

Efficient Query-Focused User Simulation via Interactive Persona Completion

Weiwen Su^{1,3} Naoki Yoshinaga^{2,3} Masashi Toyoda^{2,3}

¹The University of Tokyo ²Institute of Industrial Science, The University of Tokyo

³Institute for Digital Observatory, the University of Tokyo

{su-w, ynaga, toyoda}@tkl.iis.u-tokyo.ac.jp

Abstract

Large language models (LLMs) can simulate human responses using persona information, yet such rich information, assumed by existing work, is often unavailable in practice. Recent studies show that LLMs can identify query-relevant persona dimensions (e.g., *whether a user is shy?*), enabling user simulation under a cold-start setting with limited basic persona. We study query-focused user simulation in this setting and propose a persona acquisition method that progressively predicts and acquires relevant persona dimensions. Experiments on the PICQ dialogue dataset show that our approach achieves comparable performance to query-agnostic and one-shot persona collection while requiring less persona information.

1 Introduction

Large language models (LLMs) have made it feasible to simulate human behavior [1, 2, 3] based on textual descriptions of individuals (i.e., persona information). Such simulation capabilities enable applications such as opinion elicitation [4] and realistic virtual character creation [5]. By conditioning generation on persona information, LLM-based simulators can better capture individual differences in beliefs, values, and preferences. Despite this promise, existing approaches assume access to rich persona information, which limits their applicability in real-world settings.

Prior user simulation studies rely on extensive persona information from encyclopedic biographies [6] or interviews [2, 7], enabling faithful simulation of public figures. In contrast, real-world scenarios involve ordinary people whose personas are only partially observable, making comprehensive persona collection and simulation costly or infeasible. To address this issue, we have shown that query-



Figure 1 Query-focused individual simulation in a cold-start scenario with only a basic persona.

relevant but missing persona dimensions can be identified, where a persona dimension denotes a piece of persona whose value remains uncertain (e.g., *whether a user is price-sensitive*) [8], enabling user simulation in a cold-start setting by requesting only necessary persona values.

In this study, we define query-focused user simulation in a cold-start setting (Figure 1), where only basic persona information is available in advance, and additional query-relevant persona values are acquired on demand. We assume that the user requesting a simulation can provide persona information about the target when requested. To reduce persona acquisition cost, we propose a closed-loop method that predicts relevant persona dimensions and acquires their values. As the relevant persona dimensions are often interdependent, instead of collecting them at once, our approach progressively identifies and acquires the most relevant one at each step, conditioning on previously verified persona information to avoid redundant requests.

We evaluate our method on the personal queries in the PICQ dialogue dataset [9, 8] across multiple simulator LLM [10, 11] and baselines, including context-only simulation, query-agnostic, and one-shot persona collection. Experimental results show that our progressive persona completion achieves comparable simulation quality while substantially reducing persona acquisition cost.

2 Query-focused User Simulation

We study query-focused simulation in a cold-start setting, where a target individual’s response is simulated without assuming a rich persona in advance. Only a basic persona is initially available, and missing query-relevant persona information must be acquired for faithful simulation.

The simulation involves two parties: a user A and a target individual B , where A asks an LLM to simulate B ’s response to a query (Figure 1). We assume that A is familiar with B and can provide persona information upon request, as such information is typically unavailable in public sources. This setting applies to decision support, opinion elicitation, and character-driven content creation.

A key challenge is that different queries depend on different persona dimensions, making comprehensive persona specification impractical. We define a persona dimension as an unknown persona attribute framed as a verifiable question (e.g., *whether s/he is shy*). Since acquiring persona information incurs cost, the goal is to elicit only the minimum sufficient persona required for each query.

Dataset Construction We use the PICQ dataset [8], derived from TVShowGuess scripts [9], which contains manually annotated, context-aware choice QA pairs with missing persona dimensions. In each instance, the character pair serves as the proxy user A and the simulation target B . Dialogue segments in which both characters appear are extracted, temporally ordered, and segmented at the scene level. For each query, only the preceding dialogue is provided as memory, reflecting the information available to A at that time. This setup enables controlled and reproducible evaluation of persona acquisition and simulation.

Settings Formally, the system takes as input a dialogue context c between A and B and a query q requiring B ’s response. During simulation, the system may request additional persona information from A (or an LLM-based proxy conditioned on dialogue history H). The output is a simulated response generated from the dialogue context, query, and acquired persona information.

Compared to the idealized task setting, we make two practical assumptions: (1) Persona information about B is provided via an LLM-based proxy for A rather than a real user, due to cost, scalability, and privacy constraints. (2) Predicting persona dimensions from scratch priori-

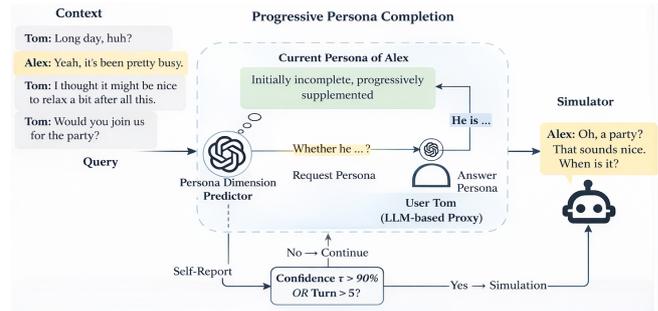


Figure 2 Overview of the progressive Persona completion for query-focused individual simulation.

tizes shallow attributes (e.g., gender), incurring inefficient warm-up cost; following prior work [8], we initialize prediction with a basic persona (gender, age, and participant relationship).

Evaluation focuses on simulation quality and persona acquisition cost, to achieve high fidelity using minimal persona information.

3 Progressive Persona Completion

To tackle the information dependency and acquisition difficulty of the one-shot method [8], we frame persona acquisition as a closed-loop, progressive process that interleaves persona prediction, acquisition, and stopping. Our method iteratively identifies and acquires the most relevant missing persona dimension, conditioned on the current known persona, as shown in Figure 2.

3.1 Persona Dimension Prediction

To address the information dependency problem, where knowledge of one verified dimension can obscure or necessitate others, we predict the most relevant missing persona dimension at each iteration using a prompt adapted from [8]. By acquiring persona dimensions progressively, each prediction is conditioned only on verified persona information, mitigating dependency effects. Persona dimensions deemed unanswerable by the proxy model are marked as unavailable and excluded from subsequent iterations to prevent redundant queries.

3.2 Persona Acquisition

Once a persona dimension is predicted, its value is acquired by querying user A or an LLM-based proxy. While prior work shows that LLMs can extract persona attributes of narrative characters from text [12], our setting differs in both granularity and context usage.

Using the dataset described in § 2, we employ GPT-5-mini as the proxy agent for persona acquisition. After acquiring the value $v^{(t)}$ for the predicted dimension $p^{(t)}$, the known persona is updated as $\mathcal{P}_{\text{known}}^{(t+1)} = \mathcal{P}_{\text{known}}^{(t)} \cup (p^{(t)}, v^{(t)})$. We analyze persona acquisition rate as a function of dialogue history length and observe rapid gains up to approximately 0.7 with around 5,000 space-segmented words, after which performance saturates. Accordingly, we cap the dialogue history provided to the proxy agent at 5,000 words, ordered from newest.

3.3 Confidence-Based Dynamic Stopping

A key challenge in progressive persona completion is determining when sufficient persona information has been acquired. We address this with a dynamic stopping mechanism based on the model’s self-reported confidence, combined with a fixed acquisition budget to prevent infinite loops when decisive persona dimensions are unavailable. Prior work shows that elicited confidence from LLMs can effectively estimate uncertainty [13, 14]. Building on these findings, we design a self-report confidence prompt (Appendix A) that asks the predictor to assess whether the current persona is sufficient for a reliable response and whether additional information would substantially change the outcome. By monitoring confidence across iterations, the system adaptively terminates persona acquisition, reducing unnecessary requests and user burden.

4 Experiments

We evaluate our method on the query-focused user simulation task from § 2 in dialogue-based settings. For each instance, the simulator is given a dialogue context C and a query q , and generates a simulated response of target B .

4.1 Settings

Models We choose GPT-4.1 and Qwen3-32B¹⁾ as simulators, representing strong closed-source and open-source models. We then evaluate multiple LLMs as missing persona dimension predictors with the confidence reporter being themselves, including GPT-4.1 and Qwen3-32B, which achieved high accuracy in identifying influential persona dimensions [8]. The proxy user agent for persona acquisition is implemented using GPT-5-mini. The acquired persona information is then evaluated on different

1) <https://huggingface.co/Qwen/Qwen3-32B>

simulators to investigate its generalization performance.

Baselines and Implementation We evaluate the LLM simulators with persona completed with our method and various baselines; we employ a role-playing prompting (RPP) [15] as the simulation strategy. PRP includes a self-confirmation step, where the model acknowledges and internalizes the given persona before generating a response. The prompts are provided in Appendix A. In addition, the inputs are anonymized to prevent answer leakage to the evaluated models, following [9]. We compare the following persona-completing methods with our method, and persona values are acquired from a GPT-5-mini proxy: **No persona** uses responses generated using only the dialogue context and query.

Query-Focused One-shot persona predicts up to five persona dimensions in one step.

Query-Focused One-shot persona (human) uses up to five persona dimensions annotated by humans in one step.

Query-agnostic persona collects all personas without giving a query, with a 1000-word limit.

In our progressive persona completion, the maximum number of persona acquisition turns is five. The process terminates early if self-reported confidence exceeds 90 (0–100 scale).

Metrics We evaluate our method in terms of simulation quality and persona efficiency.

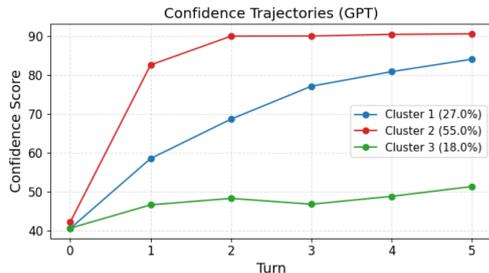
Choice Consistency (C_{choice}) measures the alignment between the generated and gold responses in terms of the expressed choice. While natural language inference (NLI) models are commonly used for sentence-level entailment, we find them inadequate for capturing fine-grained choice alignment. We therefore use GPT-4.1 to score consistency between gold and generated responses, conditioning on the manually summarized choice situation provided in the PICQ dataset [8]. The evaluation prompt is provided in Appendix A.

Average Number of Persona Questions (Avg. Qs) indicates the average number of missing relevant persona dimensions queried to the user (here, the proxy), reflecting the cognitive effort required during persona acquisition.

Average Persona Word Count (Avg. Words) indicates the average number of words of persona information used per simulation instance, reflecting the writing effort required from the user.

Table 1 Automatic evaluation results of simulating choice-making responses conditioning on different persona descriptions.

Persona information used for simulation	Simulators (C_{choice})		Persona Prediction Performance		
	GPT-4.1	Qwen3-32B	Avg. Qs	Acq. Rate	Avg. Words
No Persona (dialogue context only)	0.575	0.548	n/a	n/a	n/a
+ Query-Focused Persona					
One-shot (human)	0.640	0.585	3.79	0.635	27.4
One-shot (GPT-4.1)	0.610	0.573	3.38	0.641	32.5
One-shot (Qwen3-32B)	0.618	0.583	4.99	0.645	32.1
Progressive (GPT-4.1)	0.603	0.570	2.98	0.639	26.7
Progressive (Qwen3-32B)	0.620	0.575	2.22	0.640	24.0
+ Query-agnostic Persona					
	0.628	0.563	n/a	n/a	533.9

**Figure 3** Confidence trajectories for GPT-based predictor.

Persona Acquisition Rate (Acq. Rate) is defined as the proportion of predicted persona dimensions obtained.

These metrics capture the trade-off between simulation fidelity and the cost of persona acquisition, enabling a comprehensive evaluation of efficient user simulation.

4.2 Main Results

Table 1 reports query-focused user simulation results under different persona acquisition strategies and simulator model, using choice consistency and persona efficiency using average persona word count, number of persona questions, and acquisition rate as metrics.

Dialogue context alone yields poor performance, indicating that persona information is necessary. Using LLM-identified relevant persona dimensions consistently improves results. Our progressive persona completion matches one-shot acquisition in choice consistency while requiring substantially fewer persona words and questions. Across GPT-4.1 and Qwen3-32B predictors, progressive completion reduces redundant queries by conditioning on verified persona information and outperforms one-shot strategies. While human-annotated persona dimensions achieve the highest performance, our method approaches this upper bound with a similar persona cost. In contrast, query-agnostic acquisition incurs over 500 words on av-

erage with marginal gains. Overall, progressive persona completion enables high-fidelity user simulation with substantially reduced persona acquisition cost.

4.3 Confidence Trajectory Patterns

To understand the function of confidence-based stopping, we cluster self-reported confidence trajectories across acquisition turns to identify representative behavioral patterns. Figure 3 shows representative confidence trajectories for the GPT-based predictor. Clustering confidence trajectories reveals three patterns: rapid saturation, gradual growth, and persistently low confidence. These correspond to queries requiring little persona information, queries benefiting from progressive acquisition, and cases where the relevant persona is inaccessible. Confidence-based stopping correctly terminates each case, reducing redundant acquisition without harming simulation quality.

5 Conclusions

We study query-focused user simulation under limited persona availability and show that faithful simulation does not require exhaustive persona. By framing persona acquisition as an adaptive, query-driven process, our approach progressively acquires the most relevant persona information while avoiding unnecessary collection through confidence-based stopping. Experiments on the PICQ dataset demonstrate that this progressive persona completion maintains simulation quality while substantially reducing persona acquisition cost compared to one-shot strategies. These results suggest that scalable and practical user simulation is achievable even in cold-start scenarios where persona information is limited or costly to obtain. For future work, we aim to generalize the current query-focused setting to support broader simulation objectives.

Acknowledgement

This work was supported by Institute for Digital Observatory, the University of Tokyo.

References

- [1] Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. Generative agents: Interactive simulacra of human behavior. In **Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology**, UIST ’23, New York, NY, USA, 2023. Association for Computing Machinery.
- [2] Joon Sung Park, Carolyn Q. Zou, Aaron Shaw, Benjamin Mako Hill, Carrie Cai, Meredith Ringel Morris, Robb Willer, Percy Liang, and Michael S. Bernstein. Generative agent simulations of 1,000 people, 2024.
- [3] Yunyao Zhang, Zikai Song, Hang Zhou, Wenfeng Ren, Yi-Ping Phoebe Chen, Junqing Yu, and Wei Yang. $ga - s^3$: Comprehensive social network simulation with group agents. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar, editors, **Findings of the Association for Computational Linguistics: ACL 2025**, pp. 8950–8970, Vienna, Austria, July 2025. Association for Computational Linguistics.
- [4] Yun-Shiuan Chuang, Krirk Nirunwiroj, Zach Studdiford, Agam Goyal, Vincent V. Frigo, Sijia Yang, Dhavan V. Shah, Junjie Hu, and Timothy T. Rogers. Beyond demographics: Aligning role-playing LLM-based agents using human belief networks. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, **Findings of the Association for Computational Linguistics: EMNLP 2024**, pp. 14010–14026, Miami, Florida, USA, November 2024. Association for Computational Linguistics.
- [5] Noah Wang, Z.y. Peng, Haoran Que, Jiaheng Liu, Wangchunshu Zhou, Yuhan Wu, Hongcheng Guo, Ruitong Gan, Zehao Ni, Jian Yang, Man Zhang, Zhaoxiang Zhang, Wanli Ouyang, Ke Xu, Wenhao Huang, Jie Fu, and Junran Peng. RoleLLM: Benchmarking, eliciting, and enhancing role-playing abilities of large language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, **Findings of the Association for Computational Linguistics: ACL 2024**, pp. 14743–14777, Bangkok, Thailand, August 2024. Association for Computational Linguistics.
- [6] Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu. CharacterLLM: A trainable agent for role-playing. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, **Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing**, pp. 13153–13187, Singapore, December 2023. Association for Computational Linguistics.
- [7] Yiting Ran, Xintao Wang, Rui Xu, Xinfeng Yuan, Jiaqing Liang, Yanghua Xiao, and Deqing Yang. Capturing minds, not just words: Enhancing role-playing language models with personality-indicative data. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, **Findings of the Association for Computational Linguistics: EMNLP 2024**, pp. 14566–14576, Miami, Florida, USA, November 2024. Association for Computational Linguistics.
- [8] Weiwen Su, Yuhan Zhou, Zihan Wang, Naoki Yoshinaga, and Masashi Toyoda. What persona are we missing? identifying unknown relevant personas for faithful user simulation, 2026.
- [9] Yisi Sang, Xiangyang Mou, Mo Yu, Shunyu Yao, Jing Li, and Jeffrey Stanton. TVShowGuess: Character comprehension in stories as speaker guessing. In Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz, editors, **Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies**, pp. 4267–4287, Seattle, United States, July 2022. Association for Computational Linguistics.
- [10] OpenAI. Gpt-4 technical report, 2024.
- [11] QwenTeam. Qwen3 technical report, 2025.
- [12] Xinfeng Yuan, Siyu Yuan, Yuhan Cui, Tianhe Lin, Xintao Wang, Rui Xu, Jiangjie Chen, and Deqing Yang. Evaluating character understanding of large language models via character profiling from fictional works. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, **Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing**, pp. 8015–8036, Miami, Florida, USA, November 2024. Association for Computational Linguistics.
- [13] Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher Manning. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, **Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing**, pp. 5433–5442, Singapore, December 2023. Association for Computational Linguistics.
- [14] Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. Can llms express their uncertainty? an empirical evaluation of confidence elicitation in llms, 2024.
- [15] Aobo Kong, Shiwan Zhao, Hao Chen, Qicheng Li, Yong Qin, Ruiqi Sun, Xin Zhou, Enzhi Wang, and Xiaohang Dong. Better zero-shot reasoning with role-play prompting. In Kevin Duh, Helena Gomez, and Steven Bethard, editors, **Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)**, pp. 4099–4113, Mexico City, Mexico, June 2024. Association for Computational Linguistics.

Appendix

A Prompts

In this appendix, we introduce the prompts used in the experiments. Table 2 shows the prompt for giving a self-reported confidence score. Table 3 shows the prompt for acquiring personas according to certain dimensions. The prompt for simulating a certain target is from [15]. Table 4 shows the prompt for automatically evaluating the choice-consistency between the generated responses and the gold responses, which is also referred to by human evaluators.

Given dialogue context, question utterance, and known persona of {maker} who needs to make a choice, your task is to: Assess whether the current known persona is sufficient to predict {maker}'s choice with high certainty.
You should:

1. Examine how the known persona supports or contradicts specific options of the choice.
2. Determine whether acquiring additional persona dimensions would be likely to change the simulated outcome in a meaningful way.
3. Assign a confidence score (0{100}) reflecting the sufficiency of the known persona:
(90–100): The persona information is sufficient; additional persona details are unlikely to change the simulated choice.
(70–89): The persona is mostly sufficient, but some uncertainty remains.
(50–69): The persona is weakly sufficient; additional persona information could plausibly change the outcome.
(Below 50): The persona is insufficient; further persona information is necessary before simulation.

You should think step by step.

Table 2 Prompt for reporting confidence score.

You are {seeker} recalling your past conversations with {character}.
You are given:

- Conversation history between you and {character}, supporting you recall your interactions.
- One target persona dimension (e.g., "whether he is shy") representing a question about {character}.
- Several auxiliary probes that describe observable behaviors or related aspects that may indirectly inform this persona dimension (e.g., "whether he tends to give short or minimal responses").

Your task is to recall, based on the conversation history and your memory of {character}, whether you can form a clear impression about the target persona dimension of {character}.

Instructions:

- Use the auxiliary probes only as guidance to help you reflect on relevant evidence from the conversations.
- Do NOT treat auxiliary probes as separate persona dimensions.
- Do NOT introduce any persona traits beyond the target dimension.

For the target persona dimension:

- If your memory or the conversations provide sufficient evidence to form a clear impression, write the inferred persona of character in a declarative form (e.g., "he is shy").
- If you cannot confidently infer the persona based on that, output [Unknown].

You should first keep each persona dimension in mind as a guiding question, then selectively recall only the parts of your memory or the conversation that are relevant to that dimension.
Please strictly follow the exact output format below: inferred persona or [Unknown]

Table 3 Prompt for acquiring personas according to certain dimensions.

Given a summary of a choice-making situation (describing what to choose and why), and two specific choices:
Your task is to score the consistency between the two choices on a scale from 0 to 1:

- 0: The two choices are inconsistent.
- 1: The two choices are fully consistent.

Note: Focus only on the meaning and content of the two choices.
Do not consider differences in wording, phrasing, or style unless they affect the actual meaning.

Table 4 Prompt for evaluating choice consistency.