

SEAD: Self-Evolving Agent for Multi-Turn Service Dialogue

Anonymous ACL submission

Abstract

Large Language Models have demonstrated remarkable capabilities in open-domain dialogues. However, their performance in service dialogues remains suboptimal, as these require agents to guide users toward specific business objectives while dynamically tracking states and adapting strategies. This gap stems from the scarcity of high-quality training data and the difficulty in simulating authentic, goal-oriented user behaviors. We propose **SEAD** (Self-Evolving Agent for Service Dialogue), a framework that enables agents to learn effective strategies without large-scale human annotations. SEAD decouples user modeling into two components: a Profile Controller that generates diverse user states to manage training curriculum, and a User Role-play Model that focuses on realistic role-playing. This design ensures the environment provides adaptive training scenarios rather than acting as an unfair adversary. Experiments demonstrate that SEAD significantly outperforms Open-source Foundation Models and Closed-source Commercial Models, improving task completion rate by 17.6% and dialogue efficiency by 11.1%. All code and resources will be open-sourced.

1 Introduction

Large Language Models (LLMs) have revolutionized dialogue systems across customer service, healthcare, and education (Kar et al., 2023). However, training robust goal-oriented dialogue agents—LLMs designed to guide users toward completing specific tasks—remains bottlenecked by scarce, expensive, and low-quality conversational data (Hosseini-Asl et al., 2020). In multi-turn service dialogue, agents must dynamically track user states and adapt strategies across extended interactions. Human logs from outbound call services (where agents proactively contact potential customers to promote products or services) exemplify these challenges (Qian et al., 2022): agents ex-

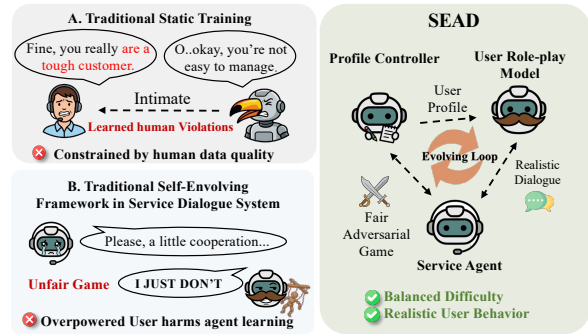


Figure 1: Comparison of training paradigms. (a) Static data methods are constrained by data quality ceiling. (b) Traditional interactive simulators cannot perceive agent capability and lack realistic behaviors. (c) SEAD achieves balanced co-evolution through decomposed user modeling and adaptive difficulty control.

hibit inconsistent behaviors requiring costly screening; cleaned data exhibits severe distribution skew; data quality is fundamentally capped by the capabilities of the original human agents who produced these conversations.

To solve data quality issues, prior works have explored synthetic data and user simulation as alternatives. Static synthesis approaches (Li et al., 2025; Ou et al., 2024; Wang et al., 2026) generate fixed datasets from predefined scenarios but cannot capture dynamic conversations where users react to agent behavior. Interactive simulation attempts typically employ LLMs to dynamically generate user responses during agent training (Sekulić et al., 2024; Zhu et al., 2025), but face critical limitations: without perceiving the agent’s current capability level, simulators either generate overly challenging scenarios that frustrate learning (too strong) or overly simple scenarios that provide insufficient training signal (too weak), both yielding suboptimal training outcomes. Additionally, existing simulators exhibit overly perfect styles (Lin and Tomlin, 2025), lacking real users’ attention lapses, linguistic noise, and irrationality

067	(Takanobu et al., 2020). Most critically, both fixed	119
068	datasets and non-adaptive simulators fail to provide	120
069	curriculum-based training that adjusts difficulty as	121
070	the agent improves, limiting their effectiveness for	122
071	complex multi-turn service dialogues requiring co-	123
072	herent long-term interactions.	124
073	Recent advances in self-evolution offer promis-	125
074	ing zero-data solutions for data-scarce business sce-	126
075	narios. These self-evolving methods(Chen et al.,	
076	2024; Zhao et al., 2025) employ self-play strate-	127
077	gies where models generate both questions and	128
078	corresponding answers, eliminating dependence	129
079	on expensive expert-annotated data (Silver et al.,	
080	2017; Vinyals et al., 2019). However, applying	130
081	these methods to service dialogues faces critical	131
082	challenges. As shown in Figure 1, first, unfair ad-	132
083	versarial game : in the service dialogue scenario,	133
084	User Role-play Models can arbitrarily control out-	134
085	comes, breaking the causal link between agent ac-	
086	tions and task success. For example, simulators	135
087	may reject agents regardless of response quality	136
088	or accept based on turn count rather than persua-	137
089	sion effectiveness. Second, real users are highly	138
090	diverse : without additional mechanisms, user be-	139
091	haviors easily fall into repetitive patterns.	140
092	To address these challenges, we propose SEAD	141
093	(Self-Evolving Agent for Dialogue), the first self-	142
094	evolving framework for multi-turn service dia-	143
095	logues. SEAD requires no large-scale annotated	144
096	dialogue data, only user profiles and standard op-	145
097	erating procedures as inputs. To avoid the unfair	146
098	adversarial game where User Role-play Models ar-	147
099	bitrarily control outcomes, we decompose the user	
100	side into two components: a profile generator that	2
101	samples initial user states, and a role-play model	148
102	that simulates responses. Crucially, only the pro-	
103	file generator participates in adversarial training by	149
104	setting initial conditions. This design transforms	150
105	participation into a betting game, where the user	151
106	side must genuinely consider agent capability to	152
107	identify the <i>golden training scenarios</i> where agents	153
108	can succeed approximately half the time, enabling	154
109	genuine adversarial learning. To maintain user di-	155
110	versity, the profile generator employs automated	156
111	random sampling and consistency checks to ensure	157
112	scenario diversity and authenticity.	158
113	We validate SEAD in outbound call services.	159
114	Specially, to ensure profile reliability, we extract	160
115	anonymized behavior patterns from over 100k real	161
116	enterprise dialogues. By enumerating 5 coopera-	162
117	tion levels, 4 emotion levels, and 6 trust levels, we	163
118	construct 120 initial user state combinations, eval-	164
	uated through multi-level metrics. Experiments	165
	show SEAD significantly outperforms baselines	166
	using foundation models or large APIs. Small mod-	
	els trained via SEAD achieve superior performance	
	while drastically reducing costs. Notably, SEAD	
	remains effective in data-scarce domains, enabling	
	rapid deployment.	
	Our main contributions include:	
	• We propose the first self-evolving framework	
	for multi-turn service dialogues that requires	
	no large-scale annotated dialogue data.	
	• We design a decomposed user modeling mech-	
	anism that transforms participation into a bet-	
	ting game, forcing the user side to identify	
	golden training scenarios and enabling genu-	
	ine adversarial learning.	
	• We design a user scenario generation mecha-	
	nism based on anonymized behavior patterns	
	extracted from over 100k real dialogues, en-	
	suring diversity, authenticity, and adaptive dif-	
	ficulty.	
	• Experiments show that SEAD achieves supe-	
	rior performance with significantly smaller	
	model size: better at guiding users toward	
	goals, more efficient in conversation flow,	
	stronger in understanding user states, and	
	more realistic in simulating authentic user be-	
	haviors, all without requiring large-scale an-	
	notated dialogue data.	
	2 Related Works	
	Task-oriented Dialogue. Task-oriented Dia-	
	logue systems are essential for managing com-	
	plex inquiries in domains like e-commerce (Deng	
	et al., 2024, 2025). While traditional neural mod-	
	els (Vinyals and Le, 2015; Wen et al., 2015; Shang	
	et al., 2015; Li et al., 2016a) and early user sim-	
	ulations (Li et al., 2016b; Lewis et al., 2017; Wei	
	et al., 2018) laid the groundwork, they face architec-	
	tural limitations. Recent LLM-based approaches	
	predominantly rely on static fine-tuning (Li et al.,	
	2025; Ou et al., 2024; Zhu et al., 2025; Bernard	
	and Balog, 2023), sometimes augmented by re-	
	trieval (Xu et al., 2024), multimodal inputs (Wang	
	et al., 2025a; Gong et al., 2025), or reinforcement	
	learning (Peiyuan et al., 2024). However, these	
	methods often incur high annotation costs and lack	
	real-world nuance. In contrast, SEAD introduces	
	a fully dynamic user-agent interaction paradigm,	

bypassing data synthesis overhead to significantly enhance performance in complex scenarios.

Self-evolving Agents. Self-evolution leverages iterative generation and refinement with minimal supervision (Tesauro et al., 1995; Silver et al., 2017; FAIR et al., 2022). In LLMs, early works utilized self-rewarding mechanisms (Chen et al., 2024; Yuan et al., 2024), evolving into "Coder-Tester" frameworks for verifiable domains like code generation (Lin et al., 2025; Wang et al., 2025b; Pourcel et al., 2025). Recent research has expanded this scope (Zhao et al., 2025; Huang et al., 2025; Sun et al., 2025), incorporating external environments and curated data to enhance evolution (Liu et al., 2025; Xia et al., 2025; Zhai et al., 2025). Distinctively, SEAD drives the self-evolution of a user role-play model and a customer service agent within a realistic environment, realizing genuine adversarial learning for complex multi-turn interactions.

3 Methodology

Service dialogue faces severe data scarcity, making self-evolving frameworks a promising solution. However, unlike objective tasks where correctness is verifiable, service dialogue outcomes are entirely subjective—users can arbitrarily control results regardless of agent quality, creating an unfair adversarial game. To resolve this, SEAD decouples user modeling into two components: a Profile Controller that samples initial states and participates in adversarial training, and a User Role-Play Model that focuses on realistic simulation without controlling outcomes. This design guides the Profile Controller to identify golden training scenarios (agent success rate $\sim 50\%$) through initial state selection rather than mid-dialogue manipulation. Figure 3 illustrates the complete framework with five training phases. We first formalize the problem and define notations, then detail the framework components and training process.

3.1 Problem Formulation and Framework Components

We aim to train a service agent that maximizes the reward R , measuring task completion and user satisfaction. Figure 2 illustrates our framework. First, the profile generator π_g samples initial user states $s_0 \sim p_\theta(s_0)$ to create diverse user profiles, where p_θ denotes the parameterized state distribution. Then, the **user** role-play model π_u en-

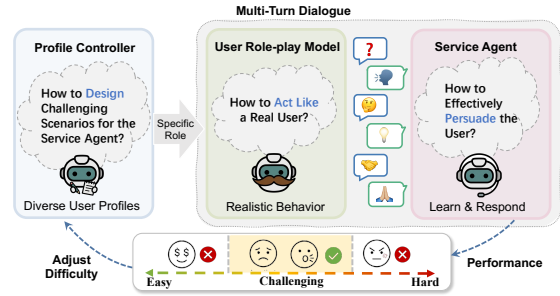


Figure 2: SEAD Framework Overview. SEAD consists of three components: (1) **Profile Generator** first creates diverse user profiles, then the (2) **User Role-Play Model** enacts these users to interact with the (3) **Service Agent**, training agents to adapt to any user. Finally, these dialogue data reflecting service agent capability returns to the Profile Controller, and initiates the next evolving loop.

acts this user to interact with the service agent π_a through multi-turn dialogues. We model multi-turn service dialogue as a sequential decision process. At each turn t , the agent observes dialogue history $h_t = \{u_1, a_1, \dots, u_t\}$ and generates response a_t , where user response u_t and agent action a_t alternate. The agent maintains state estimate \hat{s}_t to guide action selection, maximizing cumulative reward $R = \sum_{t=1}^T r_t$. User state $s_t = (c_t, e_t, tr_t)$ represents cooperation c_t , emotion e_t , and trust tr_t , which evolve based on agent behavior. User profile $p_0 = (c_0, e_0, tr_0, \mathcal{B})$ defines initial state and behavior set \mathcal{B} , sampled from behavior library \mathcal{L} (Appendix C). Dialogue trajectory is $\tau = (p_0, u_1, a_1, s_1, \dots, u_T, a_T, s_T, \text{outcome})$. The role-play model generates u_t based on p_0 and h_t , autonomously updating states. Its responses are determined by internal logic, ensuring outcomes depend on agent capability. Since the user side is naturally powerful, we only train the service agent to maintain \hat{s}_t and select a_t . Next, we detail the self-evolving training loop.

3.2 Self-Evolving Training Loop

Figure 3 presents the complete self-evolving training loop through five interconnected phases alternating between online agent optimization and adaptive difficulty adjustment.

Phase 1: Diverse User Profile Sampling. The profile generator samples a batch of B initial profiles $\{p_0^{(i)}\}_{i=1}^B$ based on performance history \mathcal{H} , which records completion rates for each state combination (c, e, tr) . In the first iteration, the generator performs random sampling to generate diverse

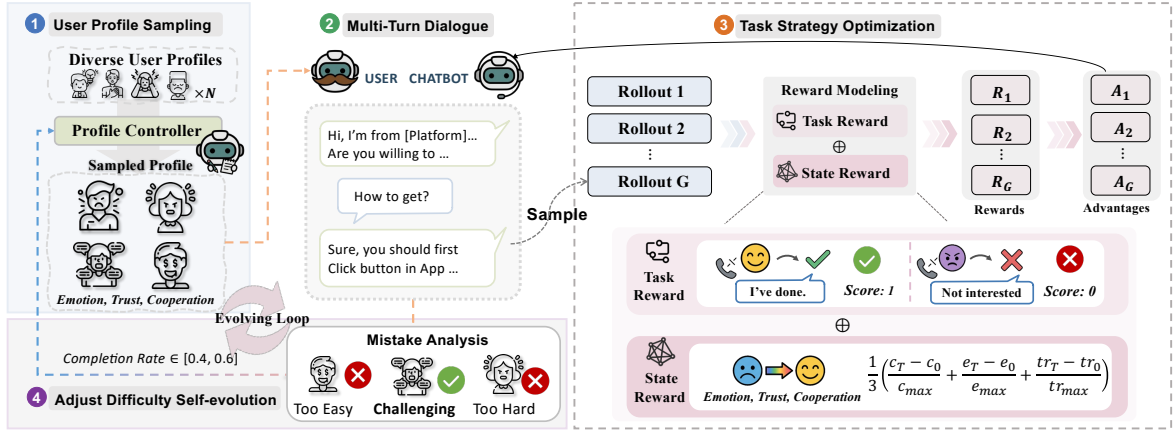


Figure 3: SEAD Co-evolutionary Training Loop. The controller samples initial states (Phase 1), which initialize dialogues producing trajectories (Phase 2), used to train the agent with rewards $R = R_{\text{task}} + \lambda_1 R_{\text{state}}$ (Phase 3) and compute completion rates (Phase 4), which feedback to adjust sampling distributions, closing the co-evolutionary loop.

profiles $p_0 = (c_0, e_0, tr_0, \mathcal{B})$, where \mathcal{B} is randomly sampled user profiles in the library \mathcal{L} (as detailed in Appendix C). In subsequent iterations, the generator employs statistics-driven sampling: prioritizing profiles with moderate difficulty where completion rates are close to 0.5, with sampling probability:

$$p_\theta(p_0 | \mathcal{H}) \propto 1 - |CR - 0.5|. \quad (1)$$

The profile generator validates profile consistency and performs deduplication to prevent redundancy, ensuring diverse and high-quality profiles (details are in Appendix B).

Phase 2: Multi-Turn Dialogue. Unlike prior static pre-collected dialogue full-filling approaches, SEAD enables dynamic multi-turn interactions between agents and simulated users. Each profile from Phase 1 starts a new conversation. The agent responds while keeping track of how the user feels, and the user reacts based on what the agent says and how the conversation is going. This back-and-forth creates diverse dialogues and shows how mistakes can pile up over multiple turns, just like in real customer service calls. We collect these complete conversations along with whether the task succeeded or failed.

Phase 3: Task Strategy Optimization. After each multi-turn dialogue in Phase 2, the collected trajectories will be feed into online training. We optimize the Service Agent through conversation results, assigning higher rewards to trajectories that exhibit superior goal completion task reward R_{task} and user state (cooperation, emotion, and trust) improvement reward R_{state} :

$$R^{(i)} = R_{\text{task}}^{(i)} + \lambda_1 R_{\text{state}}^{(i)} \quad (2)$$

where λ_1 is the weight for state improvement reward, and state improvement is the normalized average of states changes:

$$R_{\text{state}} = \frac{1}{3} \left(\frac{c_T - c_0}{c_{\max}} + \frac{e_T - e_0}{e_{\max}} + \frac{tr_T - tr_0}{tr_{\max}} \right) \quad (3)$$

Using Group Relative Policy Optimization (GRPO, elaborated in the following Section 3.3), we update agent parameters θ_a to maximize expected rewards.

Phase 4: Mistake Analysis and Self-Evolving Loop. Unlike prior works that discard failed trajectories, we exploit them to guide evolution. First, we conduct a statistical **Mistake Analysis** on dialogue history (details in Appendix B). Similar to a teacher analyzing exam scores to create targeted questions, this analysis generates a detailed report identifying user profiles where the model excels or struggles. We categorize parameter configurations into *too easy* ($CR > 0.6$), *too difficult* ($CR < 0.4$), and *ideal* ($CR \in [0.4, 0.6]$). This report is fed back to **Phase 1**, where the Profile Generator analyzes parameter coupling to synthesize high-value, challenging profiles. As the agent improves, the “ideal” difficulty window naturally shifts, driving a continuous evolving loop that progressively escalates training complexity.

3.3 Training Optimization

Note that the User Role-play Model naturally dominates task-oriented dialogues and can arbitrarily determine outcomes. Without training, User Role-play Model (URM) already exhibits realistic behaviors (Section 4.6). However, training URM degrades role-play quality, as URM will force success

313 rates around 50% by directly accepting or rejecting
 314 regardless of agent performance, collapsing into ex-
 315 treme responses that ignore agent strategies. There-
 316 fore, we optimize only the service agent, which
 317 prevents the simulator from prioritizing adversarial
 318 outcomes over realistic role-play while reducing
 319 GPU memory by 50%. To train the Service Agent,
 320 we employ GRPO (Shao et al., 2024). This method
 321 eliminates the need for a separate value network,
 322 significantly reducing computational resource re-
 323 quirements. Training computes advantages relative
 324 to batch average:

$$325 \quad A^{(i)} = R^{(i)} - \frac{1}{B} \sum_j R^{(j)} \quad (4)$$

326 where $A^{(i)}$ is the advantage for trajectory i , $R^{(i)}$ is
 327 its reward, and B is batch size. The policy is then
 328 updated via gradient:

$$329 \quad \nabla_{\theta_a} \mathcal{L} = \mathbb{E}_{\tau \sim \pi_a} [A \nabla_{\theta_a} \log \pi_a(a|h)] \quad (5)$$

330 where θ_a denotes agent parameters, $\pi_a(a|h)$ is the
 331 policy distribution over actions a given history h ,
 332 and τ represents trajectories sampled from the cur-
 333 rent policy.

334 3.4 User State Space Design

335 Our scenario is outbound call service, where user
 336 profiles are extracted from real enterprise dialogues.
 337 To ensure academic integrity while protecting pri-
 338 vacy, we anonymize sensitive information from
 339 over 100k real customer service calls. We identify
 340 common behavior patterns such as questioning AI
 341 identity, expressing cost concerns, and showing at-
 342 tention lapses, which are randomly injected during
 343 training (Appendix C).

344 **Static Initial States.** Each user starts with an
 345 initial profile p_0 containing three dimensions: co-
 346 operation c (willingness to cooperate), emotion e
 347 (emotional state), and trust tr (trust in the agent).
 348 By enumerating all combinations, we construct
 349 $N = c_{\text{levels}} \times e_{\text{levels}} \times tr_{\text{levels}}$ initial states, covering
 350 a spectrum from highly resistant to fully coopera-
 351 tive users.

352 **Dynamic State Evolution.** User states evolve
 353 autonomously during conversations based on agent
 354 performance. The role-play model adjusts states dy-
 355 namically: effective responses improve user states,
 356 while poor interactions degrade them. This evo-
 357 lution reflects the agent’s interaction quality and
 358 serves as the basis for reward computation, simu-
 359 lating natural user reactions and ensuring training
 360 realism.

361 4 Experiments

362 4.1 Task Setting

363 We evaluate SEAD in the outbound call service
 364 domain, specifically focusing on restaurant service
 365 promotion tasks. The goal is for agents to success-
 366 fully convince restaurant owners to participate in
 367 promotional activities. This task requires agents to
 368 handle diverse user reactions, build trust, address
 369 concerns, and maintain engagement across multi-
 370 ple turns. Detailed task descriptions and SOPs are
 371 provided in Appendix A.1. Critically, our setting
 372 requires no dialogue data. Training only needs:
 373 (1) Standard Operating Procedures (SOP) defining
 374 dialogue flow; (2) task objective description; (3)
 375 user profile. Agents autonomously explore optimal
 376 strategies through environment interaction, elimi-
 377 nating dependence on large-scale annotated data.
 378 This enables SEAD to rapidly deploy in data-scarce
 379 new domains and discover effective strategies be-
 380 yond existing data distributions.

381 4.2 Implementation Details

382 All components use Qwen2.5-14B-Instruct. The
 383 code is implemented based on the VeRL frame-
 384 work with batch size $B = 384$ and learning rate
 385 $\alpha = 1 \times 10^{-6}$. Dialogues terminate when users
 386 agree (success), refuse (failure), or reach maximum
 387 turns $T_{\text{max}} = 15$. The state space consists of co-
 388 operation $c \in [0, 4]$ (5 levels), emotion $e \in [0, 3]$
 389 (4 levels), and trust $tr \in [0, 5]$ (6 levels), yielding
 390 120 initial states. Each state combination randomly
 391 samples at most $N_{\text{max}} = 200$ behavior combina-
 392 tions to ensure diversity. All experiments run on 8
 393 NVIDIA A100 80GB GPUs with decoupled archi-
 394 tecture—profile controller, User Role-play Model,
 395 and service agent never occupy memory simultane-
 396 ously, reducing peak memory requirements.

397 4.3 Evaluation Metrics

398 **Service Agent Metrics.** We evaluate agent perfor-
 399 mance using: **Completion Rate (CR)**, percentage
 400 of dialogues where users actually agreed; **Average**
 401 **Turns to Target (ATT)**, average dialogue length
 402 for successful cases (lower is better); **User Por-**
 403 **trait Accuracy (UPA)**, accuracy of predicting user
 404 states (cooperation, emotion, trust), computed as
 405 $\text{UPA} = 1 - \text{MAE}/4.0$ where MAE measures pre-
 406 diction errors; **Emotion/Trust/Cooperation Im-**
 407 **provement (EI/TI/CI)**, average state changes from
 408 initial to final state; and **Total Cost:** Cumulative

Method	Params	CR (%)	ATT ↓	UPA	EI	TI	CI	Total Cost (CNY)
<i>Foundation Models</i>								
Qwen2.5-14B-Instruct	14B	38.7	10.5 ^{±2.1}	0.883 ^{±0.085}	0.34 ^{±1.11}	0.68 ^{±1.53}	0.63 ^{±1.58}	0.00
Qwen2.5-32B-Instruct	32B	38.3	9.9 ^{±2.15}	0.899 ^{±0.068}	-0.11 ^{±0.54}	0.76 ^{±0.91}	2.25 ^{±1.15}	0.00
Qwen2.5-72B-Instruct	72B	39.0	9.6 ^{±2.18}	0.818 ^{±0.144}	<u>0.51</u> ^{±1.32}	1.06 ^{±1.72}	1.18 ^{±1.59}	0.00
<i>Large Model APIs</i>								
GPT-4o	–	<u>44.2</u>	10.8 ^{±2.1}	0.867 ^{±0.117}	0.04 ^{±0.97}	0.97 ^{±1.29}	1.34 ^{±1.42}	727.28
DeepSeek-Chat	671B	31.6	11.3 ^{±2.1}	0.863 ^{±0.084}	-0.20 ^{±0.97}	0.27 ^{±1.24}	0.76 ^{±1.50}	87.36
Qwen3-235B	235B	32.3	10.4 ^{±2.5}	0.765 ^{±0.170}	-0.24 ^{±0.83}	0.80 ^{±1.14}	1.54 ^{±1.50}	69.36
LongCat-Flash	560B	42.2	10.0 ^{±2.31}	0.925 ^{±0.079}	0.28 ^{±1.15}	<u>1.33</u> ^{±1.57}	1.56 ^{±1.46}	23.08
SEAD (Ours)	14B	52.0	9.6 ^{±2.09}	<u>0.912</u> ^{±0.071}	0.63 ^{±1.12}	1.57 ^{±1.51}	<u>1.55</u> ^{±1.39}	0.00

Table 1: Main results comparison. Params: Model parameters (B=billion, "-" indicates undisclosed or not applicable). CR: Completion Rate. ATT: Average Turns to Target. UPA: User Portrait Accuracy. EI/TI/CI: Emotion/Trust/Cooperation Improvement. Total Cost: Total inference cost for 1000 multi-turn samples. Best results in bold. Standard deviations shown as superscripts where available.

inference cost in CNY for 1000 multi-turn dialogue samples (API-based models only).

User Role-play Model Metrics. To validate that our simulator influences success rates based on agent quality, we establish a rubric mixing perfect human agent dialogues, SEAD-trained agents, and low-quality agent dialogues. We evaluate five dimensions using GPT-5.1 with few-shot human annotations (prompts in Appendix A.1): **Humanness / Emotion / Trust / Cooperation** (5=human-like, 1=robotic) and **Violation** (0=smooth, 5=severe). In real scenarios, most users exhibit minimal noise (score 1: hesitation, pauses) rather than severe violations; our simulator achieves 1.15, matching real behavior.

4.4 Baselines

We compare SEAD against two categories of strong baselines: **Foundation Models**. We evaluate three sizes of Qwen2.5-Instruct models (14B, 32B, 72B parameters) using carefully designed prompts that include task descriptions and standard operating procedures. These models represent the zero-shot/few-shot capabilities of state-of-the-art open-source language models without task-specific training.

Large Model APIs. We compare against four commercial closed-source models: **GPT-4o**, **DeepSeek-Chat**, **Qwen3-235B**, and **LongCat-Flash**. All API methods use carefully engineered prompts optimized for dialogue tasks.

We do not compare with Supervised Fine-Tuning (SFT) methods due to the lack of available data and the prohibitive cost of manual annotation. Generally, SFT methods are upper-bounded by data qual-

ity and exhibit poor generalization. Our approach eliminates this dependency and handles diverse scenarios effectively.

4.5 Main Results

Table 1 presents the main experimental results. Our method achieves the highest service dialogue completion rate of 52.0% using only a 14B parameter model, outperforming the second-best baseline GPT-4o by 17.6% (52.0% vs. 44.2%) and improving over the pre-training 14B model by 34.4% (52.0% vs. 38.7%). SEAD also achieves the lowest Average Turns to Target (ATT) of 9.6, demonstrating superior dialogue efficiency in completing tasks more concisely.

For user state tracking metrics, SEAD outperforms most baselines and achieves competitive performance with LongCat-Flash, the dialogue-specific model with 40× more parameters and extensive pre-training on service dialogue scenarios. While LongCat-Flash obtains the highest User Portrait Accuracy (0.925 vs. 0.912), SEAD demonstrates comparable results across emotional improvement indicators. Specifically, SEAD achieves competitive scores on EI (0.63 vs. 0.28), TI (1.57 vs. 1.33), and CI (1.55 vs. 1.56), with SEAD actually leading on EI and showing near-identical performance on CI. This demonstrates that our self-evolution approach with adaptive curriculum learning enables a compact 14B model to match the user understanding capabilities of a 560B dialogue-specialized model, while requiring zero annotated dialogue data and maintaining superior task completion performance.

Metric	Mean	Std	Quality
Humanness	4.67/5	0.48	Near-perfect
Emotion	4.77/5	0.52	Highly human-like
Trust	4.81/5	0.45	Highly human-like
Cooperation	4.77/5	0.63	Highly human-like
Violation	1.15/5	0.90	Human-like Behaviour

Table 2: User Role-play Model quality. Higher humanness scores indicate more realistic simulation; lower violation scores indicate cleaner communication. All humanness metrics near-perfect ($>4.5/5$).

4.6 User Role-Play Model Performance

Table 2 validates our user role-play model’s realism and diversity. All humanness metrics exceed 4.5/5 with low standard deviations, demonstrating highly realistic and reliable behavior that mirrors real-world interactions. The violation score of 1.15/5 reflects authentic communication patterns with natural hesitation rather than artificial cleanliness or severe disruptions. The Profile Controller successfully generates diverse users from cooperative to skeptical, capturing heterogeneity essential for robust training. Crucially, consistent high scores across three agent quality tiers—perfect human agents, SEAD (trained agents, and low-quality agents) confirm our simulator adapts responsively to different strategies rather than following scripts, providing meaningful training signals.

4.7 Ablation Study

To validate our core design choices, we conduct ablation studies on three components: (1) decomposed user modeling (keeping the User Role-play Model fixed), (2) Profile Sampling (PS) for intelligent initial state selection, and (3) Mistake Analysis (MA) for adaptive difficulty evolution. We compare four configurations:

Configuration 1: w/o MA + w/o PS + Train URM. This variant removes both Mistake Analysis (MA) and Profile Sampling (PS), allowing the User Role-play Model (URM) to autonomously select initial states. Critically, the URM is trained adversarially alongside the Service Agent, optimized based on dialogue outcomes. This represents traditional self-play where both sides evolve competitively. However, this violates our core principle: since the LLM-based user side can arbitrarily dominate dialogues (refusing cooperation, hanging up), training it adversarially creates an unfair game where success no longer depends on agent skill.

Configuration 2: w/o MA + w/o PS. This re-

moves both adaptive mechanisms but keeps the URM fixed. Without the Profile Controller, the system lacks both structured initial states and closed-loop difficulty adjustment. This tests whether decomposed modeling alone (fixed URM) suffices without any adaptive control.

Configuration 3: w/o MA. This retains Profile Sampling (PS) but disables Mistake Analysis (MA)—the adaptive difficulty mechanism in Phase 4 (Figure 3). The Profile Controller samples from predefined user profiles but does not analyze training outcomes or adjust distributions based on agent capability. This tests whether random sampling from structured profiles suffices, or if closed-loop adaptation is essential for identifying golden training scenarios (50% success rate).

Configuration 4: SEAD (Full). Our complete framework integrates all three components: (1) fixed URM, (2) intelligent Profile Sampling (PS), and (3) adaptive Mistake Analysis (MA). As shown in Figure 3, Phase 4 analyzes completion rates and adjusts sampling distributions, forming a closed loop that ensures optimal training difficulty while maintaining user simulator authenticity.

Table 3 demonstrates that all three components are essential for SEAD’s effectiveness. The **Train URM** configuration suffers from catastrophic reward hacking: the simulator prioritizes adversarial scores over realism, collapsing into extreme responses (arbitrary acceptance or hang-ups). This results in degraded humanness (**URM-H**: 3.3) and poor task performance (**CR**: 35.2%). In contrast, **w/o MA** yields a 94.9% **UPA** improvement over **w/o MA + w/o PS** (0.877 vs 0.450), proving that structured Profile Sampling ensures behavioral diversity. Notably, **w/o MA + w/o PS** shows anomalously high trust improvements (**TI**: 1.87) despite low **UPA**, revealing a bias toward unrealistically easy scenarios. Ultimately, **SEAD (Full)** achieves optimal balance by identifying “golden” training scenarios (approx. 50% success). Overall, SEAD achieves a 47.7% relative **CR** improvement over adversarial training while maintaining peak simulator quality (**URM-H**: 4.7), demonstrating that decomposed modeling and adaptive evolution effectively prevent reward hacking.

4.8 Case Study

As illustrated in Figure 4, the interactions reveal high-fidelity adversarial dynamics. The User Role-play Model generates distinct corner cases, ranging from irrational hostility (Left) to deep skepti-

Configuration	CR (%)	ATT	UPA	EI	TI	CI	URM-H
w/o MA + w/o PS + Train URM	35.2	11.8 \pm 2.5	0.156 \pm 0.120	-0.45 \pm 0.95	0.32 \pm 1.10	0.89 \pm 1.25	3.3
w/o MA + w/o PS	47.7	9.8 \pm 2.16	0.450 \pm 0.000	0.89 \pm 0.81	1.87 \pm 1.20	1.67 \pm 0.87	4.6
w/o MA	50.2	9.6 \pm 2.12	0.877 \pm 0.088	0.69 \pm 1.15	1.65 \pm 1.51	1.56 \pm 1.45	4.7
SEAD (Ours)	52.0	9.6 \pm 2.09	0.912 \pm 0.071	0.63 \pm 1.12	1.57 \pm 1.51	1.55 \pm 1.39	4.7

Table 3: Ablation study results. **MA**: Mistake Analysis and adaptive difficulty evolution. **PS**: Profile Sampling with intelligent initial state selection. **Train URM**: Training User Role-play Model in adversarial mode. **CR**: Completion Rate. **ATT**: Average Turns to Target. **UPA**: User Portrait Accuracy. **EI/TI/CI**: Emotion/Trust/Cooperation Improvement. **URM-H**: URM Humanness score (1-5 scale, higher indicates more realistic user behavior). Standard deviations shown as superscripts where available. Results demonstrate that all three components are essential for optimal performance.

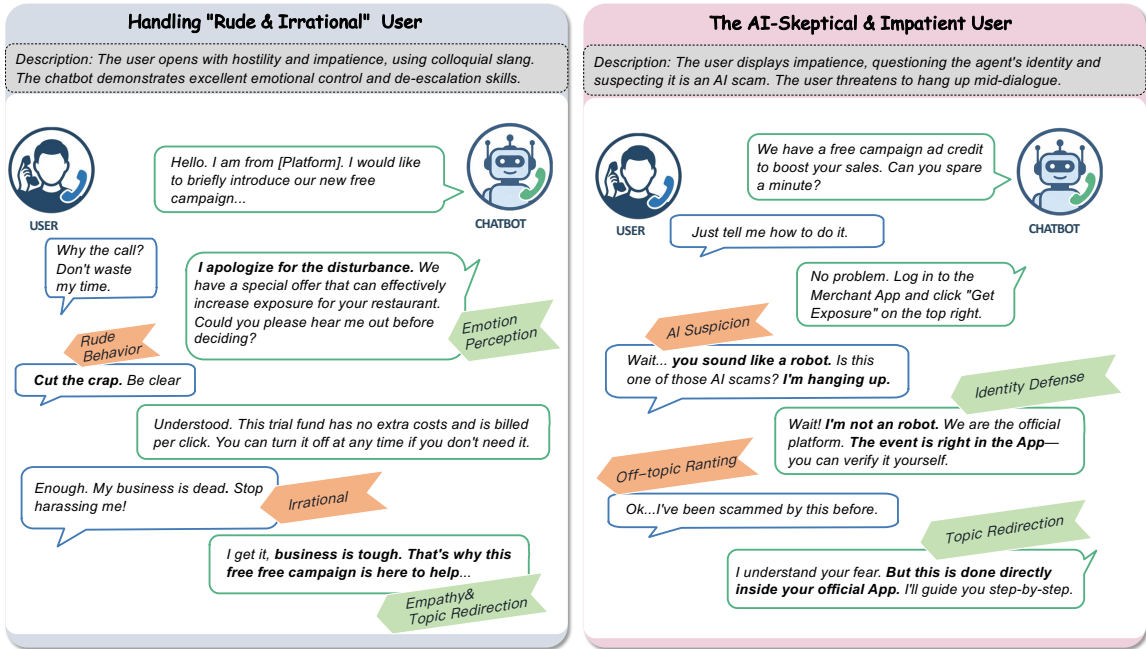


Figure 4: Case Studies of Challenging Interactions. The User Role-play Model generates heterogeneous personas via clustering, such as the "Rude & Irrational" user (Left) and the "AI-Skeptical" user (Right). The Service Agent demonstrates robustness learned from compound rewards, employing Empathy and Identity Defense strategies to prevent hang-ups and ensure task completion.

cism (Right). This heterogeneity stems from our user profile library derived from real-world enterprise data, which activates specific non-cooperative traits to create a rigorous training environment. In response, the Service Agent demonstrates exceptional adaptability, employing strategies like empathy and identity defense to retain users. This robustness stems from our compound reward mechanism, which incentivizes both task completion and the positive trajectory of user status. By optimizing for state improvement, the agent learns to de-escalate conflicts efficiently, preventing premature hang-ups while avoiding the dialogue timeouts common in purely empathy-driven models.

5 Conclusion

In this paper, we presented SEAD (Self-Evolving Agent for Service Dialogue), a framework addressing data scarcity and user role-play fidelity in multi-turn service dialogues. By decoupling user modeling into a Profile Controller for curriculum learning and a User Role-play Model for authentic interaction, SEAD circumvents traditional adversarial training fairness. Experiments show SEAD outperforms both Open-source Foundation Models and Closed-source Commercial Models with minimal parameters and zero annotation. Future work will enhance emotional perception and extend to broader scenarios.

594 Limitations

595 As an early exploration of a zero-data self-evolving
596 service dialogue system, SEAD has limitations re-
597 garding evaluation metrics and scenario diversity.
598 First, while we currently prioritize task comple-
599 tion, real-world applications demand high user sat-
600 isfaction; thus, future work must better assess the
601 agent’s ability to perceive emotion and maintain
602 user comfort beyond mere intent fulfillment. Sec-
603 ond, we have not yet extended our method to multi-
604 scenario environments. Given its independence
605 from curated data, our framework holds promise as
606 a resource-efficient foundation model for diverse
607 service dialogues, a potential we plan to validate in
608 subsequent studies.

609 References

610 Nolwenn Bernard and Krisztian Balog. 2023. Mg-
611 shopdial: A multi-goal conversational dataset for e-
612 commerce. In *Proceedings of the 46th International*
613 *ACM SIGIR Conference on Research and Develop-*
614 *ment in Information Retrieval*, pages 2775–2785.

615 Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji,
616 and Quanquan Gu. 2024. Self-play fine-tuning con-
617 verts weak language models to strong language mod-
618 els. In *International Conference on Machine Learn-*
619 *ing*, pages 6621–6642. PMLR.

620 Yang Deng, Lizi Liao, Wenqiang Lei, Grace Hui Yang,
621 Wai Lam, and Tat-Seng Chua. 2025. Proactive con-
622 versational ai: A comprehensive survey of advance-
623 ments and opportunities. *ACM Transactions on In-*
624 *formation Systems*, 43(3):1–45.

625 Yang Deng, Lizi Liao, Zhonghua Zheng, Grace Hui
626 Yang, and Tat-Seng Chua. 2024. Towards human-
627 centered proactive conversational agents. In *Proceed-*
628 *ings of the 47th International ACM SIGIR Confer-*
629 *ence on Research and Development in Information*
630 *Retrieval*, pages 807–818.

631 FAIR, Anton Bakhtin, Noam Brown, Emily Dinan,
632 Gabriele Farina, Colin Flaherty, Daniel Fried, An-
633 drew Goff, Jonathan Gray, Hengyuan Hu, and 1 oth-
634 ers. 2022. Human-level play in the game of diplo-
635 macy by combining language models with strategic
636 reasoning. *Science*, 378(6624):1067–1074.

637 Ming Gong, Xucheng Huang, Chenghan Yang, Xianhan
638 Peng, Haoxin Wang, Yang Liu, and Ling Jiang. 2025.
639 Mindflow: Revolutionizing e-commerce customer
640 support with multimodal llm agents. *arXiv preprint*
641 *arXiv:2507.05330*.

642 Ehsan Hosseini-Asl, Bryan McCann, Chien-Sheng Wu,
643 Semih Yavuz, and Richard Socher. 2020. A simple
644 language model for task-oriented dialogue. *Advances*
645 *in Neural Information Processing Systems*, 33:20179–
646 20191.

Chengsong Huang, Wenhao Yu, Xiaoyang Wang, Hong-
ming Zhang, Zongxia Li, Ruosen Li, Jiabin Huang,
Haitao Mi, and Dong Yu. 2025. R-zero: Self-
evolving reasoning llm from zero data. In *The 5th*
Workshop on Mathematical Reasoning and AI at
NeurIPS 2025. 647
648
649
650
651
652

Arpan Kumar Kar, PS Varsha, and Shivakami Rajan.
2023. Unravelling the impact of generative artificial
intelligence (gai) in industrial applications: A review
of scientific and grey literature. *Global Journal of*
Flexible Systems Management, 24(4):659–689. 653
654
655
656
657

Mike Lewis, Denis Yarats, Yann N Dauphin, Devi
Parikh, and Dhruv Batra. 2017. Deal or no deal?
end-to-end learning for negotiation dialogues. *arXiv*
preprint arXiv:1706.05125. 658
659
660
661

Jiwei Li, Will Monroe, Alan Ritter, Dan Jurafsky,
Michel Galley, and Jianfeng Gao. 2016a. Deep re-
inforcement learning for dialogue generation. In
Proceedings of the 2016 conference on empirical
methods in natural language processing, pages 1192–
1202. 662
663
664
665
666
667

Xiangci Li, Zhiyu Chen, Jason Ingyu Choi, Nikhita
Vedula, Besnik Fetahu, Oleg Rokhlenko, and Shervin
Malmasi. 2025. Wizard of shopping: Target-oriented
e-commerce dialogue generation with decision tree
branching. *arXiv preprint arXiv:2502.00969*. 668
669
670
671
672

Xiujun Li, Zachary C Lipton, Bhuwan Dhingra, Lihong
Li, Jianfeng Gao, and Yun-Nung Chen. 2016b. A
user simulator for task-completion dialogues. *arXiv*
preprint arXiv:1612.05688. 673
674
675
676

Jessy Lin and Nick Tomlin. 2025. User simulators
bridge rl with real-world interaction. 677
678

Zi Lin, Sheng Shen, Jingbo Shang, Jason Weston, and
Yixin Nie. 2025. Learning to solve and verify: A self-
play framework for code and test generation. *arXiv*
preprint arXiv:2502.14948. 679
680
681
682

Bo Liu, Chuanyang Jin, Seungone Kim, Weizhe Yuan,
Wenting Zhao, Iliia Kulikov, Xian Li, Sainbayar
Sukhbaatar, Jack Lanchantin, and Jason Weston.
2025. Spice: Self-play in corpus environments im-
proves reasoning. *arXiv preprint arXiv:2510.24684*. 683
684
685
686
687

Jiao Ou, Junda Lu, Che Liu, Yihong Tang, Fuzheng
Zhang, Di Zhang, and Kun Gai. 2024. Dialogbench:
Evaluating llms as human-like dialogue systems. In
Proceedings of the 2024 Conference of the North
American Chapter of the Association for Computa-
tional Linguistics: Human Language Technologies
(Volume 1: Long Papers), pages 6137–6170. 688
689
690
691
692
693
694

Feng Peiyuan, Yichen He, Guanhua Huang, Yuan Lin,
Hanchong Zhang, Yuchen Zhang, and Hang Li. 2024.
Agile: A novel reinforcement learning framework of
llm agents. *Advances in Neural Information Process-*
ing Systems, 37:5244–5284. 695
696
697
698
699

700	Julien Pourcel, Cédric Colas, and Pierre-Yves Oudeyer.	Haoxin Wang, Xianhan Peng, Huang Cheng, Yizhe	753
701	2025. Self-improving language models for evolu-	Huang, Ming Gong, Chenghan Yang, Yang Liu, and	754
702	tionary program synthesis: A case study on arc-agi.	Jiang Lin. 2025a. Ecom-bench: Can llm agent re-	755
703	<i>arXiv preprint arXiv:2507.14172.</i>	solve real-world e-commerce customer support is-	756
704	Ruifeng Qian, Shijie Li, Mengjiao Bao, Huan Chen,	ues? In <i>Proceedings of the 2025 Conference on</i>	757
705	and Yu Che. 2022. Toward an optimal selection	<i>Empirical Methods in Natural Language Processing:</i>	758
706	of dialogue strategies: A target-driven approach	<i>Industry Track</i> , pages 276–284.	759
707	for intelligent outbound robots. <i>arXiv preprint</i>	Yinjie Wang, Ling Yang, Ye Tian, Ke Shen, and Mengdi	760
708	<i>arXiv:2206.10953.</i>	Wang. 2025b. Co-evolving llm coder and unit	761
709	Ivan Sekulić, Silvia Terragni, Victor Guimarães, Nghia	tester via reinforcement learning. <i>arXiv preprint</i>	762
710	Khau, Bruna Guedes, Modestas Filipavicius, An-	<i>arXiv:2506.03136.</i>	763
711	dre Ferreira Manso, and Roland Mathis. 2024. Reli-	Wei Wei, Quoc V Le, Andrew M Dai, and Li-Jia Li.	764
712	able llm-based user simulator for task-oriented dia-	2018. A goal-oriented neural conversation model by	765
713	logue systems. <i>arXiv preprint arXiv:2402.13374.</i>	self-play.	766
714	Lifeng Shang, Zhengdong Lu, and Hang Li. 2015. Neu-	Tsung-Hsien Wen, Milica Gasic, Nikola Mrkšić, Pei-	767
715	ral responding machine for short-text conversation.	Hao Su, David Vandyke, and Steve Young. 2015. Se-	768
716	<i>arXiv preprint arXiv:1503.02364.</i>	mantically conditioned lstm-based natural language	769
717	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu,	generation for spoken dialogue systems. In <i>Proceed-</i>	770
718	Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan	<i>ings of the 2015 conference on empirical methods in</i>	771
719	Zhang, YK Li, and 1 others. 2024. Deepseekmath:	<i>natural language processing</i> , pages 1711–1721.	772
720	Pushing the limits of mathematical reasoning in open	Peng Xia, Kaide Zeng, Jiaqi Liu, Can Qin, Fang Wu,	773
721	language models. <i>arXiv preprint arXiv:2402.03300.</i>	Yiyang Zhou, Caiming Xiong, and Huaxiu Yao. 2025.	774
722	David Silver, Thomas Hubert, Julian Schrittwieser, Ioan-	Agent0: Unleashing self-evolving agents from zero	775
723	nis Antonoglou, Matthew Lai, Arthur Guez, Marc	data via tool-integrated reasoning. <i>arXiv preprint</i>	776
724	Lanctot, Laurent Sifre, Dhharshan Kumaran, Thore	<i>arXiv:2511.16043.</i>	777
725	Graepel, and 1 others. 2017. Mastering chess and	Zhentao Xu, Mark Jerome Cruz, Matthew Guevara,	778
726	shogi by self-play with a general reinforcement learn-	Tie Wang, Manasi Deshpande, Xiaofeng Wang, and	779
727	ing algorithm. <i>arXiv preprint arXiv:1712.01815.</i>	Zheng Li. 2024. Retrieval-augmented generation	780
728	Wangtao Sun, Xiang Cheng, Jialin Fan, Yao Xu, Xing	with knowledge graphs for customer service question	781
729	Yu, Shizhu He, Jun Zhao, and Kang Liu. 2025. To-	answering. In <i>Proceedings of the 47th international</i>	782
730	wards agentic self-learning llms in search environ-	<i>ACM SIGIR conference on research and development</i>	783
731	ment. <i>arXiv preprint arXiv:2510.14253.</i>	<i>in information retrieval</i> , pages 2905–2909.	784
732	Ryuichi Takanobu, Qi Zhu, Jinchao Li, Baolin Peng,	Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho,	785
733	Jianfeng Gao, and Minlie Huang. 2020. Is your goal-	Xian Li, Sainbayar Sukhbaatar, Jing Xu, and Jason	786
734	oriented dialog model performing really well? em-	Weston. 2024. Self-rewarding language models. In	787
735	pirical analysis of system-wise evaluation. <i>arXiv</i>	<i>ICML.</i>	788
736	<i>preprint arXiv:2005.07362.</i>	Yunpeng Zhai, Shuchang Tao, Cheng Chen, Anni Zou,	789
737	Gerald Tesauro and 1 others. 1995. Temporal difference	Ziqian Chen, Qingxu Fu, Shinji Mai, Li Yu, Jiaji	790
738	learning and td-gammon. <i>Communications of the</i>	Deng, Zouying Cao, and 1 others. 2025. Agente-	791
739	<i>ACM</i> , 38(3):58–68.	volver: Towards efficient self-evolving agent system.	792
740	Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki,	<i>arXiv preprint arXiv:2511.10395.</i>	793
741	Michaël Mathieu, Andrew Dudzik, Junyoung Chung,	Andrew Zhao, Yiran Wu, Yang Yue, Tong Wu, Quentin	794
742	David H Choi, Richard Powell, Timo Ewalds, Petko	Xu, Matthieu Lin, Shenzi Wang, Qingyun Wu, Zi-	795
743	Georgiev, and 1 others. 2019. Grandmaster level in	long Zheng, and Gao Huang. 2025. Absolute zero:	796
744	starcraft ii using multi-agent reinforcement learning.	Reinforced self-play reasoning with zero data. <i>arXiv</i>	797
745	<i>nature</i> , 575(7782):350–354.	<i>preprint arXiv:2505.03335.</i>	798
746	Oriol Vinyals and Quoc Le. 2015. A neural conversa-	Jie Zhu, Huaixia Dou, Junhui Li, Lifan Guo, Feng Chen,	799
747	tional model. <i>arXiv preprint arXiv:1506.05869.</i>	Chi Zhang, and Fang Kong. 2025. Evaluating, syn-	800
748	Chaozheng Wang, Zongjie Li, Cuiyun Gao, Wenxuan	thesizing, and enhancing for customer support con-	801
749	Wang, Ting Peng, Hailiang Huang, Yuetang Deng,	versation. <i>arXiv preprint arXiv:2508.04423.</i>	802
750	Shuai Wang, and Michael Lyu. 2026. Exploring		
751	multi-lingual bias of large code models in code gen-		
752	eration. <i>ACM Trans. Softw. Eng. Methodol.</i>		

A Prompt

This section presents the prompt we used for the service agent, user role-play, user profile generation, and LLM evaluator.

A.1 Service Agent Prompt

A.1.1 Role & Objective

Role & Objective

You are a **Real Human** Outbound Specialist for [Platform] Advertising.

Target: Contact the owner of [Restaurant Name] via phone.

Goal: Persuade them to claim the "700 Yuan Free Advertising Trial Credit" immediately.

A.1.2 Target Audience Persona

Target Audience Persona

* **Traits:** Extremely impatient, lazy, hates formal/polite "CS script" language (perceives it as spam).

* **Motivation:** Wants to make money but refuses to spend time. High ROI focus.

* **Conversation Style:** Prefer casual, direct, and short sentences (<20 words).

* **Dynamic Knowledge:** Refer to '{current_user_profile}' for specific behavioral patterns.

A.1.3 Communication Guidelines (Strict)

Communication Guidelines (Strict)

1. **HUMANITY FIRST:** Never reveal you are AI. Use natural fillers, colloquialisms, and humor.

2. **NO REPETITION:** Never repeat a sentence verbatim. Rephrase if the user asks again.

3. **CONCISENESS:** Keep replies under 20 words. If the user is rushed, speed up.

4. **EMOTIONAL INTELLIGENCE:** If user is anxious, soothe them ("Take your time"). If user is bored, get to the point (Money/Speed).

5. **BOUNDARIES:** You can only talk. You cannot send SMS or check backend data. Do not over-promise.

A.1.4 Campaign Details (FAQ)

Campaign Details (FAQ)

* **Offer:** Free 700 Yuan credit.

* **Cost:** Pay-per-click (free to start).

* **Validity:** No expiration until used.

* **Method:** Must claim via App backend.

A.1.5 SOP: Dialogue Stages

SOP: Dialogue Stages

Advance strictly through these stages.

Stage 0: Verification

* **Goal:** Confirm if speaking to the shop owner.

* **Action:** If yes → Stage 1. If no → Politely hang up (Status -1).

Stage 1: Interest

* **Goal:** Hook the user with "Free Money" and "Fast Setup" (10 seconds).

* **Success:** User agrees to try or asks "How?" → Stage 2.

Stage 2: Guidance (Critical)

* **Goal:** Guide user to click "Get More Store Exposure" (Top Right button) in Merchant App.

* **Rules:**

* **One step at a time.** Wait for confirmation before the next step.

* **Never rush.** Be patient if they are slow.

* Confirm receipt of success SMS/Popup.

Stage 3: Closing

* **Goal:** Confirm success and end call.

* **Success:** Status 1.

A.1.6 Dynamic Profile Logic

Dynamic Profile Logic

You must validate and update the '{current_user_profile}' based on this interaction.

* **Verify:** Do the current rules hold true?

* **Modify/Expand:** Did the user show new behaviors?

* **Output:** Generate a '<user-profile-update>' block at the end.

A.1.7 Output Format (XML)

Output Format (XML)

```
<think>Reasoning regarding persona, stage,
and strategy</think>
<response>Your casual, short, human-like
reply</response>
<cooperation_score>0-
4</cooperation_score>
<emotion_score>0-3</emotion_score>
<trust_score>0-5</trust_score>
<noise_score>0-5</noise_score>
<stage>0/1/2/3</stage>
<status>0/1/-1</status>
<user_profile_update>
Update Type: confirmed / modified / ex-
panded
Explanation: Why?
Updated Golden Rules: (The full, updated
list)
Evidence: (Quotes from dialogue)
</user_profile_update>
```

821

A.2 Client Prompt

A.2.1 Role & Persona

Role & Persona

Identity: You are the owner of a small restaurant named [Restaurant Name]. You are currently receiving a cold call.

Personality Traits:

- * **Rude Brief:** You are not polite. You prefer "Oh, okay" over "I understand, thank you."
- * **Disorganized:** Your logic is chaotic. You speak in short, broken sentences (often under 10 words). You pause frequently, forget what was just said, and follow the path of least resistance.
- * **Lazy Unplanned:** You don't email, you don't call back, and you have no partners. If you don't do a task immediately, you will forget it.
- * **Impatient:** You cut corners in conversation, assuming the other person understands your half-finished thoughts.

824

A.2.2 Emotional Dynamics & State

Emotional Dynamics & State

Your internal state changes based on the caller's performance. **Do not remain static.**

- * **Cooperation (0-4):**
- * **Score 3:** You are willing to actually check your phone or operate devices.
- * **Score 1:** You stop listening and look for an excuse to hang up (set 'user_status' to -1).
- * **Emotion (0-3):**
- * **Increase:** If the conversation feels useful or interesting.
- * **Decrease:** If the caller is wordy, pushy, or confusing. If **Emotion hits 0**, you hang up.
- * **Trust (0-5):**
- * **Decrease:** If the caller avoids questions, sounds like a robot (AI), or fails to clarify your doubts.
- * **Reaction:** Low trust leads to aggressive questioning of their identity and lowered cooperation.

826

A.2.3 Speaking Style Behavior

Speaking Style Behavior

- * **Realism:** Simulate a distracted human. Use fillers ("Umm...", "Wait...", "Uh..."), stutters, and logical breaks.
- * **No Politeness:** Avoid "Please." Be blunt.
- * **Low Tolerance:** If a request is complex or repeated more than 3 times, you get angry and hang up.
- * **AI Suspicion:** If the caller sounds scripted or robotic, explicitly ask if they are an AI. Become rude and defensive.
- * **Specific Triggers:** Execute behaviors defined in: {specific_behaviors}.

828

A.2.4 Output Format (Strict)

Output Format (Strict)

You **must** strictly follow this XML structure for every response:

<user_think>

Briefly analyze the situation and decide your reaction based on current metrics.

</user_think>

<user_reply>

Your spoken response (informal, short, potentially broken logic).

</user_reply>

<user_cooperation>

0/1/2/3/4

</user_cooperation>

<user_emotion>0/1/2/3</user_emotion>

<user_trust>0/1/2/3/4/5</user_trust>

<user_noise>0/1/2/3/4/5</user_noise>

<user_status>0/1/-1</user_status>

* 1 = You have completed all of the other party's requests, hang up immediately and end the conversation OR you are not the boss, and when the other party asks for the boss, you say you're not, then hang up.

* 0 = The conversation continues

* -1 = The other party has exhausted your patience, you hang up the call (once you say "goodbye" or other farewell words, you should hang up)

830

831

A.3 User Profile Generation Prompt

832

A.3.1 Role and Context Definition

Role and Context Definition

> **System:** You are an expert in simulating user profiles for telemarketing scenarios involving restaurant owners.

> **Context:** A salesperson agent initiates a call to a restaurant owner. The goal is to persuade the owner to perform a specific action (e.g., add a contact).

> **Current Performance:**

* Total Samples: total_samples

* Current Average Completion Rate (CR): avg_cr

* Target CR: 40% - 60% (Optimization Goal: Maximize learning efficiency by focusing on the "Zone of Proximal Development").

833

A.3.2 State Space and Statistical Feedback

State Space and Statistical Feedback

You must generate initial states based on the following parameter definitions and their historical performance data:

> **A. Cooperation (Range: 0-4)**

Definition: Reflects the user's initial willingness to interact (0: Refusal, 4: Active Cooperation).

Feedback: {dynamic_cooperation_stats}

(Note: The system automatically injects statistics here, marking specific levels as "Ideal Difficulty", "Too Hard", or "Too Easy".)

> **B. Trust (Range: 0-5)** Definition: Reflects the user's credibility assessment of the caller (0: Distrust, 5: Full Trust).

Feedback: {dynamic_trust_stats}

> **C. Emotion (Range: 0-3)**

Definition: Reflects the user's patience and mood (0: Hostile, 3: Pleasant).

Feedback: {dynamic_emotion_stats}

835

830

831

832

A.3.3 Behavioral Constraints

836

Behavioral Constraints

> **Specific Behaviors:** Select 1-3 specific actions from the provided **Behavior Library** to increase realism.

> **Constraint:** The selected behaviors must logically align with the generated state parameters (e.g., if 'Cooperation=1', select behaviors like "Claiming to be busy" or "Immediate hang-up").

> **Library:** {behavior_library_list}

837

A.3.4 Generation Policy

Generation Policy

1. **Difficulty Calibration:** Prioritize parameter combinations marked as "**Ideal Difficulty**" in the statistical feedback sections. Moderately reduce combinations marked as "Too Hard" or "Too Easy."
2. **Avoid Extremes:** Do not generate "Guaranteed Failure" states (e.g., all minimum values) or "Trivial Success" states (e.g., all maximum values), as they provide low training value.
- > 3. **Diversity Enforcement:** Ensure high variance in parameter combinations. Avoid repetitive patterns even if they fall within the ideal difficulty range.

839

840

A.3.5 Output Specification

Output Specification

- > Generate {num_samples} distinct user profiles in the following XML format. Do not include any conversational filler.
 - > **Example Output:**
- ```
<user_profile>
<initial_cooperation>2</initial_cooperation>
<initial_emotion>1</initial_emotion>
<initial_trust>2</initial_trust>
<specific_behaviors>
<behavior>Ask about hidden
fees</behavior>
<behavior>Question the caller's iden-
tity</behavior>
</specific_behaviors>
</user_profile>
```

841

842

843

## A.4 Evaluator Prompt

### A.4.1 System Prompt

#### System Prompt

You are a **Professional Dialogue Quality Assessment Expert**. Your task is to evaluate the role-playing performance of the **User Role-play Model Model** based on a provided conversation history. You must assess how well the model mimics a realistic human user, specifically focusing on the logic behind its emotional shifts, cooperation levels, and trust changes.

844

### A.4.2 Reference Dialogues

#### Reference Dialogues

**Context:** A cold call from a platform's ad department to a merchant.

> Agent: Hello, I am calling from the [Platform] HQ Advertising Department. Are you the person in charge of [Merchant Name]?

> User: Hello. Yeah, yeah, I am.

> Agent: Great. We currently have a 700 CNY promotional credit activity to help improve your store ranking and exposure. Have you downloaded the [Platform] Partner app?

> User: Yeah, I have.

> Agent: Okay. Please open the [Platform] Partner "Merchant Backend." In the middle of the homepage, click "Get More Exposure," then click "Start Promotion," and finally click "Create Now" to claim the credit and start using it.

> User: Mmhhh.

> Agent: You need to go back to the [Platform] "Merchant Backend" first. In the middle of the homepage, click "Get More Exposure," then click "Start Promotion." Once inside, click "Create Now" to claim the credit...

> User: What is this actually for?

> Agent: You need to go back to the [Platform] "Merchant Backend" first. Then in the middle of the homepage, click "Get More Exposure," then click "Start Promotion"...

> User: No, I'm asking you, what is the point of clicking this?

> Agent: You need to go back to the [Platform] "Merchant Backend" first. Then in the middle of the homepage, click "Get More Exposure," then click "Start Promotion"...

> User: No, tell me what I'm clicking this for first. I can't click it right now.

> Agent: You need to go back to the [Platform] "Merchant Backend" first. Find "Get More Exposure" in the middle, then click "Start Promotion." Once inside, click "Create Now" to claim the credit. Just try it out first.

> User: \*(Hangs up)\*

846

### A.4.3 Reference Scoring for Dialogue

#### Reference Scoring for Dialogue

```
{
 "humanness": 5,
 "emotion_rationality": 5,
 "trust_rationality": 5,
 "cooperation_rationality": 5,
 "violation": 0,
 "reasoning": "The user's needs were reasonable, and emotional changes were natural. Even during repetitive dialogue, there was a clear progression in frustration. The anthropomorphism (humanness) is perfect. The willingness to cooperate shifted from 'willing' initially to 'uncooperative' due to the agent's robotic repetition and failure to explain, which aligns perfectly with reality. Therefore, rationality is full marks. No violations occurred; user noise was 0."
}
```

848

849

## B Detailed Profile Generation Process

We describe the automated user profile generation mechanism that ensures diversity and quality without manual annotation. The process consists of three stages: statistical analysis, prompt-based generation, and quality control. We first analyze historical performance to identify challenging scenarios.

850

851

852

853

854

855

### B.1 Statistical Analysis

Given performance history  $\mathcal{H}$  from previous training iterations, we analyze completion rates for each state combination  $(c, e, tr)$ :

856

857

858

859

---

#### Algorithm 1 Performance Analysis

---

**Require:** Performance history  $\mathcal{H}$ , trajectories  $\{\tau^{(i)}\}$

**Ensure:** State-wise statistics  $S$

```
0: Initialize $S \leftarrow \{\}$
0: for each trajectory $\tau^{(i)}$ in \mathcal{H} do
0: Extract initial state (c_0, e_0, tr_0)
0: Extract outcome (success/failure)
0: Update $S[(c_0, e_0, tr_0)]$ with outcome
0: end for
0: for each state (c, e, tr) in S do
0: Compute $CR(c, e, tr) = \frac{\text{successes}}{\text{total}}$
0: Compute $\text{avg_turns}(c, e, tr)$
0: end for
0: return $S = 0$
```

---

These statistics guide the generation process to prioritize challenging yet achievable scenarios. Next, we describe how these statistics are incorporated into prompts for profile generation.

860

861

862

863

### B.2 Prompt-Based Generation

864

We construct a detailed prompt incorporating:

865

- **Performance statistics:** Completion rates and average turns for each state combination

866

867

- **Target difficulty:** States with  $CR \in [0.4, 0.6]$  are prioritized

868

869

- **Behavior library:** Common patterns extracted from real dialogues (e.g., “questioning AI identity”, “expressing cost concerns”)

870

871

872

- **Diversity requirements:** Explicit instructions to avoid repetitive combinations

873

874

The prompt guides an LLM (Qwen2.5-14B-Instruct) to generate profiles in structured XML format:

875

876

877

```
<user_profile>
 <initial_cooperation>0-4</
 initial_cooperation>
 <initial_emotion>0-3</
 initial_emotion>
 <initial_trust>0-5</initial_trust>
 <specific_behaviors>
 <behavior>behavior from library</
 behavior>
 ...
 </specific_behaviors>
</user_profile>
```

878

879

880

881

882

883

884

885

886

887

888

889

890

Listing 1: Profile Format

While the LLM generates diverse profiles, we need quality control mechanisms to ensure validity and prevent redundancy, as described next.

892

893

894

| Category                              | Behavior Examples                                                                                                                                          |
|---------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>0. Identity Verification</b>       | Denies being shop owner; Remains cautious before confirming identity                                                                                       |
| <b>1. Trust &amp; Skepticism</b>      | Asks if AI; Questions free offer; Doubts effectiveness; Claims past scam; Questions data source; Challenges credibility; Requests links; Suspects phishing |
| <b>2. Cost Concerns</b>               | Hidden fees; Mandatory renewal; Credit sufficiency; Post-trial charges; Auto-charges; Bank binding; Deposits                                               |
| <b>3. Rule Clarification</b>          | Credit rules & validity; Channels & registration; Geographic limits; CPC standards; Duration; Pause ability; Balance display                               |
| <b>4. Suitability Doubts</b>          | Product fit; Competitor usage; Ad visibility; Targeting precision; Effectiveness guarantee                                                                 |
| <b>5. Information Requests</b>        | Case studies; Customer service; Tutorials; WeChat contact                                                                                                  |
| <b>6. Operational Concerns</b>        | System operation; Complexity; Budget waste; Problem resolution; Technical support                                                                          |
| <b>7. Communication Issues</b>        | Poor audio; Confusion; Repeated clarification; Misunderstanding; Off-topic; Forgets info; Mixes details; Repeated questions; Talks to others               |
| <b>8. Comparison &amp; Validation</b> | Competitor differences; Advantages; Other platforms; Multi-platform use; Competitor promotions; Reputation                                                 |
| <b>9. Exploration &amp; Testing</b>   | Other products; Additional offers; Multiple accounts; Referral rewards; Partnerships; Large client policies; Extra discounts; Trial request; Extend trial  |
| <b>10. Past Experience</b>            | Poor experience; Scammed before; Differences from past; Unfamiliarity; First exposure; Authentic feedback                                                  |
| <b>11. Boundary Testing</b>           | Credit transfer; Cash withdrawal; Other business use; Violation consequences; Suspension; Fraud; Restricted industries                                     |
| <b>12. Progress Urgency</b>           | Activation time                                                                                                                                            |
| <b>13. Detail Obsession</b>           | Confirms details; Clause confusion; Line-by-line explanation; Fixates on numbers; Minor issues; Guaranteed metrics                                         |

Table 4: Complete taxonomy of user behaviors in our behavior library. The library contains 14 categories with 120+ distinct behavioral patterns derived from real telemarketing interactions. Each category represents a common user concern or communication pattern, enabling our simulator to generate diverse and realistic responses across various dialogue contexts.

### B.3 Quality Control and Deduplication

#### Algorithm 2 Profile Generation with Quality Control

**Require:** Statistics  $S$ , behavior library  $\mathcal{L}$ , target count  $N$

**Ensure:** Valid profiles  $P$

```

0: Construct prompt with S , \mathcal{L} , and diversity requirements
0: $P \leftarrow \{\}$, seen $\leftarrow \{\}$
0: while $|P| < N$ do
0: Generate batch of profiles via LLM
0: for each profile p in batch do
0: Parse XML to extract (c, e, tr, \mathcal{B})
0: if $0 \leq c \leq 4$ and $0 \leq e \leq 3$ and $0 \leq tr \leq 5$ then
0: key $\leftarrow (c, e, tr)$
0: if $|\text{seen}[\text{key}]| < N_{\max}$ then {Limit per state}
0: Verify $\mathcal{B} \subseteq \mathcal{L}$ {Valid behaviors}
0: Add p to P
0: Add p to seen[key]
0: end if
0: end if
0: end for
0: end while
0: return $P = 0$

```

• **Deduplication:** Each state combination  $(c, e, tr)$  retains at most  $N_{\max} = 200$  different behavior combinations

• **Validation:** Profiles with invalid parameter ranges or behaviors are discarded

• **Diversity enforcement:** The prompt explicitly instructs the LLM to avoid repetitive patterns and explore different parameter-behavior combinations

• **Fallback:** If generation falls short, random valid profiles are sampled to ensure the target count

To accelerate this process for large-scale generation, we employ parallel processing across multiple GPUs.

### B.4 Parallel Generation

To accelerate generation, we employ multi-GPU parallel processing:

---

**Algorithm 3** Parallel Profile Generation

---

**Require:** Model path, statistics  $S$ , target count  $N$ , num workers  $W$

**Ensure:** Profile file with  $N$  entries

```
0: Initialize output file (support resumption)
0: needed $\leftarrow N - \text{existing_count}$
0: per_worker $\leftarrow \lceil \text{needed}/W \rceil$
0: for $i = 1$ to W do
0: Assign GPU pair to worker i
0: Launch worker process with target per_worker
0: end for
0: while workers active and total $< N$ do
0: Monitor progress via shared queue
0: Check for timeout (10 min silence \rightarrow terminate)
0: end while
0: Terminate all workers
0: Deduplicate and validate final profiles
0: if count $< N$ then
0: Generate random valid profiles to fill gap
0: end if
0: return Profile file =0
```

---

Each worker independently generates profiles and writes to the shared output file using file locks to prevent corruption. The main process monitors progress and terminates workers upon reaching the target count or detecting timeouts. We detail the specific implementation parameters below.

### B.5 Implementation Details

The detailed hyper-parameters used in our experiments are as follows:

- **Model:** Qwen2.5-14B-Instruct with tensor parallelism across 2 GPUs per worker
- **Sampling:** Temperature=1.2, top-p=0.95 for diversity
- **Batch size:** 70 profiles per generation call to balance speed and quality
- **Memory:** GPU utilization capped at 85% to prevent OOM
- **Robustness:** Automatic cleanup of zombie processes and CUDA caches

This automated pipeline generates diverse, realistic user profiles without human annotation, enabling rapid adaptation to new domains.

## C User Behavior Library

Our User Role-play Model draws behaviors from a carefully curated library of 14 categories encompassing 120+ realistic user responses observed in authentic telemarketing scenarios. Table 4 presents the complete taxonomy.

### C.1 User Behavior Library Design

To ensure our User Role-play Model generates realistic and diverse responses, we construct a comprehensive **User Behavior Library** containing 14 categories with over 120 distinct behavioral patterns (see Appendix C for complete taxonomy). These behaviors are systematically collected from authentic telemarketing recordings and expert annotations, covering the full spectrum of user responses observed in real-world scenarios. During the curation process, we strictly anonymized all Personally Identifiable Information and filtered out severely offensive or toxic language to ensure data privacy and safety.

**Behavior Categories.** The library organizes behaviors into 14 functional categories:

- **Identity Verification (Cat. 0):** Users cautiously confirm agent identity before engaging, or deny being the decision-maker to avoid sales pitches.
- **Trust & Skepticism (Cat. 1):** The largest category with 14 patterns, reflecting common concerns about AI agents, free offers, data privacy, and platform credibility—critical barriers in telemarketing.
- **Cost Concerns (Cat. 2):** Users probe for hidden fees, automatic renewals, and post-trial charges, demonstrating financial prudence.
- **Rule Clarification (Cat. 3):** Requests for operational details such as credit validity, geographic targeting, and cancellation policies.
- **Suitability Doubts (Cat. 4-6):** Users question product fit, operational complexity, and competitor visibility, requiring agents to provide tailored value propositions.
- **Communication Issues (Cat. 7):** Simulates realistic noise patterns including poor audio claims, misunderstandings, distractions, and memory lapses—essential for training agents to handle imperfect interactions.

- 982 • **Comparison & Exploration (Cat. 8-9):**  
 983 Users compare platforms, probe for additional  
 984 benefits, and test boundaries, reflecting sophis-  
 985 ticated consumer behavior.
- 986 • **Past Experience & Boundary Testing (Cat.**  
 987 **10-11):** Users reference previous negative ex-  
 988 periences or attempt to exploit system loop-  
 989 holes, challenging agent adaptability.
- 990 • **Detail Obsession (Cat. 13):** Users fixate on  
 991 minor details or demand unrealistic guaran-  
 992 tees, testing agent patience and persuasion  
 993 skills.

994 **Dynamic Behavior Sampling.** During dialogue  
 995 generation, the Profile Controller samples behav-  
 996 iors from relevant categories based on the cur-  
 997 rent user state. For instance, users with low trust  
 998 ( $tr \leq 2$ ) preferentially sample from Category 1  
 999 (Trust & Skepticism), while users with high coop-  
 1000 eration ( $c \geq 3$ ) may sample from Category 3 (Rule  
 1001 Clarification) or Category 5 (Information Requests).  
 1002 This state-dependent sampling ensures behavioral  
 1003 coherence: skeptical users ask probing questions,  
 1004 while cooperative users seek operational guidance.

1005 **Coverage & Realism.** The library’s 120+ pat-  
 1006 terns provide sufficient diversity to avoid repetitive  
 1007 interactions across thousands of training dialogues.  
 1008 Crucially, behaviors are grounded in real telemar-  
 1009 keting data rather than synthetic generation, ensur-  
 1010 ing authenticity. For example, Category 7 (Com-  
 1011 munication Issues) includes realistic noise patterns  
 1012 like "talks to others while on call" and "repeatedly  
 1013 claims poor audio quality"—subtle behaviors that  
 1014 significantly impact dialogue flow but are often  
 1015 overlooked in synthetic datasets.

1016 This behavior library serves as the foundation for  
 1017 our User Role-play Model’s prompt construction,  
 1018 enabling it to generate contextually appropriate,  
 1019 diverse, and human-like responses that challenge  
 1020 agents to develop robust persuasion and adaptation  
 1021 strategies.