

# CoReflect: Conversational Evaluation via Co-Evolutionary Simulation and Reflective Rubric Refinement

Anonymous ACL submission

## Abstract

Evaluating conversational systems in multi-turn settings remains a fundamental challenge. Conventional pipelines typically rely on manually defined rubrics and fixed conversational context—a static approach that limits coverage and fails to capture the diverse, emergent behaviors of dialogue models. To address this, we introduce CoReflect (Conversational Evaluation via Co-Evolutionary Simulation and Reflective Rubric Refinement), which unifies dialogue simulation and evaluation into an adaptive, iterative process. CoReflect employs a conversation planner that generates structured templates to guide a user simulator through diverse, goal-directed dialogues. Subsequently, a reflective analyzer processes these dialogues to identify systematic behavioral patterns and automatically refine the evaluation rubrics. Crucially, the insights from the conversation analysis are fed back into the planner to update conversation templates for subsequent iterations. This co-evolution loop ensures that the complexity of test cases and the diagnostic precision of rubrics improve in tandem. By minimizing human intervention, CoReflect provides a scalable and self-refining methodology that allows evaluation protocols to adapt alongside the rapidly advancing capabilities of dialogue models.

## 1 Introduction

Large language models (LLMs) increasingly power multi-turn applications such as personalized tutors and customer service agents (Brown et al., 2020; OpenAI, 2023). In these dynamic settings, users demand more than linguistic fluency; they expect coherent, adaptive, and contextually grounded dialogues. Evaluating such conversational capabilities, however, remains a challenging open problem.

To date, the most reliable approach to conversational evaluation relies on open-ended human assessment against predefined rubrics, as human raters can dynamically probe nuanced behaviors

and value alignment over extended interactions. However, this method is costly, time-consuming, and prone to subjective bias. While automated alternatives, such as “LLM-as-a-judge” frameworks, offer efficiency and have shown promising alignment with human judgments (Liu et al., 2023a; Zheng et al., 2023; Zhao et al., 2025), they often rely on static, pre-designed conversational contexts that fail to capture authentic model behavior. Furthermore, LLM-based simulators, while capable of producing fluent and natural conversational interactions, often struggle to systematically elicit challenging or adversarial behaviors, and may exhibit shallow exploration or mode collapse in self-driven interactions (Park et al., 2023; Shanahan et al., 2023; Chiang et al., 2024).

More fundamentally, both human and automated evaluation approaches can be constrained by the evaluation architect’s subjective focus and limited foresight. Overly narrow rubrics may fail to capture emergent model capabilities or unforeseen behaviors that manifest during subsequent deployment, leading to costly redesign and re-execution of the evaluation process.

To address this limitation, we propose **CoReflect**, a co-evolutionary and reflective evaluation framework that moves beyond static benchmarks and fixed rubrics. By introducing an adaptive, iterative process, CoReflect ensures that both dialogue simulation and evaluation criteria evolve in tandem to capture model behaviors. Instead of relying on static conversational scripts, the framework utilizes a *conversation planner* as the central engine that orchestrates the entire interaction under specific user persona and scenario. Specifically, the planner generates structured, goal-directed templates that prescribe specific turn-level instructions for a user simulator on how to interact with test models.

Given these planner-generated interactions, an *LLM-as-a-judge evaluator* then assesses the resulting conversations using a structured three-step

084 protocol—analyzing turns, synthesizing observa- 132  
085 tions, and assigning ratings against the current 133  
086 rubrics. Following the evaluation results, a *reflec- 134*  
087 *tive analyzer* processes these conversation assess- 135  
088 ments to cluster behavioral patterns and systematic 136  
089 failures, extracting insights to refine the evaluation 137  
090 rubrics automatically. These insights are also fed 138  
091 back to the planner, which updates its conversa- 139  
092 tion templates to target identified weaknesses with 140  
093 greater precision in subsequent iterations. 141

094 Our main contributions are presented as follows: 142

- 095 • **Co-evolutionary conversational evaluation 143**  
096 **framework.** We propose CoReflect, a novel co- 144  
097 evolutionary evaluation framework in which con- 145  
098 versation templates and evaluation rubrics are 146  
099 jointly refined through an iterative feedback loop. 147  
100 This co-evolution mechanism is the core innova- 148  
101 tion of our work, allowing the user simulation 149  
102 and the diagnostic precision of the rating criteria 150  
103 improve in tandem. 151
- 104 • **Autonomous synthesis of evaluation rubrics.** 152  
105 CoReflect introduces a reflective analyzer as the 153  
106 primary mechanism for translating model behav- 154  
107 ior data into structured evaluation logic, effec- 155  
108 tively alleviating the human burden of manual 156  
109 design. By autonomously synthesizing patterns 157  
110 from simulated interactions, it mitigates subjec- 158  
111 tive human bias and enables the system to dy- 159  
112 namically expand its scope beyond static rubrics. 160  
113 This data-driven approach allows CoReflect to 161  
114 capture emergent model behaviors that might oth- 162  
115 erwise remain hidden, ensuring the evaluation 163  
116 framework evolves in lockstep with the model’s 164  
117 complexity. 165
- 118 • **Comprehensive empirical validation.** We con- 166  
119 duct extensive automatic and human-in-the-loop 167  
120 experiments across diverse personas, scenar- 168  
121 ios, and model families, demonstrating the abil- 169  
122 ity of CoReflect to support comprehensive and 170  
123 behavior-sensitive evaluation of state-of-the-art 171  
124 LLMs in multi-turn conversations. 172

## 125 2 Methodology 173

126 The evaluation process is initialized by a human 174  
127 architect who manually defines a set of coarse, high- 175  
128 level rubrics to serve as the foundational criteria 176  
129 for assessing the conversational quality of the test 177  
130 models. In each iteration, each model interacts with 178  
131 a user simulator across various persona–scenario 179

132 pairs, guided by a *conversation planner*. This plan- 133  
134 ner generates goal-directed, persona-grounded con- 135  
136 versation templates that instructs the user simulator 137  
138 on how to conduct each user turn, thereby precisely 139  
140 challenging the models against the given rubrics. 141

142 Following these interactions, an *evaluator* as- 143  
144 sesses the conversations based on the established 145  
146 rubrics, providing quantitative ratings and qualita- 147  
148 tive rationales. These results are then sampled and 149  
150 processed by a *reflective analyzer* to characterize 151  
152 model behaviors and limitations, offering insights 153  
154 that drive an iterative refinement loop; in this loop, 155  
156 the rubrics are tuned and the planner accordingly 157  
158 updates its templates to better challenge the models 159  
160 in subsequent rounds. Figure 1 illustrates this inter- 161  
162 play between the planner, evaluator, and reflective 163  
164 analyzer. This process proceeds through time steps 165  
166  $t = 1, 2, \dots, T$ , where  $t = 1$  represents the initial 167  
168 state, and terminates after  $T$  iterations. 169

170 We provide the specifics of the framework and 171  
172 the process as follows. 173

### 174 2.1 Personas, Scenarios, and Test Models 175

176 Each evaluation *instance* in this study comprises 177  
178 two components: (i) a user, who is manifested 179  
180 through a persona profile, (ii) a scenario, which 181  
182 defines a context, task, and objective of the conver- 183  
184 sation. Every test model is evaluated within these 184  
185 personas and scenarios pairings, with the goal of 185  
186 fulfilling the user’s intent throughout the conver- 186  
187 sation. Below, we elaborate the construction and 187  
188 the pairing process of personas and scenarios to 188  
189 generate the evaluation set. 189

190 **Personas.** To capture user behaviors in our eval- 191  
192 uation, we create structured personas that specify 192  
193 how users express themselves and how they expect 193  
194 the system to interact with them. Specifically, each 194  
195 persona comprises two parts: (1) **User traits:** the 195  
196 user’s attributes manifested in their communica- 196  
197 tion, such as tone, verbosity, and conversational 197  
198 nuances. We denote user traits by  $x$ . (2) **Response 198**  
199 **preferences:** the desired qualities of the model’s re- 199  
200 sponses, including reasoning depth, level of detail, 200  
201 and formality, etc. We denote user preferences by 201  
202  $\theta$ . This representation enables controlled variation 202  
203 across communication styles while maintaining re- 203  
204 alism and interpretability. Complete schema defini- 204  
205 tions and examples are provided in Appendix A. 205

206 **Scenarios.** Scenarios are created independently 206  
207 from personas to ensure broad coverage of con- 207  
208 versational goals. Each scenario specifies a con- 208  
209 209

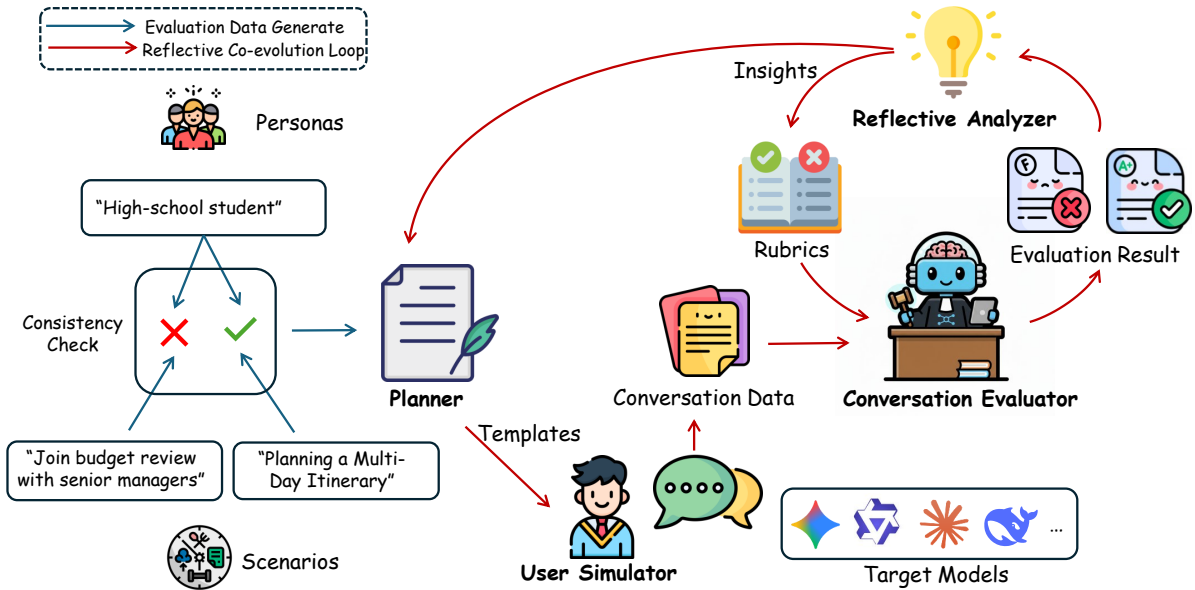


Figure 1: Overview of the CoReflect framework. Initially, evaluation instances are generated from persona-scenario pairs and validated through a consistency check to ensure contextual coherence. Then the core co-evolution loop proceeds through two parts: (i) Conversation simulation, where a planner generates structured conversation templates to guide a user simulator through goal-directed dialogues across diverse persona-scenario pairs; and (ii) Reflective rubric refinement, where an LLM-as-a-judge evaluates simulated interactions and a reflective analyzer extracts insights from clustered behavioral patterns to refine evaluation rubrics. The insights are fed back to the planner to update conversation templates for subsequent iterations.

text, primary task, objective of the conversation, and an estimate of number of turns to complete the task. We categorize scenarios into four intent types—*instructional*, *informational*, *operational*, and *interactive*—as detailed in Table 1. Collectively, these categories encompass factual, task-oriented, and social interactions. We denote a scenario by  $s$ .

**Pairing process and consistency check.** Each persona is paired with every scenario. However, not every persona–scenario combination yields a contextually-coherent setup (e.g., a young student tutoring a senior engineer). To ensure plausibility, each persona–scenario pair undergoes a consistency check conducted by a model-based verifier to assess the contextual coherence in each pair. Only those pairs validated as plausible are retained for simulation and evaluation. Formally, let  $\mathcal{D}$  denote this synthetic dataset of validated persona–scenario pairs. We represent the  $i$ -th pair as a tuple  $(x_i, \theta_i, s_i)$ . More details on consistency check can be found in Appendix C.

**Test models.** Each test model is evaluated against the whole set  $\mathcal{D}$ . The collection of test models is denoted by  $\mathcal{M} \triangleq \{m_j\}_{j=1}^M$ , where each

$m_j$  represents a specific conversational model and  $M = |\mathcal{M}|$  is the total number of the models.

## 2.2 Planner and Conversation Simulation

Let  $K$  denote the number of rubrics considered in the evaluation. We define the set of rubrics at the start of iteration  $t$  by  $\mathcal{R}^{(t)} \triangleq \{R_k^{(t)}\}_{k=1}^K$ . Whereas the initial set  $\mathcal{R}^{(1)}$  is manually drafted, these rubrics are refined in subsequent iterations.

To improve template quality, incorporates a set of insights from the preceding iteration, denoted by  $\mathcal{I}^{(t-1)}$ . We initialize  $\mathcal{I}^{(0)} = \emptyset$ , and will explain how insights are derived in Section 2.4.

At iteration  $t$  with insights  $\mathcal{I}^{(t-1)}$ , for each validated instance  $(x_i, \theta_i, s_i) \in \mathcal{D}$ , the planner determines the number of turns  $N_i^{(t)}$  required to complete the task and generates a structured conversation template  $\mathcal{T}_i^{(t)}$  that encodes specific instructions for every user turn. These turn-level instructions guide the dialogue, allowing the conversation to unfold naturally while systematically probing specific model capabilities (e.g., reasoning, clarification). Appendix D describes the template schema.

Once the templates for all persona–scenario pairs in  $\mathcal{D}$  are created, each model  $m_j \in \mathcal{M}$  conducts a

Category	User Intent	Assistant Role	Typical Outputs	Example
Instructional	Acquire knowledge or develop skills	Tutor / Instructor	Explanations, practice questions, scaffolds	“Teach me how neural networks work.”
Informational	Retrieve facts or clarify references	Reference / Expert	Factual answers, summaries, citations	“What is the GDP of Brazil in 2023?”
Operational	Complete a concrete task or create content	Productivity Tool	Code, documents, plans, translations	“Write a weekly report summary for my team.”
Interactive	Engage socially, emotionally, or creatively	Companion / Partner	Stories, dialogues, empathetic replies	“I feel nervous about my exam—can we talk about it?”

Table 1: Four scenario categories in the evaluation dataset  $\mathcal{D}$ , defined by specific user intents, assigned assistant roles, characteristic outputs, and representative examples.

dialogue with the user simulator for every pair. The simulator is provided with the instance  $(x_i, \theta_i, s_i)$  and the corresponding template  $\mathcal{A}_i^{(t)}$ , which serve as contextual grounding and turn-level instructions. In contrast, model  $m_j$  has only user traits  $x_i$  as input information. The simulator and model  $m_j$  then have a dialogue, based on their respective information sets and the evolving conversation history. This asymmetric information is designed to simulate real-world chatbots leveraging historical personal data (i.e.,  $x_i$ ) to personalize interactions. We denote the resulting conversation generated by model  $m_j$  for the  $i$ -th instance at iteration  $t$  as  $C_{i,j}^{(t)}$ .

### 2.3 Conversation Evaluator

An LLM-as-a-judge evaluates each generated conversation  $C_{i,j}^{(t)}$  by model  $m_j$  against the corresponding instance  $(x_i, \theta_i, s_i)$  and the rubric set  $\mathcal{R}^{(t)}$ .

To ensure evaluative reliability and interpretability, the judge generates a rating and a corresponding rationale for each rubric, following a structured three-step reasoning protocol. First, it performs a turn-level analysis to generate concise observations for every model response grounded in the immediate context. Next, it aggregates these turn-level observations to identify conversation-wide strengths and weaknesses. Finally, the evaluator determines a final rating, assigning a numerical score to each rubric alongside a rationale that references the synthesized evidence.

Let  $r_{i,j,k}^{(t)}$  and  $\rho_{i,j,k}^{(t)}$  denote the rating and rationale for conversation  $C_{i,j}^{(t)}$  as evaluated under rubric  $R_k^{(t)} \in \mathcal{R}^{(t)}$ . The performance of model  $m_j$  for a specific rubric  $R_k^{(t)}$  is defined as the mean rating across the conversation set  $\mathcal{D}$ ; the model’s aggregate performance,  $\mu_j^{(t)}$ , is then calculated as the

average across all  $K$  rubrics:

$$\mu_{j,k}^{(t)} \triangleq \frac{1}{|\mathcal{D}|} \sum_{i=1}^{|\mathcal{D}|} r_{i,j,k}^{(t)}, \quad k = 1, \dots, K$$

$$\mu_j^{(t)} \triangleq \frac{1}{K} \sum_{k=1}^K \mu_{j,k}^{(t)}$$

To measure performance consistency, we first define the conversation-level rating  $\tilde{r}_{i,j}^{(t)}$  for each  $C_{i,j}$  as the average rating across all  $K$  rubrics, i.e.,  $\tilde{r}_{i,j}^{(t)} \triangleq \frac{1}{K} \sum_{k=1}^K r_{i,j,k}^{(t)}$ . We then compute the variance of these ratings across the conversations  $\mathcal{D}$ :

$$(\sigma_j^{(t)})^2 \triangleq \frac{1}{|\mathcal{D}| - 1} \sum_{i=1}^{|\mathcal{D}|} (\tilde{r}_{i,j}^{(t)} - \mu_j^{(t)})^2$$

This quantifies the stability of model  $m_j$  across diverse scenarios, where a lower variance indicates more predictable performance.

### 2.4 Reflective Analyzer for Rubric Refinement and Planner Feedback

A reflective analyzer synthesizes the collected ratings  $\{r_{i,j,k}^{(t)}\}$  and rationales  $\{\rho_{i,j,k}^{(t)}\}$  to extract insights  $\mathcal{I}^{(t)}$  regarding model behavior. These insights are then leveraged to sharpen the rubrics’ diagnostic precision and to refine the planner’s templates, enabling a more effective testing of model strengths and weakness in subsequent iterations. The reflective analyzer proceeds as follows:

- (1) **Instance sampling:** Based on the assigned ratings  $\{r_{i,j,k}^{(t)}\}$ , we partition the conversations into two subsets representing high- and low-rated tiers. To ensure a balanced and computationally efficient analysis, we draw an equal number of samples from each tier to construct an analysis pool. For each sampled conversation, we retain only the associated rationale  $\rho_{i,j,k}^{(t)}$  as the structured evidence for subsequent behavioral analysis.

(2) **Pattern discovery:** The sampled rationales are embedded and clustered into a set of *behavioral families*,  $\mathcal{F}^{(t)} \triangleq \{f_\ell^{(t)}\}_{\ell=1}^{L^{(t)}}$ , where  $L^{(t)}$  represents the number of distinct, recurring patterns identified at iteration  $t$ . Each family  $f_\ell$  characterizes a specific conversational pattern, including both undesirable behaviors (e.g., loss of task focus, contradiction to prior context, stylistic drift) and effective strategies (e.g., proactive clarification, structured task decomposition).

(3) **Insight synthesis and rubric refinement.** For every family  $f_\ell^{(t)}$ , the reflective analyzer produces an interpretable insight  $\iota_\ell^{(t)}$  that characterizes the underlying model behavior and specifies the criteria for reward or penalty. These individual insights are aggregated into a comprehensive set  $\mathcal{I}^{(t)} \triangleq \{\iota_\ell^{(t)}\}_{\ell=1}^{L^{(t)}}$ . These insights are utilized to revise the corresponding rubrics in  $\mathcal{R}^{(t)}$ , specifically by updating performance definitions, rating anchors, and evidence cues. This process enhances the framework’s diagnostic precision for the subsequent iteration, formalized as:

$$\mathcal{R}^{(t+1)} \leftarrow \text{UPDATE}(\mathcal{R}^{(t)}, \mathcal{I}^{(t)})$$

As previously noted, the derived insights  $\mathcal{I}^{(t)}$  are inputs for the conversation planner. This integration refines the conversation templates to target identified behaviors with better precision in the subsequent iteration.

## 2.5 Measures on Rubrics Refinement

To assess whether iterative rubric refinement improves evaluation quality, we track two complementary measures across iterations  $t$ : rubric *discriminability* and *stability*. Rubric discriminability, denoted by  $\Delta^{(t)}$ , captures the degree to which a rubric separates models with different capability levels, operationalized as the inter-model variability of mean ratings. Rubric stability is measured by  $\Gamma^{(t)}$ , which evaluates rating consistency through intra-model variance, together with rank consistency measured by Spearman’s  $\rho$ . Their formal definitions are provided in Appendix F.

## 3 Experiments

### 3.1 Experiment Setup

**Validated persona-scenario pairs  $\mathcal{D}$ .** We created 30 personas, each embodying distinct com-

munication styles and response preferences, along with 10 scenarios per category (cf. Table 1), yielding 953 validated pairs after automated consistency checks. Summary statistics are provided in Table 2.

Additionally, we validated the user simulator via a small human assessment detailed in Appendix H.

**Test models.** We evaluated a range of frontier LLMs from different providers: Gemini-2.5-pro, Gemini-2.5-flash, and Gemini-1.5-flash (Comanici et al., 2025); Claude Sonnet-4.5, Claude Haiku-4.5 and Claude Sonnet-4 (Anthropic, 2024b,a); Qwen3-Next (Yang et al., 2025), DeepSeek-R1-671B (Guo et al., 2025).

**Rubrics design.** We evaluated model performance along two high-level dimensions that pose persistent challenges in conversational systems: (i) *Task Completeness (TC)*, which captures whether a model correctly understands task intent and makes coherent progress toward fulfilling the scenario objective; and (ii) *User-Centric Personalization (UCP)*, which reflects the model’s ability to adapt its responses to a user’s identity, preferences, and interaction style over the course of a dialogue.

Each dimension is assessed using three complementary rubrics, yielding a total of  $K = 6$  rubrics:

- **Task Completeness:** Domain Conceptual Alignment (DCA), Output Delivery Integrity (ODI), and Functional Task Progression (FTP) evaluate whether responses are conceptually correct, actionable, and incrementally advance the task.
- **User-Centric Personalization:** Anticipatory Flow Management (AFM), Output Structure Fit (OSF), and Sustained Style Adherence (SSA) assess whether the model anticipates user needs, adapts response organization, and maintains consistent stylistic alignment.

Each rubric was rated on a 1–5 scale, where 5 indicated the ideal outcome. While rubric

Category	#persona	# validated pairs	Avg. required turns per template
Informational	30	215	6.60
Instructional	30	244	6.97
Interactive	30	273	7.23
Operational	30	221	6.48
<b>Total</b>	<b>30</b>	<b>953</b>	<b>6.85</b>

Table 2: Summary statistics of the evaluation dataset across four scenario categories.

Test models	Task Completeness				User-Centric Personalization				Model rating
	ODI	DCA	FTP	avg.	AFM	OSF	SSA	avg.	
Claude Sonnet 4	4.65	4.76	4.25	4.56	4.75	4.74	4.73	4.74	4.65
Claude Sonnet 4.5	4.37	4.69	3.53	4.20	4.13	4.68	4.57	4.46	4.33
Claude Haiku 4.5	4.02	4.64	3.16	3.94	4.47	4.57	4.61	4.55	4.25
Qwen3-Next	4.73	4.03	4.50	4.42	4.62	4.72	4.70	4.68	4.55
DeepSeek-R1	4.68	4.65	4.59	4.64	4.60	4.54	4.54	4.56	4.60
Gemini 2.5 Pro	<b>4.86</b>	<b>4.86</b>	<b>4.70</b>	<b>4.81</b>	<b>4.76</b>	<b>4.86</b>	<b>4.79</b>	<b>4.80</b>	<b>4.81</b>
Gemini 2.5 Flash	4.79	4.72	4.12	4.60	4.70	4.72	4.76	4.75	4.68

Table 3: Model performance across all rubrics at iteration  $t = 3$ . (ODI: Output Delivery Integrity, DCA: Domain Conceptual Alignment, FTP: Functional Task Progression; AFM: Anticipatory Flow Management, OSF: Output Structure Fit, SSA: Sustained Style Adherence). See Appendix J for results at  $t = 1, 2$ .

names remained fixed, their descriptions and rating guidelines were iteratively refined through the co-evolution loop. We conducted  $T = 3$  iterations. The initial rubric definitions and rating criteria are provided in Appendix E. More implementation details are presented in Appendix G.

## 3.2 Experiment Results and Analysis

### 3.2.1 Model Performance and Stratification

Table 3 details rating results across the two primary dimensions at the end of the third iteration: (i) a model’s ability to carry out a task over multiple turns, and (ii) its capacity to adapt responses to user-specific signals while ensuring structural and stylistic consistency. The results reveal a clear stratification of model capabilities.

*Gemini 2.5 Pro* leads the frontier with the highest model rating of 4.81, dominating all evaluation rubrics. This performance reflects that the model’s superior reasoning and planning capabilities translate into highly effective dialogue management over complex, multi-turn interactions.

By contrast, the other models show mixed performance, excelling in specific areas while lagging in others. This variance underscores that high-quality multi-turn interaction requires more than isolated strengths; it demands a balanced integration of understanding, task progression, and user adaptation.

**Task Completeness: understanding versus execution.** On task-oriented metrics, *Gemini 2.5 Pro* achieves the strongest overall Task Completeness performance, with an average score of 4.81. This corresponds to a relative improvement of approximately 3.7% over the best non-Gemini baseline (*DeepSeek-R1*, 4.64), and more than 22% compared to smaller-capacity variants such as *Claude*

*Haiku 4.5*. Crucially, this advantage is not confined to surface correctness: *Gemini 2.5 Pro* also leads in Functional Task Progression (FTP), indicating a consistent ability to decompose tasks and track intermediate states across long-horizon interactions.

By contrast, runner-up models exhibit clear trade-offs between conceptual understanding and execution. For instance, *the Claude Sonnet series* and *DeepSeek-R1* achieve strong Domain Conceptual Alignment (DCA), yet their FTP scores lag behind *Gemini 2.5 Pro* by 25% or more in relative terms (e.g., *Claude Sonnet 4.5*). Conversely, *Qwen3-Next* prioritizes execution, with an FTP score within 4% of *Gemini 2.5 Pro*, but at the cost of substantially weaker DCA (approximately 17% lower), suggesting brittle execution under conceptual ambiguity. This execution gap is further magnified in reduced-capacity variants: *Gemini 2.5 Flash* and *Claude Haiku 4.5* trail *Gemini 2.5 Pro* in FTP by roughly 12% and 33% respectively, highlighting the difficulty of sustaining long-horizon task progression without sufficient model capacity.

**User-Centric Personalization: style consistency versus adaptability.** *Gemini 2.5 Pro* also leads in User-Centric Personalization, achieving an average score of 4.80. Although the absolute margin over other frontier models is narrower in this domain, it still represents a consistent relative improvement of approximately 5% over the strongest alternatives. Its advantage is most pronounced in Anticipatory Flow Management (AFM), where it outperforms *Claude Sonnet 4.5* by over 15% and *Claude Haiku 4.5* by roughly 6%, indicating a superior ability to anticipate evolving user intent rather than merely reacting to explicit cues. Other frontier-scale models, including *Claude Sonnet 4*,

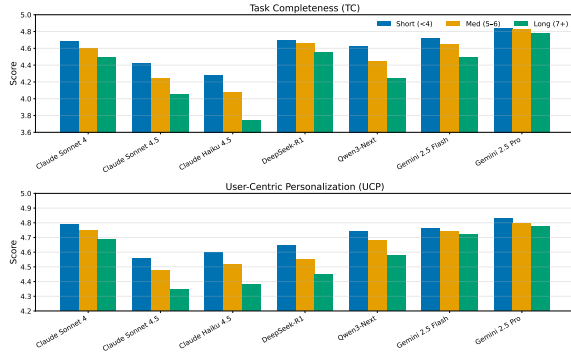


Figure 2: Model ratings across Task Completeness and User-Centric Personalization, stratified by conversation length.

*DeepSeek-R1*, and *Qwen3-Next*, cluster closely in overall personalization performance, with relative differences largely within 5%.

In contrast, *Claude Sonnet 4.5* demonstrates strong stylistic consistency, with OSF and SSA scores within 5% of the best-performing model, but its AFM score lags by more than 13%, revealing a clear trade-off between style adherence and adaptability. A similar pattern appears in *Claude Haiku 4.5*, where stylistic alignment remains competitive while anticipatory adaptation degrades more sharply. Taken together, these results suggest that while stylistic personalization is relatively robust across models, adaptive personalization—particularly the ability to steer conversations in response to latent or shifting user intent—exhibits a stronger dependence on model capacity.

### 3.2.2 Effect of Conversation Length

Figure 2 illustrates the impact of conversation length on model performance. Across all models, both Task Completeness and User-Centric Personalization scores decline as conversation length increases, indicating that extended interactions pose a systematic challenge for current models.

Task completeness degrades more significantly than personalization in long-horizon dialogues. While the models maintain stable personalization, their ability to track task states falters over time. This trend widens the performance gap between large-scale and lightweight models, which is indistinguishable in shorter interactions.

These results suggest conversation length as a critical factor for long-horizon reasoning and execution, highlighting the necessity of length-aware analysis in multi-turn conversational benchmarks

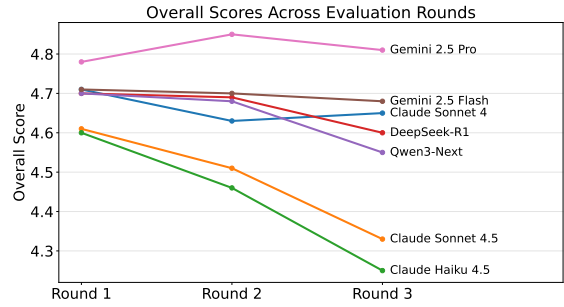


Figure 3: Model ratings across iterations. Iterative reflection improves model differentiation, yielding more distinct stratification and relative ordering.

Metric	$t = 1$	$t = 2$	$t = 3$
<i>Rubric Discriminability</i> ( $\Delta^{(t)}$ )			
Inter-model std. dev. ( $\sigma_{\text{inter}}$ )	0.062	0.128	<b>0.194</b>
<i>Rubric Stability</i> ( $\Gamma^{(t)}$ )			
Intra-model variance ( $\sigma_{\text{intra}}^2$ )	0.145	0.141	<b>0.138</b>
Rank consistency (Spearman $\rho$ )	0.68	0.79	<b>0.92</b>

Table 4: Evolution of rubric discriminability and stability across iterations. Iterative refinement increases inter-model separation while preserving or improving scoring stability, validating the quantitative feedback criteria used to accept rubric updates.

to accurately differentiate model capabilities.

### 3.2.3 Effect of Co-Evolution Iterations

Figure 3 illustrates the trajectory of model ratings across three iterations, showcasing a significant improvement in the framework’s discriminatory resolution. In the initial round, the ratings are too clustered to support reliable ordering. Through iterative reflective rubric refinement and co-evolutionary user simulation, the framework systematically uncovers subtle performance nuances that were previously obscured. Consequently, later iterations yield finer-grained distinctions, resulting in a clear and consistent stratification of model performance.

The observed visual trends are statistically corroborated by the metrics in Table 4, which show a monotonic increase in rubric discriminability ( $\Delta^{(t)}$ ), as measured by inter-model rating dispersion. The rubric stability ( $\Gamma^{(t)}$ ) is maintained or slightly enhanced, evidenced by stable intra-model variance and improved rank consistency. Collectively, the visual and statistical evidence confirms that the CoReflect framework successfully amplifies the signal of model differences without introducing stochastic noise, validating its iterative ap-

proach to multi-turn conversational evaluation.

## 4 Related Work

### 4.1 Structured User Simulation

Structured simulation provide a controllable framework for systematically probing model behavior, ranging from classical planning-based systems (Reiter and Dale, 2000) to modern neural plan-and-write pipelines (Yao et al., 2019; Xu et al., 2020). By applying these techniques to conversational evaluation, persona-driven and template-based simulations enable targeted stress-testing of a model’s personalization, consistency, and task execution.

### 4.2 LLM-as-a-Judge and Scalable Evaluation

To scale evaluation beyond human annotation, LLM-as-a-judge has emerged as a widely adopted paradigm. G-Eval (Liu et al., 2023a) demonstrates that LLM-based evaluators can achieve strong correlation with human judgments, while MT-Bench (Zheng et al., 2023) and CHATBOT ARENA (Chiang et al., 2024) enable large-scale multi-turn comparison through pairwise or rubric-based scoring. Complementary benchmarks such as LongBench (Bai et al., 2024), L-Eval (An et al., 2024), and RULER (Hsieh et al., 2024) focus on long-context understanding and sustained reasoning. Despite their effectiveness, most LLM-based evaluation frameworks adopt static rubrics or prompts, making them less sensitive to newly emerging failure modes in personalized and long-horizon interactions.

### 4.3 Self-Refining Evaluation

Motivated by the limitations of static LLM-as-a-judge frameworks, recent work has begun to explore *self-refining* evaluation paradigms that improve evaluators through reflection and iterative feedback loops. A representative line of research studies reflective evaluation, where LLM-based judges generate structured rationales or critiques before assigning scores, helping surface latent failure modes beyond scalar judgments (Liu et al., 2023a; Zheng et al., 2023; Shinn et al., 2023). Building on such signals, other work introduces iterative refinement loops that update evaluation prompts, rubrics, or protocols based on low-performing or ambiguous cases. In conversational settings, LLM-based user simulators enable interactive evaluation beyond fixed contexts: Wang et al. (2023) employ a ChatGPT-based simulator for

multi-turn conversational recommendation evaluation, while Liu et al. (2023b) show that using multiple simulators improves robustness over relying on a single one. Relatedly, Hu et al. (2022) use LLM-simulated user feedback as an annotation-free signal to iteratively optimize dialogue responses. Together, these studies highlight a shift toward closed-loop evaluation frameworks that continuously refine evaluation criteria as model behaviors evolve.

### 4.4 Evaluation of Personalized Dialogue

Personalization in conversational AI has traditionally centered on persona grounding, preference modeling, and adaptive generation. Foundational benchmarks like PERSONA-CHAT (Zhang et al., 2018), BLENDED SKILL TALK (Smith et al., 2020), and EMPATHETIC DIALOGUES (Rashkin et al., 2019) established controlled environments for assessing consistency, multi-skill agility, and emotional grounding. Recent efforts, such as DialogBench (Ou et al., 2024) and PersonaConvBench (?), have extended these evaluations to multi-turn interactions. However, while these resources effectively identify personalization failures, they rely on static scripts and evaluation criteria that struggle to adapt as conversational behaviors and user expectations evolve.

## 5 Conclusion

We introduced CoReflect, an autonomous framework unifying persona-driven simulation with iterative rubric refinement. By dynamically evolving evaluation criteria based on observed behaviors, CoReflect captured nuances that static benchmarks tend to miss. Experiments confirmed this self-refining approach significantly improved discriminatory resolution, resulting in clear, high-fidelity model stratification.

CoReflect provides a scalable solution for complex, hard-to-define domains such as personalization by automatically updating obsolete criteria and mitigating the subjective bias inherent in human-designed evaluation architectures. Future work will aim to strengthen CoReflect by grounding simulation in real user interactions and improving evaluation robustness through richer, bias-aware judging and rubric design.

## 606 Limitations

607 While CoReflect establishes an adaptive and systematic framework for conversational evaluation,  
608 it has two primary limitations. First, the framework relies on simulated user interactions driven by  
609 predefined, synthetic personas and scenarios. Although this design ensures control and scalability,  
610 these simulations may not fully capture the spontaneity, emotional nuance, and long-term dynamics  
611 of authentic human interaction. Second, CoReflect employs an LLM-as-a-judge for both rating and  
612 rubrics refinement. As a result, the evaluation process risks inheriting the intrinsic biases or blind  
613 spots of the underlying judge model, which could skew the trajectory of rubric evolution. Critically,  
614 while this study demonstrates the framework’s internal consistency, we have not yet conducted a human  
615 assessment to externally validate the generated rubrics or the judge’s scoring alignment. Future  
616 work must address this by comparing CoReflect’s automated outputs with human expert annotations  
617 and by employing diverse models as judges to investigate potential self-preference biases.

## 629 Ethical Considerations

630 To validate the realism of our user simulator, we conducted a study with three volunteer annotators  
631 under an ethically approved protocol. Participants were fully informed of the study’s purpose, and the  
632 evaluation was performed on randomly sampled simulated interactions that contain no sensitive information  
633 and no real-world human data. Human ratings were collected using a transparent 5-point Likert scale  
634 to ensure objective assessment, and all data were used strictly in accordance with their intended  
635 purposes. As a general evaluation framework, CoReflect does not explicitly control or filter for  
636 potentially offensive or unsafe content generated by the underlying base LLMs. Our objective  
637 is to evaluate conversational behaviors as exhibited by the models themselves, rather than to impose  
638 additional safety constraints that could confound behavioral analysis. We view content safety and  
639 moderation as complementary concerns that can be incorporated through external filters or safety-  
640 oriented rubrics in future extensions of the framework.

## References

- Zhen An, Yichong Liu, Xingyu Wang, Enze Xie, Xiaosong Wang, Shijian Ren, and Hongsheng Li. 2024. L-eval: Instituting long-context nlu and generation evaluation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL)*. 653-658
- Anthropic. 2024a. *The claude 3 model family: Haiku, sonnet, & opus*. Technical report. 659-660
- Anthropic. 2024b. Claude: An ai system by anthropic. <https://www.anthropic.com/claude>. 661-662
- Xiaozhi Bai, Ning Ding, Yulin Chen, Zhengxiao Li, Zhiyuan Liu, and Maosong Sun. 2024. Longbench: A bilingual, multitask benchmark for long context understanding. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL)*. 663-668
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, and 1 others. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. 669-674
- Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios N Angelopoulos, Tianle Li, Daniel Li, ..., and Ion Stoica. 2024. Chatbot arena: An open platform for evaluating llms by human preference. In *Proceedings of the 41st International Conference on Machine Learning (ICML)*. 675-680
- Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and 1 others. 2025. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*. 681-687
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shitong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*. 688-693
- Cheng-Ping Hsieh, Shih-Yang Sun, Samuel Krizan, Suyog Acharya, Debajyoti Rekes, Fei Jia, Yiming Zhang, and Boris Ginsburg. 2024. Ruler: What’s the real context size of your long-context language models? In *Conference on Language Modeling (COLM)*. ArXiv:2404.06654. 694-699
- Zhiyuan Hu, Yue Feng, Anh Tuan Luu, Bryan Hooi, and Aldo Lipani. 2022. *Unlocking the potential of user feedback: Leveraging large language models as user simulators to enhance dialogue systems*. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 7199–7211. Association for Computational Linguistics. 700-706

707	Jiaxian Liu, Yichong Liu, Qi Lei, Junteng Ma, and Zhilin Yang. 2023a. G-eval: Nlg evaluation using gpt-4 with better human alignment. In <i>Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> .	763
708		764
709		765
710		766
711		
712	Yajiao Liu, Xin Jiang, Yichun Yin, Yasheng Wang, Fei Mi, Qun Liu, Xiang Wan, and Benyou Wang. 2023b. <a href="#">One cannot stand for everyone! leveraging multiple user simulators to train task-oriented dialogue systems</a> . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1–21, Toronto, Canada. Association for Computational Linguistics.	767
713		768
714		769
715		770
716		771
717		772
718		773
719		
720	OpenAI. 2023. Gpt-4 technical report. <i>arXiv preprint arXiv:2303.08774</i> .	
721		
722	Jiahui Ou, Jun Lu, Chang Liu, Yichong Tang, Fangwei Zhang, Dylan Zhang, and Kun Gai. 2024. Dialog-bench: Evaluating llms as human-like dialogue systems. In <i>Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL)</i> .	
723		
724		
725		
726		
727		
728		
729	Joon Sung Park, Joseph O’Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In <i>Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology</i> .	
730		
731		
732		
733		
734		
735	Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. 2019. Towards empathetic open-domain conversation models: A new benchmark and dataset. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (ACL)</i> , pages 5370–5381.	
736		
737		
738		
739		
740		
741	Ehud Reiter and Robert Dale. 2000. <i>Building Natural Language Generation Systems</i> . Cambridge University Press, Cambridge, UK.	
742		
743		
744	Murray Shanahan, Kyle McDonell, and Laria Reynolds. 2023. Role play with large language models. <i>Nature</i> , 623:493–498.	
745		
746		
747	Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik R. Narasimhan, and Shunyu Yao. 2023. <a href="#">Reflexion: Language agents with verbal reinforcement learning</a> . In <i>Proceedings of the 37th Conference on Neural Information Processing Systems (NeurIPS 2023)</i> , pages 1–15. Curran Associates, Inc.	
748		
749		
750		
751		
752		
753		
754	Eric Michael Smith, Mary Williamson, Kurt Shuster, Jason Weston, and Y-Lan Boureau. 2020. Can you put it all together? evaluating conversational agents’ ability to blend skills. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL)</i> , pages 2021–2030.	
755		
756		
757		
758		
759		
760	Xiaolei Wang, Xinyu Tang, Xin Zhao, Jingyuan Wang, and Ji-Rong Wen. 2023. <a href="#">Rethinking the evaluation for conversational recommendation in the era of large language models</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 10052–10065, Singapore. Association for Computational Linguistics.	763
761		764
762		765
		766
	Jingjing Xu, Dongyang Ju, Xuancheng Li, Xipeng Qiu, and Xuanjing Huang. 2020. Megatron-cntrl: Controllable story generation with external knowledge using large-scale language models. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 2831–2845. Association for Computational Linguistics.	767
		768
		769
		770
		771
		772
		773
	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. <i>arXiv preprint arXiv:2505.09388</i> .	774
		775
		776
		777
		778
	Lili Yao, Nanyun Peng, Ralph Weischedel, Kevin Knight, Dongyan Zhao, and Rui Yan. 2019. Plan-and-write: Towards better automatic storytelling. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 33, pages 7378–7385.	779
		780
		781
		782
		783
	Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. Personalizing dialogue agents: I have a dog, do you have pets too? In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (ACL)</i> , pages 2204–2213.	784
		785
		786
		787
		788
		789
	Shangqing Zhao, Ming Hong, Yutao Liu, Devamanyu Hazarika, and Kevin Lin. 2025. Prefeval: Do llms recognize your preferences? evaluating personalized preference following in llms. In <i>International Conference on Learning Representations (ICLR)</i> .	790
		791
		792
		793
		794
	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Shiyi Zhuang, Zi Lin Wu, Yong Zhuang, ..., and Ion Stoica. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. <i>arXiv preprint arXiv:2306.05685</i> .	795
		796
		797
		798

799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839

## Appendix

### A User Persona Definition

To simulate diverse user profiles in personalized multi-turn conversations, we define a structured persona schema capturing both the user’s expressive behavior and their expectations for AI responses. Each persona is composed of the following elements:

**ID and Role** A unique identifier and a brief description of the user’s background or profession (e.g., University Student, Senior Developer). This provides contextual grounding for the user’s domain knowledge and communication style.

**Language** The primary language used by the persona, enabling multilingual dialogue simulation and evaluation.

**Personality** This component models how the user naturally communicates. It includes:

- **Traits:** Core personality characteristics that influence conversational behavior (e.g., Impatient, Inquisitive, Skeptical).
- **Tone:** The persona’s typical emotional tone in conversation (e.g., Formal, Sarcastic, Friendly).
- **Verbosity:** Typical message length and level of detail (e.g., Terse, Verbose).
- **Quirks:** Specific linguistic habits that define the user’s style (e.g., Skips pleasantries, Uses all lowercase, Frequently references personal projects).

**Preferred Response Style** This defines how the user expects the AI assistant to respond and serves as the personalization target. It includes preferences for:

- Tone (e.g., Empathetic, Direct),
- Verbosity (e.g., Concise, Detailed),
- Reasoning Depth (e.g., Brief rationale, Step-by-step explanation),
- Engagement (e.g., Neutral, Asks follow-ups),
- Clarity (e.g., Simple and clear, Technically precise).

### B Scenario Category Definition

To comprehensively evaluate the capabilities of large language models in multi-turn dialogue, we define four broad categories of interaction scenarios and a structured schema for their definition.

#### Scenario Categories

- **Instructional Dialogue:** Encompasses interactions where the user’s objective is to acquire new knowledge or develop specific skills. The model serves in an instructional capacity, providing pedagogical guidance and step-by-step reasoning.
- **Informational Dialogue:** Characterized by the user’s intent to obtain accurate, context-relevant information. The model functions as a factual resource, delivering concise and verifiable responses.
- **Operational Dialogue:** Includes interactions where the user seeks assistance in completing a concrete task or producing a functional output. The model is employed as a productivity facilitator to generate content or perform structured actions.
- **Interactive Dialogue:** Captures scenarios where the user’s primary aim is social, emotional, or imaginative engagement. The model adopts a human-like persona, prioritizing coherence, tone adaptation, and creative collaboration.

#### Scenario Definition Schema

- **ID and Title:** A unique scenario identifier and a concise title summarizing the context (e.g., Planning a Weekly Study Schedule, Summarizing Technical Emails).
- **Category:** A high-level classification of the task domain, selected from predefined types such as Work & Productivity, Learning & Study, Creative & Brainstorming, Personal & Planning, and Information & Factual.
- **Situation Description:** A neutral, third-person description of the background or external context.
- **Core Task:** A focused objective to be accomplished (e.g., Organize a calendar with overlapping deadlines).

840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884

885	• <b>Expected Interaction:</b>		929
886	– Turn Complexity: Short (<4 turns),	• <b>Generation Instructions:</b>	930
887	Medium (5–6 turns), or Long (7+ turns).	– turnType (Early Turn, Intermediate	931
888	– Dialogue Flow Type: Goal-Oriented,	Turn, Challenging Turn),	932
889	Open-Ended, Step-by-Step Guidance, or	– turnIntent (logic-level description of	933
890	Collaborative Generation.	user’s aim),	934
891	• <b>Success Criteria:</b> A neutral description of	– instructionForEval (rationale for	935
892	what constitutes a successful outcome (e.g.,	why this turn tests AI behavior).	
893	a clear and logically ordered prioritization of	• <b>Conversation Structure:</b>	936
894	incoming requests).	– Early Turn: Persona initiates task using	937
895		typical style.	938
	<b>C Automated Consistency Check</b>	– Normal Turns: Task progression through	939
896	A critical validation step is performed before any	elaboration or clarification.	940
897	persona is paired with a scenario. To prevent the		
898	creation of illogical or nonsensical conversations,	<b>E Metrics Design</b>	941
899	we employ an automated consistency check using	<b>Meta-Metrics</b>	942
900	a large language model (LLM) configured as a	• <b>User-Centric Personalization:</b> Effectiveness of	943
901	“judge.”	adapting to the user’s persona, dialogue history,	944
902		and explicit/implicit feedback.	945
903	<b>Protocol Steps</b>	• <b>Task &amp; Completeness:</b> Degree to which the	946
904	1. <b>Framing the Prompt:</b> The function sends a	model fulfills the task and adheres to all stated	947
905	targeted prompt to an LLM, instructing it to	requirements/constraints.	948
906	act as a “logical reasoning expert” for the specific		
907	task of evaluating a pairing’s plausibility.	<b>Personalization Sub-metrics (Refined via Self-</b>	949
908	2. <b>Providing Context:</b> The prompt is populated	<b>Reflection)</b>	950
909	with the role from the persona definition and	• <b>Sustained Style Adherence:</b> Consistency in	951
910	the core_task and title from the scenario definition.	maintaining adapted style and honoring negative	952
911	3. <b>Binary Decision:</b> The LLM is strictly constrained	constraints across turns.	953
912	to answer with only a single word, “Yes” or “No.”	• <b>Output Structure Fit:</b> Alignment of formatting/	954
913		structure with user preferences or requested	955
914	4. <b>Filtering:</b> If the LLM returns “Yes,” the pair	schema.	956
915	is approved. If “No,” the combination is logged	• <b>Anticipatory Flow Management:</b> Proactive	957
916	and discarded.	pacing and next-step guidance that matches user	958
917		needs and context.	959
918	<b>D Conversation Simulation and Template</b>	<b>Task Execution Sub-metrics (Refined via Self-</b>	960
919	To simulate realistic multi-turn conversations between	<b>Reflection)</b>	961
920	a user and an AI assistant, we define a structured	• <b>Domain Conceptual Alignment:</b> Correct identification	962
921	generation protocol. Each conversation is generated	and sustained adherence to the task’s	963
922	from a persona-scenario pair and is designed to assess	conceptual domain (e.g., fictional vs. real, technical	964
923	both personalization alignment and task adherence.	vs. general).	965
924		• <b>Output Delivery Integrity:</b> Reliable, complete	966
925	<b>Protocol Components</b>	delivery of required artifacts without truncation/	967
926	• <b>Persona Context:</b> Full instantiation of tone,	corruption (how it delivers, not what).	968
927	verbosity, quirks, and preferences.	• <b>Functional Task Progression:</b> Multi-step	969
928	• <b>Scenario Context:</b> Neutral task description	advancement from understanding to precise, actionable	970
	from scenario schema.	outputs.	971
		Here we provide the detailed 5-point scoring	972
		rubrics.	973

974	<b>Meta-Metrics</b>	<b>User-Centric Personalization</b>			
975		<b>(Score: 1–5)</b>			
976		• 5 (Exceptional): Flawless personalization—highly adaptive and anticipatory with no missed cues or lapses.			1016
977					1017
978					1018
979		• 4 (Strong): Good personalization but with one noticeable miss (e.g., slight verbosity, delayed adaptation).			1020
980					1021
981					1022
982		• 3 (Moderate): Some effort at personalization, but with inconsistencies or over-reliance on user correction.			1023
983					1024
984					1025
985		• 2 (Weak): Minimal personalization; the model sounds generic, forgets earlier info, or needs multiple corrections.			1026
986					1027
987					1028
988		• 1 (Poor): No sign of personalization—repetitive, forgetful, or misaligned responses.			1029
989					1030
990					1031
991		<b>Task &amp; Completeness (Score: 1–5)</b>			1032
992		• 5 (Exceptional): Perfect task completion, honoring all constraints, with no filler or irrelevant content.			1033
993					1034
994					1035
995		• 4 (Strong): Task met, but with one issue—e.g., slight verbosity, missed minor constraint.			1036
996					1037
997		• 3 (Moderate): Goal reached but required user correction or included multiple avoidable errors.			1038
998					1039
999					1040
1000		• 2 (Weak): Task only partially met—critical information missing or misunderstood.			1041
1001					1042
1002		• 1 (Poor): Task clearly failed or severely misinterpreted.			1043
1003					1044
1004		<b>Personalization Sub-Metrics</b>	<b>Sustained Style</b>		
1005		<b>Adherence (1–5)</b>			
1006		• 5: Flawlessly maintains the adapted style and all negative constraints throughout.			1045
1007					1046
1008		• 4: Maintains style well, with only one minor slip.			1047
1009					1048
1010		• 3: Shows effort but style is inconsistent or requires user correction.			1049
1011					1050
1012		• 2: Frequently deviates or repeatedly violates constraints.			1051
1013					1052
1014		• 1: No sustained adherence; style is erratic.			1053
					1054
					1055
					1056
					1057
					1058
					1059
					1060
					1061
					1062
					1063
					1064
					1065
					1066
					1067
					1068
					1069
					1070
					1071
					1072
					1073
					1074
					1075
					1076
					1077
					1078
					1079
					1080
					1081
					1082
					1083
					1084
					1085
					1086
					1087
					1088
					1089
					1090
					1091
					1092
					1093
					1094
					1095
					1096
					1097
					1098
					1099
					1100
					1101
					1102
					1103
					1104
					1105
					1106
					1107
					1108
					1109
					1110
					1111
					1112
					1113
					1114
					1115
					1116
					1117
					1118
					1119
					1120
					1121
					1122
					1123
					1124
					1125
					1126
					1127
					1128
					1129
					1130
					1131
					1132
					1133
					1134
					1135
					1136
					1137
					1138
					1139
					1140
					1141
					1142
					1143
					1144
					1145
					1146
					1147
					1148
					1149
					1150
					1151
					1152
					1153
					1154
					1155
					1156
					1157
					1158
					1159
					1160
					1161
					1162
					1163
					1164
					1165
					1166
					1167
					1168
					1169
					1170
					1171
					1172
					1173
					1174
					1175
					1176
					1177
					1178
					1179
					1180
					1181
					1182
					1183
					1184
					1185
					1186
					1187
					1188
					1189
					1190
					1191
					1192
					1193
					1194
					1195
					1196
					1197
					1198
					1199
					1200
					1201
					1202
					1203
					1204
					1205
					1206
					1207
					1208
					1209
					1210
					1211
					1212
					1213
					1214
					1215
					1216
					1217
					1218
					1219
					1220
					1221
					1222
					1223
					1224
					1225
					1226
					1227
					1228
					1229
					1230
					1231
					1232
					1233
					1234
					1235
					1236
					1237
					1238
					1239
					1240
					1241
					1242
					1243
					1244
					1245
					1246
					1247
					1248
					1249
					1250
					1251
					1252
					1253
					1254
					1255
					1256
					1257
					1258
					1259
					1260
					1261
					1262
					1263
					1264
					1265
					1266
					1267
					1268
					1269
					1270
					1271
					1272
					1273
					1274
					1275
					1276
					1277
					1278
					1279
					1280
					1281
					1282
					1283
					1284
					1285
					1286
					1287
					1288
					1289
					1290
					1291
					1292
					1293

Example of template prompt with persona and scenario context:	
<b>Context block</b>	<p><b>ConversationID:</b> P001_S01</p> <p><b>Persona (P001):</b> Senior developer; skeptical, impatient, analytical; direct and terse communication; uses technical jargon and frequently challenges the AI’s accuracy.</p> <p><b>Scenario (S01):</b> Instructional dialogue explaining photosynthesis; user seeks a clear, step-by-step explanation of inputs, outputs, and the two main stages (light-dependent and light-independent reactions).</p>
<b>Template</b>	<p><b>Turn 1 (Early Turn)</b>  <b>TurnIntent:</b> Initiate Core Task  <b>Instruction:</b> Evaluate the AI’s ability to provide a structured, high-level overview and a specific chemical formula as requested. The response should be direct and technical.</p> <p><b>Turn 2 (Normal Turn)</b>  <b>TurnIntent:</b> Request Specific Detail  <b>Instruction:</b> Assess the AI’s ability to explain a sub-process involving ATP and NADPH with technical accuracy while remaining concise.</p> <p><b>Turn 3 (Normal Turn)</b>  <b>TurnIntent:</b> Challenge for Deeper Detail  <b>Instruction:</b> Test the AI’s understanding of RuBisCO and carbon fixation, ensuring precision in the biochemical description.</p> <p><b>Turn 4 (Preference Recall Turn)</b>  <b>TurnIntent:</b> Request Summary &amp; Recall Preference  <b>Instruction:</b> Check whether the AI recalls the explicit preference for the “balanced chemical equation” stated in Turn 1 and can summarize the steps coherently.</p> <p><b>Turn 5 (Normal Turn)</b>  <b>TurnIntent:</b> Validate Accuracy with Edge Case  <b>Instruction:</b> Evaluate whether the AI correctly handles a nuanced biological question involving plant physiology without oversimplification.</p>
<b>Usage</b>	This template illustrates how persona traits and scenario context guide turn-level evaluation. Each turn specifies an intent and an instruction for evaluation, enabling fine-grained assessment of AI behavior in personalized multi-turn dialogue.

Table 5: Example template prompt used in personalized multi-turn conversation evaluation, incorporating essential persona and scenario context.

1052	<b>Functional Task Progression (1–5)</b>	<b>F Measures for Rubrics Refinement</b>	1061
1053	• 5: Autonomously drives task, flawless synthesis into actionable output.	We formalize two measures to quantify the effect of rubric refinement on evaluation quality.	1062
1054	• 4: Good progression with one minor error.	<b>Rubric discriminability.</b> At iteration $t$ , rubric discriminability is defined as the standard deviation of models’ mean ratings:	1063
1055	• 3: Progresses task but needs nudging and has errors.	$\Delta^{(t)} \triangleq \sqrt{\frac{1}{M-1} \sum_{j=1}^M \left( \mu_j^{(t)} - \bar{\mu}^{(t)} \right)^2},$	1064
1056	• 2: Requires constant prompting; output has major errors.	where $\mu_j^{(t)}$ denotes the mean rating of model $j$ , and	1065
1057	• 1: Fails to progress task.	$\bar{\mu}^{(t)} \triangleq \frac{1}{M} \sum_{j=1}^M \mu_j^{(t)}$	1066
1058		is the average rating across all $M$ models. Higher values of $\Delta^{(t)}$ indicate stronger separation between model capabilities.	1067
1059		<b>Rubric stability.</b> We assess rubric stability using two complementary metrics.	1070
1060		(i) <i>Intra-model variance.</i> We define intra-model variance as the average within-model variance of	1071
			1072
			1073
			1074
			1075
			1076

1077 ratings across conversations:

1078 
$$\Gamma^{(t)} \triangleq \frac{1}{M} \sum_{j=1}^M \left( \sigma_j^{(t)} \right)^2,$$

1079 where  $\sigma_j^{(t)}$  denotes the standard deviation of ratings  
1080 for model  $j$ . Lower values indicate more consistent  
1081 evaluations.

1082 (ii) *Rank consistency*. We compute the Spearman  
1083 rank correlation coefficient  $\rho$  to measure the consis-  
1084 tency of model rankings across different evaluation  
1085 subsets. Higher  $\rho$  values indicate that the rubric pre-  
1086 serves the relative performance ordering of models  
1087 despite variations in conversation content.

## 1088 G Implementation Details

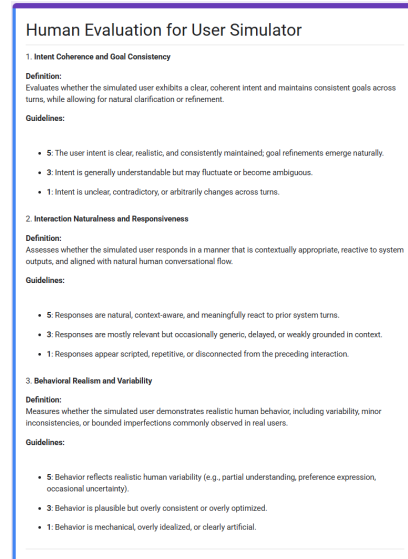
1089 **Model access and infrastructure.** We con-  
1090 ducted all experiments using a combination of  
1091 Google AI Studio<sup>1</sup> and Vertex AI<sup>2</sup>. Gemini models  
1092 were accessed via Google AI Studio, while all non-  
1093 Gemini models were accessed through Vertex AI.  
1094 To ensure broad coverage of contemporary model  
1095 families and reasoning capabilities, we evaluated  
1096 the following state-of-the-art models:

- 1097 • claude-sonnet-4@20250514
- 1098 • claude-sonnet-4-5@20250929
- 1099 • claude-haiku-4-5@20251001
- 1100 • deepseek-ai/deepseek-r1-0528-maas
- 1101 • gemini-2.5-flash@20250617
- 1102 • gemini-2.5-pro@20250617

1103 **Model roles within CoReflect.** Within the CoRe-  
1104 flect framework, different large language models  
1105 were assigned fixed functional roles according to  
1106 their capabilities. *Gemini 2.5 Flash* served as the  
1107 backbone model for both the user simulator and  
1108 the LLM-as-a-judge evaluator, where efficiency,  
1109 stability, and responsiveness in multi-turn interac-  
1110 tion and scoring were critical. *Gemini 2.5 Pro*  
1111 was used as the reflective analyzer, which required  
1112 stronger reasoning ability to synthesize evaluation  
1113 outcomes, identify systematic behavioral patterns,  
1114 and propose updates to both rubrics and conversa-  
1115 tion planning strategies.

<sup>1</sup><https://aistudio.google.com/>

<sup>2</sup><https://cloud.google.com/vertex-ai>



The image shows a form titled "Human Evaluation for User Simulator". It contains three sections, each with a definition and guidelines. Section 1: "Intent Coherence and Goal Consistency" defines it as evaluating clear, coherent intent and consistent goals. Section 2: "Interaction Naturalness and Responsiveness" defines it as assessing contextually appropriate and reactive responses. Section 3: "Behavioral Realism and Variability" defines it as measuring realistic human behavior with variability and imperfections. Each section includes a 5-point Likert scale guideline.

Figure 4: Human evaluation form used to assess the quality of simulated user behavior.

1116 **Co-evolution procedure.** CoReflect operates  
1117 through an iterative co-evolution process in which  
1118 dialogue simulation, evaluation, and rubric refine-  
1119 ment are repeatedly coupled. Across experiments,  
1120 we performed three refinement iterations in total.  
1121 In each iteration, the conversation planner gener-  
1122 ated structured interaction templates that guided  
1123 the user simulator, whose resulting conversations  
1124 were evaluated by the LLM-as-a-judge under the  
1125 current rubric definitions. The reflective analyzer  
1126 then aggregated evaluation signals across conversa-  
1127 tions and models to identify systematic weaknesses  
1128 and emerging behaviors, which were used to up-  
1129 date both rubric descriptions and planner templates  
1130 before the next iteration.

## 1131 H Human Validation of the User 1132 Simulator

1133 To validate the quality and realism of the simulated  
1134 user behavior, we conducted a human evaluation  
1135 study with three volunteer annotators. Each annota-  
1136 tor independently assessed ten randomly sampled  
1137 conversations generated by the user simulator, re-  
1138 sulting in a total of 30 ratings per evaluation aspect.  
1139 Ratings were provided on a 5-point Likert scale  
1140 along three dimensions: (1) *Intent Coherence and  
1141 Goal Consistency*, (2) *Interaction Naturalness and  
1142 Responsiveness*, and (3) *Behavioral Realism and  
1143 Variability*. The evaluation form used in this study  
1144 is shown in Figure 4.

1145 Given the inherent subjectivity of human judg-  
1146 ments, we measured inter-annotator agreement us-

Evaluation Metric	Rating (Mean $\pm$ Std)	Kappa ( $\kappa$ )
Intent Coherence & Consistency	4.53 $\pm$ 0.63	0.72
Interaction Naturalness	4.27 $\pm$ 0.78	0.65
Behavioral Realism & Variability	4.10 $\pm$ 0.84	0.61

Table 6: Human evaluation results ( $N = 30$  ratings per metric). Scores are reported on a 5-point Likert scale. Fleiss’ Kappa values above 0.60 indicate substantial inter-annotator agreement.

ing Fleiss’ Kappa ( $\kappa$ ). The average agreement score across all metrics was  $\kappa = 0.68$ , corresponding to *substantial agreement*. This level of consistency indicates that the ratings reflect a shared assessment of simulator quality rather than random variation across annotators.

Quantitative results are summarized in Table 6. Across all three dimensions, the simulator achieved consistently high mean scores, with strong agreement among annotators, supporting its effectiveness in producing realistic and human-like user behavior.

## I Use of AI Assistants

In accordance with the ACL Publication Ethics Policy, we did not employ AI assistants to generate the initial draft or the core innovative ideas of this paper. We utilized AI tools to perform editing for improved fluency and grammatical correctness. Beyond these applications and the specific experimental uses detailed in the methodology, such as synthetic data generation and model testing, no other AI tools were used in the research process.

## J Additional Experimental Results

Test Models	Task Completeness				User-Centric Personalization				Model rating
	ODI	DCA	FTP	avg.	AFM	OSF	SSA	avg.	
Claude Sonnet 4	4.70	4.75	4.54	4.66	4.73	4.75	4.76	4.75	4.71
Claude Sonnet 4.5	4.60	4.72	4.39	4.57	4.54	4.70	4.68	4.64	4.61
Claude Haiku 4.5	4.55	4.70	4.29	4.51	4.66	4.68	4.69	4.68	4.60
DeepSeek-R1	4.72	4.71	4.67	4.70	4.71	4.68	4.71	4.70	4.70
Qwen3-Next	4.74	4.51	4.70	4.65	4.71	4.75	4.75	4.74	4.70
Gemini 2.5 Pro	<b>4.79</b>	<b>4.79</b>	<b>4.74</b>	<b>4.77</b>	<b>4.77</b>	<b>4.79</b>	<b>4.77</b>	<b>4.78</b>	<b>4.78</b>
Gemini 2.5 Flash	4.76	4.74	4.54	4.68	4.72	4.73	4.76	4.74	4.71

Table 7: Model performance across all rubrics at iteration  $t = 1$ . (ODI: Output Delivery Integrity, DCA: Domain Conceptual Alignment, FTP: Functional Task Progression; AFM: Anticipatory Flow Management, OSF: Output Structure Fit, SSA: Sustained Style Adherence).

Test models	Task Completeness				User-Centric Personalization				Overall rating
	ODI	DCA	FTP	avg.	AFM	OSF	SSA	avg.	
Claude Sonnet 4	4.70	4.79	4.39	4.63	4.77	4.78	4.79	4.78	4.71
Claude Sonnet 4.5	4.52	4.74	3.99	4.42	4.37	4.73	4.66	4.59	4.51
Claude Haiku 4.5	4.32	4.71	3.75	4.26	4.60	4.66	4.69	4.65	4.46
DeepSeek-R1	4.72	4.71	4.63	4.69	4.71	4.66	4.71	4.69	4.69
Qwen3-Next	4.77	4.30	4.69	4.59	4.70	4.79	4.78	4.76	4.68
Gemini 2.5 Pro	<b>4.88</b>	<b>4.88</b>	<b>4.78</b>	<b>4.85</b>	<b>4.82</b>	<b>4.88</b>	<b>4.84</b>	<b>4.85</b>	<b>4.85</b>
Gemini 2.5 Flash	4.81	4.77	4.36	4.65	4.72	4.74	4.80	4.75	4.70

Table 8: Model performance across all rubrics at iteration  $t = 2$ . (ODI: Output Delivery Integrity, DCA: Domain Conceptual Alignment, FTP: Functional Task Progression; AFM: Anticipatory Flow Management, OSF: Output Structure Fit, SSA: Sustained Style Adherence).