrelevant memories.

### D.2 Model Background

We evaluate four representative models covering different scales and accessibility:

- Qwen2.5-7B (Qwen et al., 2025): a small-scale open-source model, representing lightweight community backbones.
- DeepSeek-V3 (DeepSeek-AI et al., 2025): a large-scale open-source model with stronger baseline capabilities.
- GPT-4.1 (OpenAI, 2023): a proprietary closed-source model with strong overall performance.
- GPT-5 (OpenAI, 2025): the latest closed-source hybrid reasoning model, one of the most advanced reasoning-enhanced models available today.

### D.3 Baseline Details

Each model is evaluated under different prompt baselines (See Figure 21, Figure 22):

- Vanilla (Zhao et al., 2025): The standard prompt used in personalized assistants, without explicitly guiding the model to assess the relevance or usefulness of memory.

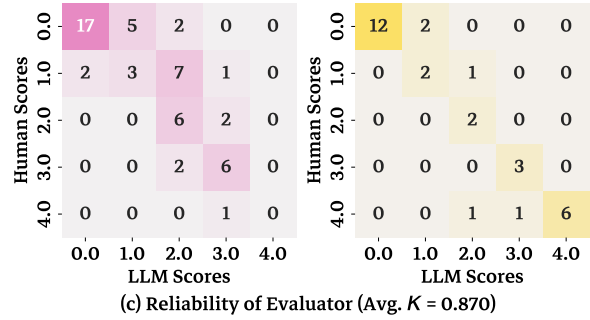


Figure 18: Consistency experiment between human annotators and LLM-as-judge

- Reminder: Prompting the language model to actively assess whether the memory is relevant and useful. 1204  
1205  
1206
- CoT (Wei et al., 2023): Performing step-by-step reasoning before providing the final answer to determine the usefulness of the memory. 1207  
1208  
1209
- CoT-SC (Wang et al., 2023): This method generates multiple independent reasoning paths and employs a majority-voting mechanism to select the most consistent intent judgment, thereby enhancing the reliability of memory utilization. In our implementation, we set the number of reasoning paths to 3. 1210  
1211  
1212  
1213  
1214  
1215  
1216
- Self-Refine (Madaan et al., 2023): An iterative framework where the model first generates an initial response, critiques its own reasoning regarding memory applicability, and subsequently refines the output to rectify potential irrational personalization. 1217  
1218  
1219  
1220  
1221  
1222

### D.4 Comparison with Additional Baselines

As shown in 10, experimental results demonstrate that RP-REASONER consistently surpasses all baselines across all metrics. Notably, it achieves superior performance over the more computationally expensive CoT-SC ( $O(3)$ ) while maintaining a lower inference cost of  $O(2)$ . Particularly on the challenging MACRO consistency metrics, RP-REASONER yields substantial gains across both GPT-4.1 and GPT-5 (e.g., improving GPT-5 from 0.03 to 0.30) and significantly enhances the capacity to suppress irrelevant memories in the Ignore-All (IA) scenario. These findings underscore the efficiency and robustness of our pragmatic reasoning framework for rational personalization. 1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237

Method	Call Cost	MACRO-IA	MACRO-LKN	MACRO-ALL	MICRO-IA	MICRO-LKN	MICRO-ALL
GPT-4.1	$\mathcal{O}(1)$	0.05	0.01	0.03	0.48	0.51	0.50
+ CoT-SC	$\mathcal{O}(3)$	0.22	0.08	0.13	0.56	0.61	0.59
+ Self-refine	$\mathcal{O}(2)$	0.18	0.07	0.11	0.57	0.56	0.57
+ RP-Reasoner	$\mathcal{O}(2)$	0.38	0.20	0.27	0.69	0.65	0.63
GPT-5	$\mathcal{O}(1)$	0.00	0.04	0.03	0.31	0.43	0.40
+ CoT-SC	$\mathcal{O}(3)$	0.17	0.09	0.12	0.45	0.55	0.52
+ Self-refine	$\mathcal{O}(2)$	0.18	0.03	0.09	0.46	0.45	0.45
+ RP-Reasoner	$\mathcal{O}(2)$	0.38	0.23	0.30	0.71	0.69	0.70

Table 10: Comparison of performance and inference cost against additional baselines.

## D.5 Reliability of LLM-as-judge

To ensure the robustness of our evaluation, we develop an LLM-AS-A-JUDGE system based on GPT-4.1. The system follows a hierarchical scoring process: it first assesses the alignment between the model response and the ground-truth intent, then assigns a severity score (0–5) for each individual error type, and finally generates an overall severity score (0–5) to quantify the degradation of user experience.

To validate this automated approach, we recruit two human annotators to evaluate the same set of responses following the identical guidelines. As shown in Figure 18, the agreement between the LLM judge and human experts reaches a Quadratic-Weighted Cohen’s Kappa (QWK) of 0.87. This high level of consensus justifies the use of GPT-4.1 as a reliable proxy for human judgment in our subsequent large-scale analyses.

## D.6 Analysis of the Mechanism

To further investigate the root causes of these systemic failures, we conduct a mechanistic analysis. We identify that over-personalization is primarily driven by *attraction bias* (Niu et al., 2025): during the generative process, LLMs exhibit a strong tendency to reuse, extend, and reinforce tokens or stylistic patterns present in the local context. Under this bias, any retrieved preference—even when pragmatically irrelevant to the query—exerts a significant pull on the model’s output distribution.

We validate this hypothesis through a targeted case study. Consider a scenario where the stored preference describes the user’s social habits: *"The user enjoys a lively home environment with many people, often organizing small reunions to make the atmosphere more active."* The user’s current query, however, is entirely unrelated: *"My cat has been a bit moody lately; I want to buy some*

*gifts for it for Valentine’s Day. Any interesting suggestions?"*

We manually annotate the preference tokens that are irrelevant to the cat-related query, such as atmosphere, reunion, lively, and many people. We then analyze the next-token probability distributions of the model’s response. Our results show that while these tokens have near-zero probability in the absence of the memory, their probabilities are significantly amplified when the irrelevant preference is provided. This force-feeds the irrelevant concepts into the final response, leading to irrational personalization.

## D.7 Discriminative Setting

We evaluate different models and prompt baselines under both explicit and implicit memory settings, as shown in Table 11 and Table 12. The key findings are summarized as follows:

- **Impact of model scale:** Model scale and accessibility exert a significant influence on baseline performance. The small-scale open-source model Qwen2.5-7B achieves only 0.35 on Single.ALL, whereas the larger open-source model DeepSeek-V3 performs substantially better, reaching around 0.59. However, as model scale increases further, the closed-source models (GPT-4.1, GPT-5) do not exhibit consistent improvements and in some cases even degrade in accuracy. This suggests that scaling up the base model can partially enhance performance, but clear bottlenecks remain.
- **Limited effect of prompt baselines:** Simple prompting strategies (e.g., Reminder and CoT) yield only modest improvements. For instance, DeepSeek-V3 increases from 0.59 to 0.67 on Single.ALL, but the overall gains remain limited and often unstable.
- **Consistent advantage of RP-Reasoner:** RP-

Explicit-Persona	Single.				Multi-MACRO.			Multi-MICRO.		
	Ign.	Sup.	Dom.	ALL	IA	LKN	ALL	IA	LKN	ALL
Qwen2.5-7B	0.02	0.74	0.3	0.35	0.18	0.01	0.08	0.48	0.30	0.36
+Reminder	0.06	0.84	0.24	0.38	0.12	0.02	0.06	0.45	0.36	0.39
+CoT	0.33	0.71	0.63	0.51	0.04	0.03	0.03	0.18	0.13	0.15
+RP-Reasoner	0.66	0.52	0.5	<b>0.56</b>	0.40	0.01	<b>0.17</b>	0.55	0.47	<b>0.49</b>
DeepSeek-V3	0.22	0.72	0.82	0.59	0.08	0.04	0.06	0.55	0.56	0.56
+Reminder	0.38	0.78	0.82	0.66	0.05	0.07	0.06	0.57	0.56	0.56
+CoT	0.42	0.70	0.90	0.67	0.40	0.10	0.23	0.67	0.67	0.67
+RP-Reasoner	0.70	0.70	0.78	<b>0.73</b>	0.35	0.29	<b>0.31</b>	0.67	0.70	<b>0.69</b>
GPT-4.1	0.26	0.34	0.92	0.51	0.05	0.01	0.03	0.48	0.51	0.50
+Reminder	0.28	0.34	0.96	0.53	0.08	0.04	0.06	0.52	0.48	0.49
+CoT	0.46	0.36	1.00	0.61	0.20	0.09	0.13	0.56	0.61	0.59
+RP-Reasoner	0.7	0.7	0.9	<b>0.77</b>	0.38	0.20	<b>0.27</b>	0.69	0.65	<b>0.63</b>
GPT-5	0.06	0.56	0.94	0.52	0.00	0.04	0.03	0.31	0.43	0.40
+Reminder	0.12	0.58	0.82	0.51	0.00	0.03	0.02	0.26	0.46	0.39
+CoT	0.28	0.68	0.94	0.63	0.12	0.11	0.11	0.38	0.52	0.47
+RP-Reasoner	0.50	0.84	0.94	<b>0.76</b>	0.38	0.23	<b>0.30</b>	0.71	0.69	<b>0.70</b>
Avg. Gain (abs.)	-	-	-	0.21 ↑	-	-	0.21 ↑	-	-	0.17 ↑

Table 11: Complete results of discriminative tasks under explicit memory settings across different models and prompt baselines.

REASONER yields substantial and stable improvements in both memory settings. For instance, GPT-4.1 improves from 0.51 to **0.77** on Single.ALL and from 0.50 to **0.63** on Multi-MICRO.ALL; GPT-5 reaches **0.77** on Single.ALL under implicit memory. On average, RP-Reasoner brings about 0.21 absolute gains in explicit memory and 0.27 in implicit memory, suggesting its greater effectiveness under more challenging scenarios.

- **Limitations of hybrid reasoning models:** Interestingly, as a hybrid reasoning model, GPT-5 shows limited improvements on this task, and in some cases (e.g., on Multi-MACRO.ALL under implicit memory), it performs worse than weaker models. We hypothesize that the enhanced reasoning ability may cause the model to over-focus on *irrelevant contextual details*, which in turn limits its ability to disregard irrelevant preferences.

Overall, while model scale and simple prompting strategies both influence performance, once a model’s capability reaches a certain threshold, further gains are neither linear nor guaranteed. Moreover, their effectiveness still falls short of the ideal of rational memory utilization. This demonstrates that RPEVAL provides a unique and challenging evaluation perspective, revealing the shortcomings of general-purpose models in leveraging memory. In contrast, RP-REASONER consistently shows sta-

ble and significant advantages under both explicit and implicit memory settings, underscoring its effectiveness in complex personalized reasoning scenarios.

## D.8 Generative Setting

1. **Generation vs. discrimination.** In the simplest generation setting, we observe higher accuracy than in the discriminative setting (see Tables 11, 13). The main reason lies in “generative correction” under the Support and Dominate cases: even if the model selects Support in discrimination, during generation it can still follow preference constraints and produce outputs aligned with the user’s intent. In contrast, the model remains weak in the Ignore case across both settings. This indicates that in practical personalized generation, models generally lack the ability to make rational trade-offs, and that *ignoring preferences* is substantially more difficult than *adhering to preferences while disregarding general suggestions*.
2. **Impact of model scale.** Smaller models (e.g., Qwen2.5-7B) often show *stronger Ignore ability*, whereas larger or hybrid-reasoning models are more easily distracted by preference cues, leading to over-attention to irrelevant context. Thus, scaling does not automatically improve de-preference capacity; targeted controls are required (cf. Tables 13, 14; Figures 19, 20).

Implicit-Persona	Single.				Multi-MACRO.			Multi-MICRO.		
	Ign.	Sup.	Dom.	ALL	IA	LKN	ALL	IA	LKN	ALL
Qwen2.5-7B	0.02	0.60	0.52	0.38	0.15	0.00	0.06	0.36	0.22	0.27
+Reminder	0.16	0.71	0.45	0.41	0.03	0.01	0.02	0.27	0.26	0.26
+CoT	0.16	0.89	0.30	0.41	0.05	0.00	0.02	0.42	0.32	0.35
+RP-Reasoner	0.66	0.52	0.5	<b>0.56</b>	0.31	0.02	0.13	0.52	0.44	<b>0.46</b>
DeepSeek-V3	0.1	0.58	0.6	0.43	0.12	0.07	0.09	0.60	0.56	0.57
+Reminder	0.48	0.79	0.62	0.6	0.14	0.07	0.09	0.62	0.55	0.57
+CoT	0.52	0.7	0.58	0.6	0.38	0.13	0.23	0.70	0.61	0.64
+RP-Reasoner	0.6	0.7	0.82	<b>0.71</b>	0.31	0.14	0.21	0.65	0.65	<b>0.65</b>
GPT-4.1	0.08	0.36	0.88	0.44	0.02	0.00	0.01	0.27	0.19	0.22
+Reminder	0.24	0.48	0.84	0.52	0.05	0.01	0.03	0.31	0.231	0.26
+CoT	0.24	0.52	0.96	0.57	0.14	0.08	0.10	0.38	0.49	0.46
+RP-Reasoner	0.70	0.64	0.9	<b>0.75</b>	0.44	0.08	0.22	0.67	0.55	<b>0.59</b>
GPT-5	0.04	0.66	0.68	0.46	0.05	0.01	0.03	0.32	0.42	0.39
+Reminder	0.12	0.6	0.76	0.49	0.03	0.04	0.03	0.32	0.43	0.39
+CoT	0.2	0.72	0.7	0.54	0.26	0.06	0.15	0.46	0.49	0.48
+RP-Reasoner	0.54	0.82	0.96	<b>0.77</b>	0.19	0.09	0.13	0.53	0.54	<b>0.53</b>
Avg. Gain (abs.)	-	-	-	0.27 ↑	-	-	0.13 ↑	-	-	0.20 ↑

Table 12: Complete results of discriminative tasks under implicit memory settings across different models and prompt baselines.

- 1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381
3. **Fine-grained analysis.** RP-REASONER consistently outperforms Vanilla, Reminder, and CoT, with limited gains from the latter two. At a fine-grained level, RP-Reasoner markedly reduces error severity on strategy-level FB and RII, while paying only a small cost on UPB; at the response level it indirectly lowers LF and **VG** errors via better strategy control (Figures 19, 20).
- 1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389
4. **Multi-preference difficulty.** Multi-memory/multi-preference settings substantially raise task complexity: models struggle on *Macro-ACC*, and even with RP-Reasoner the global accuracy is only about  $\sim 20\%$ , though still clearly above other methods. This underscores that RPEval under multiple preferences is highly challenging. (Tables 14, Figure 20).

Table 13: Complete results of generative tasks under *single-preference* setting across different models and prompt baselines.

	ACC $\uparrow$				Overall Judge $\downarrow$			
	Ign.	Sup.	Dom.	ALL	Ign.	Sup.	Dom.	ALL
Qwen2.5-7B	0.46	0.94	0.60	0.67	1.94	1.08	1.96	1.66
+Reminder	0.38	0.96	0.66	0.67	2.20	1.04	1.72	1.65
+CoT	0.02	0.98	0.88	0.63	3.74	1.20	1.74	2.23
+RP-Reasoner	0.78	0.78	0.62	<b>0.73</b>	0.66	0.76	1.62	<b>1.01</b>
DeepSeek-V3	0.02	1.00	0.96	0.66	3.98	0.98	0.96	1.97
+Reminder	0.04	0.98	0.98	0.67	3.72	1.06	0.92	1.9
+CoT	0.00	1.00	1.00	0.67	3.84	1.1	1.64	2.19
+RP-Reasoner	0.72	0.94	0.98	<b>0.88</b>	1.06	0.32	0.28	<b>0.55</b>
GPT-4.1	0.06	1.00	1.00	0.69	3.28	0.90	0.98	1.72
+Reminder	0.02	1.00	1.00	0.67	3.24	1.00	0.94	1.73
+CoT	0.00	1.00	0.98	0.66	3.90	1.24	1.52	2.22
+RP-Reasoner	0.70	0.98	1.00	<b>0.89</b>	1.1	0.96	0.9	<b>0.99</b>
GPT-5	0.07	1.00	1.00	0.61	3.22	1.02	1.07	1.79
+Reminder	0.08	1.00	1.00	0.69	3.08	1.00	1.04	1.71
+CoT	0.00	1.00	1.00	0.67	3.88	1.12	1.50	2.16
+RP-Reasoner	0.72	1.00	1.00	<b>0.91</b>	1.52	0.98	1.06	<b>1.19</b>
Avg. Gain (abs.)	-	-	-	0.20 $\uparrow$	-	-	-	-0.85 $\downarrow$

Table 14: Complete results of generative tasks under *multi-preference* setting across different models and prompt baselines.

	MACRO-ACC $\uparrow$			MICRO-ACC $\uparrow$			Overall Judge $\downarrow$		
	IA	LKN	ALL	IA	LKN	ALL	IA	LKN	ALL
Qwen2.5-7B	0.18	0.01	0.08	0.22	0.57	0.45	3.05	2.69	2.83
+Reminder	0.08	0.02	0.05	0.10	0.52	0.38	3.52	2.72	3.04
+CoT	0.00	0.00	0.00	0.00	0.53	0.35	4.25	3.64	3.89
+RP-Reasoner	0.42	0.03	<b>0.19</b>	0.45	0.51	<b>0.49</b>	2.13	2.88	<b>2.58</b>
DeepSeek-V3	0.00	0.02	0.01	0.00	0.57	0.38	4.39	2.87	3.48
+Reminder	0.00	0.02	0.01	0.00	0.58	0.39	4.47	2.79	3.46
+CoT	0.00	0.03	0.02	0.00	0.57	0.38	4.12	3.09	3.52
+RP-Reasoner	0.42	0.04	<b>0.19</b>	0.42	0.58	<b>0.53</b>	2.23	2.67	<b>2.49</b>
GPT-4.1	0.00	0.01	0.01	0.02	0.59	0.40	4.13	2.60	3.21
+Reminder	0.00	0.01	0.01	0.02	0.56	0.38	4.05	2.84	3.33
+CoT	0.00	0.01	0.01	0.00	0.62	0.42	4.23	3.19	3.61
+RP-Reasoner	0.63	0.01	<b>0.26</b>	0.59	0.59	<b>0.59</b>	1.43	2.63	<b>2.15</b>
GPT-5	0.00	0.01	0.01	0.03	0.59	0.40	4.10	2.76	3.30
+Reminder	0.00	0.02	0.01	0.02	0.60	0.39	4.03	2.63	3.23
+CoT	0.00	0.01	0.01	0.00	0.60	0.37	4.26	3.31	3.73
+RP-Reasoner	0.55	0.04	<b>0.25</b>	0.59	0.60	<b>0.60</b>	1.88	2.27	<b>2.11</b>
Avg. Gain (abs.)	-	-	0.20 $\uparrow$	-	-	0.15 $\uparrow$	-	-	-0.87 $\downarrow$

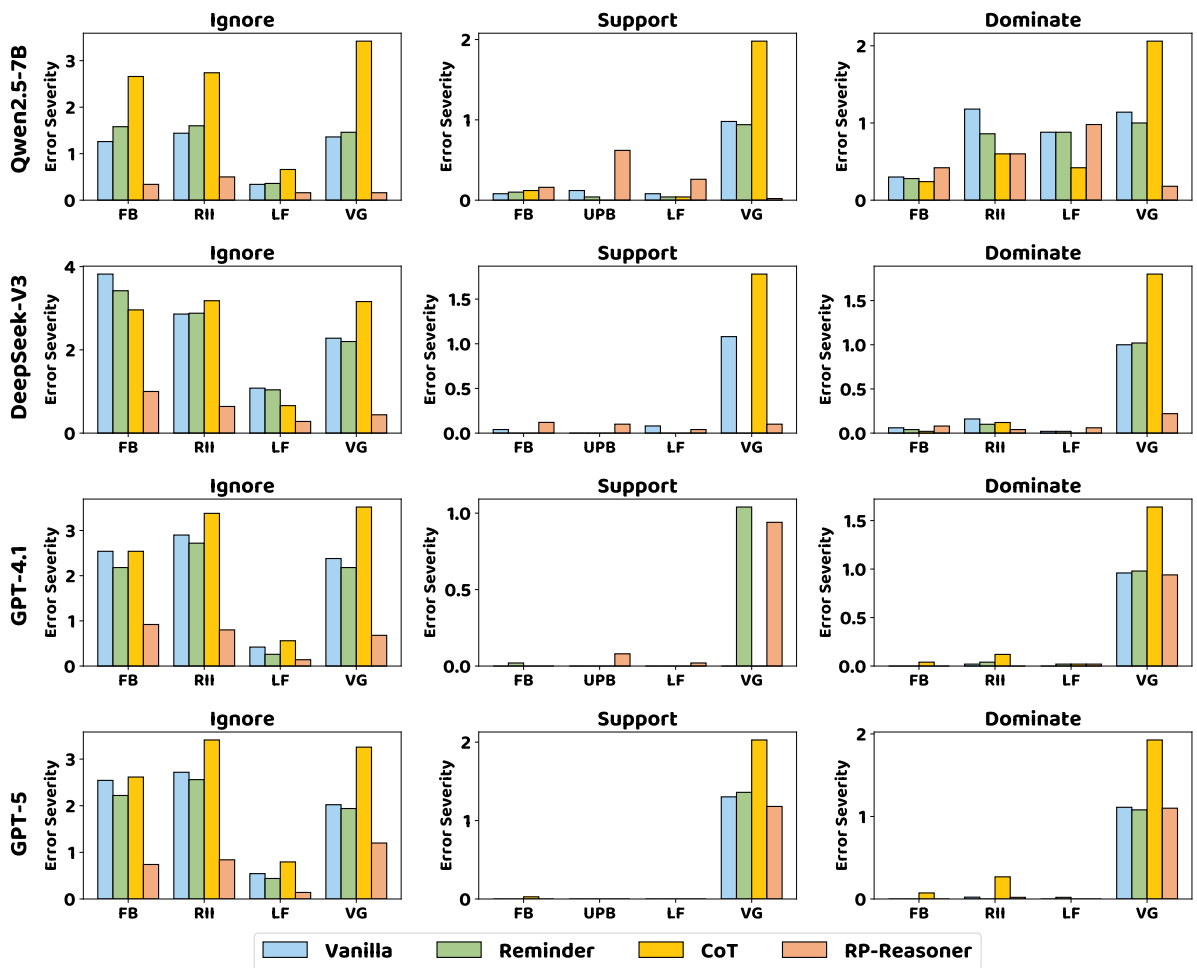


Figure 19: Complete results under the *single-preference* generative setting.

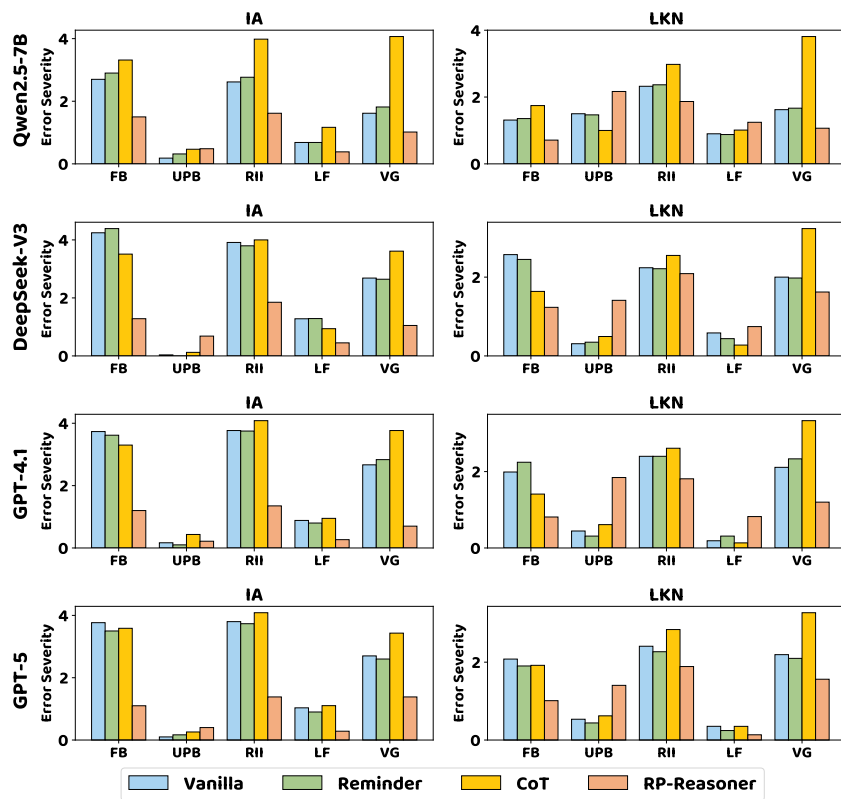


Figure 20: Complete results under the *multi-preference* generative setting.

## Prompts of Baselines: MCQ

### Vanilla

You are a personalized assistant, and you need to appropriately reference the user's persona to determine the most suitable answer strategy from the following three options:

- (A) **Ignore**: Act solely based on the current task objective, disregarding past preferences .
  - (B) **Support**: Attempt to fulfill the current task while integrating or partially retaining past preferences.
  - (C) **Dominate**: The current behavior is strongly driven by preferences, with the task focused around those preferences.
- 

### Reminder

You are a personalized assistant, and you need to decide whether and how to use preferences in this scenario. Choose the most appropriate answer strategy from the following three:

- (A) **Ignore**: Act solely based on the current task objective, disregarding past preferences .
  - (B) **Support**: Attempt to fulfill the current task while integrating or partially retaining past preferences.
  - (C) **Dominate**: The current behavior is strongly driven by preferences, with the task focused around those preferences.
- 

### CoT

You are a personalized assistant and need to refer to the user profile appropriately to determine the most appropriate answer strategy from the following three options:

- (A) **Ignore**: Act solely based on the current task objective, disregarding past preferences .
- (B) **Support**: Attempt to fulfill the current task while integrating or partially retaining past preferences.
- (C) **Dominate**: The current behavior is strongly driven by preferences, with the task focused around those preferences.

<Example>  
<Persona>  
<Question>  
<Chain of Thought>  
<Example End>

---

Output format:

```
{
  "persona": "<persona>",
  "question": "<question>",
  "reason": "<Your reasoning process>",
  "policy": "(A/B/C)"
}
```

Figure 21: Single-Preference Baseline Implementation in Discriminative Tasks.

### Vanilla

You are a personalized assistant, and your task is to appropriately reference the user's persona to determine the answer strategy. You need to assess the role of each preference in the current task and decide the corresponding strategy for each preference. Each preference should be categorized into one of the following three options:

- (A) **Ignore:** Act solely based on the current task objective, disregarding past preferences .
  - (B) **Support:** Attempt to fulfill the current task while integrating or partially retaining past preferences.
  - (C) **Dominate:** The current behavior is strongly driven by preferences, with the task focused around those preferences.
- 

### Reminder

You are a personalized assistant. You need to appropriately reference the user's persona to determine the response strategy, and decide whether and how preferences should be applied in this scenario. Each preference must be classified into one of the following three categories:

- (A) **Ignore:** Act solely based on the current task objective, disregarding past preferences .
  - (B) **Support:** Attempt to fulfill the current task while integrating or partially retaining past preferences.
  - (C) **Dominate:** The current behavior is strongly driven by preferences, with the task focused around those preferences.
- 

### CoT

You are a personalized assistant, and your task is to appropriately reference the user's persona to determine the answer strategy. You need to assess the role of each preference in the current task and decide the corresponding strategy for each preference. Each preference should be categorized into one of the following three options:

- (A) **Ignore:** Act solely based on the current task objective, disregarding past preferences .
- (B) **Support:** Attempt to fulfill the current task while integrating or partially retaining past preferences.
- (C) **Dominate:** The current behavior is strongly driven by preferences, with the task focused around those preferences.

```
<Example>
<Persona 0>... <Persona N>
<Question>
<Chain of Thought 0>... <Chain of Thought N>
<Example End>
```

---

Please analyze the role of each preference in the current task and output a string corresponding to the strategy for each preference. The length of the string should exactly match the number of preferences, and each character in the string must be A, B, or C.

Output format:

```
{
  "personas": "<personas>",
  "question": "<question>",
  "reason": "<Your reasoning process>",
  "policy": "(such as AABCC) "
}
```

### Prompts of Baselines in Generative Tasks

#### **Vanilla**

You are a personalized assistant, you need to refer to the user's persona to answer questions.

---

#### **Reminder**

You are a personalized assistant, and you should appropriately refer to the user's persona when answering questions. Note that some preferences may be inappropriate in this context and can be ignored.

---

#### **CoT**

You are a thoughtful and professional assistant who understands the user's habits and preferences. Please combine the user's preferences with the situational context to reason step by step (chain-of-thought), explain your recommendation logic or decision rationale, and provide a clear and tailored final answer.

Please structure your response as follows: 1. Clarify the user's intent 2. Analyze the user's personalized preferences 3. Reason and filter based on the scenario 4. Provide a personalized suggestion or answer.

Figure 23: Baselines Implementation in Generative Tasks.

## **E Future Works**

(1) Integrating with existing work on long-context modeling to investigate how retrieval- and generation-side filtering of irrelevant memories can be coordinated to achieve more rational personalization; (2) While RP-REASONER substantially improves rational memory utilization, it still falls short of human-level performance. Exploring how to train models to appropriately disregard preferences when necessary may represent a promising future direction.