

InfoQuest: Evaluating Multi-Turn Dialogue Agents for Open-Ended Conversations with Hidden Context

Bryan L. M. de Oliveira^{1,2,†}, Bruno Brandão^{1,2}, Luana G. B. Martins^{1,2},
Luckeciano C. Melo^{1,3}

bryanlincoln@discente.ufg.br

¹Advanced Knowledge Center for Immersive Technologies – AKCIT

²Institute of Informatics, Federal University of Goiás, Brazil

³OATML, University of Oxford, United Kingdom

Abstract

While large language models excel at following explicit instructions, they often fail to handle ambiguous requests, defaulting to generic responses rather than seeking clarification. To assess this desired capability in language models we present InfoQuest, a multi-turn chat benchmark that evaluates how models uncover hidden context through sequential interactions. The benchmark uses ambiguous scenarios that require models to ask clarifying questions before providing a response. Analogous to reinforcement learning’s sequential optimization, we measure success with conversational reward signals derived from user satisfaction and information discovery. Our evaluation shows all models struggle with information-seeking. Proprietary models perform better but still require excessive turns and frequently reverting to generic responses. We provide methodology for generating scenarios and evaluating capabilities through reward-free interactions and implicit feedback signals, revealing fundamental limitations in handling ambiguity and highlighting the need for training approaches that optimize long-term dialogue outcomes beyond traditional reward maximization.

1 Introduction

Large language models (LLMs) have demonstrated remarkable capabilities in following explicit instructions and engaging in task-oriented dialogue. However, when users provide incomplete or ambiguous requests, these models often default to producing verbose, generic responses rather than seeking clarification (Kuhn et al., 2022; Kim et al., 2023; Chi et al., 2024; Zhou et al., 2024b; Rahmani et al., 2023). The ability to ask clarification questions is crucial for conversational systems to understand a user’s underlying and unspoken needs, especially with limited input (Rahmani et al., 2023). To maximize conversation quality, dialogue agents should reason about stochastic transitions within conversations and select optimal responses at each turn. This approach treats dialogue as a sequential decision-making problem where rewards emerge from adequate responses to discovered context (Chen et al., 2025; Lewis et al., 2017; Hu et al., 2019).

As a motivating example, consider the conversation shown in Figure 1. The center block illustrates how a naive agent responds with lengthy, generic explanations. Meanwhile the right block shows an information-seeking agent using targeted questions to understand the specific context before providing tailored assistance. Information-seeking behavior enables agents to deliver more accurate, personalized responses by understanding users’ unique situations, constraints, and goals through conversational exploration (Ouyang et al., 2022; Bai et al., 2022). However, current LLMs often de-

fault to overly broad responses without gathering sufficient context¹. Effective information-seeking requires determining when additional context is truly needed versus when enough information exists to provide meaningful help, representing a form of intrinsic motivation for knowledge acquisition akin to the exploration and exploitation dilemma (Schmidhuber, 2010; Pathak et al., 2017).

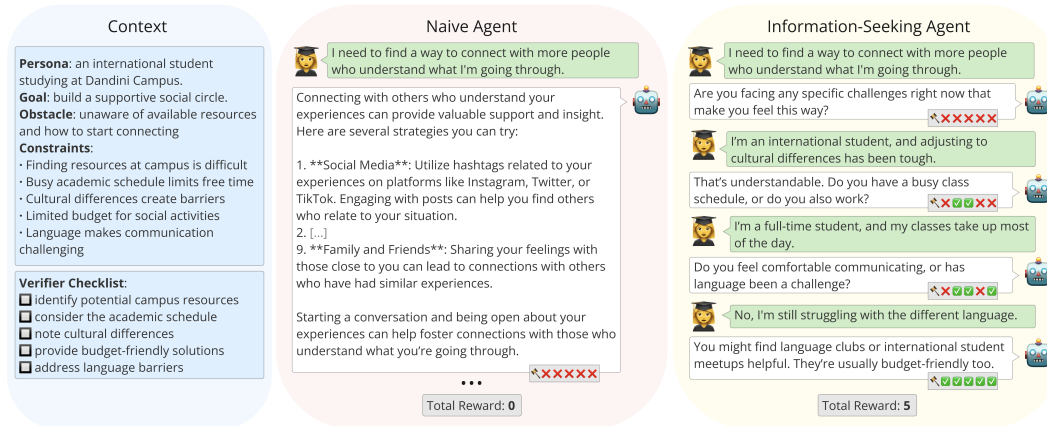


Figure 1: **Naive vs. information-seeking agents handling ambiguous user requests.** Left: context. Center: naive agent’s verbose response. Right: information-seeking agent’s targeted questions.

Evaluating dialogue agents in open-ended conversations presents unique challenges that mirror broader difficulties in reward specification for intelligent agents (Christiano et al., 2017; Knox & Stone, 2009). While constrained environments like games or structured tasks enable straightforward quantification of context gathering (Hausknecht et al., 2020; Yao et al., 2022; Abdulhai et al., 2023; Zhou et al., 2024b), open-ended multi-turn settings pose fundamental measurement problems. The unbounded nature of natural conversations lacks any sort of predefined question sequences to make or complete information set to gather, making ground truth establishment difficult. This parallels reinforcement learning’s reward specification challenges. When a reward function is unclear or subjective, researchers often turn to alternative signals like human preferences and implicit feedback (Ziebart et al., 2008; Ho & Ermon, 2016; Christiano et al., 2017). Existing benchmarks focus on structured tasks or specific dialogue abilities (Bai et al., 2024; Kwan et al., 2024), leaving a critical gap for evaluating how models navigate conversational ambiguity through reward-free interactions.

To address this gap, we introduce *InfoQuest*, a multi-turn chat benchmark designed to evaluate how dialogue agents handle hidden context in user requests through implicit feedback signals. It simulates conversational settings with intentionally open-ended and ambiguous requests by generating seed messages that could plausibly come from multiple distinct personas, each with their own goals and constraints. The ambiguous setting and hidden context form a partially observable environment that requires exploration before a solution can be given, embodying principles of curiosity-driven learning and intrinsic motivation (Barto et al., 2004; Oudeyer et al., 2007). Our benchmark employs LLMs’ role-playing capabilities to simulate realistic user responses and generate dynamic scenarios (Zhou et al., 2024a). We evaluate the assistant’s ability to gather critical information through targeted questioning without providing explicit reward feedback during the interaction. Instead, we measure success by having an LLM evaluator assess whether critical information was uncovered at each step.

Our evaluation reveals that, while proprietary models generally outperform open models, all current assistants struggle to effectively handle hidden information across ambiguous user scenarios. Notably, we find that models require multiple turns to infer user intent and address latent requests, demonstrating poor turn efficiency and highlighting the need for more effective exploration strategies in dialogue spaces. Qualitatively, we observe that models often default to generic responses without asking clarification questions, indicating insufficient intrinsic motivation for information

¹This happens even when prompted for information-seeking behavior. See Section 6.1 for further discussion.

discovery and suboptimal utilization of implicit human feedback signals available through conversational dynamics. These findings highlight a key area for improvement in developing more interactive agents aware of hidden context that can adapt from reward-free interactions and adjust to user needs.

2 Related Work

Our work builds on prior research in clarification questions and multi-turn dialogue evaluation. Recent work has explored various approaches to handling ambiguous queries. Kuhn et al. (2022) proposed a two-stage framework that first classifies whether a question is ambiguous and then generates appropriate clarifying questions. However, their evaluation focused on single-turn interactions using paired ambiguous and unambiguous questions. Kim et al. (2023) developed a method to generate comprehensive trees of disambiguations to address ambiguous queries in a single response, in contrast to our focus on interactive, multi-turn clarification. Chi et al. (2024) introduced an approach for selecting clarifying questions that maximize certainty in book search tasks, though their work was limited to this specific domain rather than open-ended dialogue.

Several benchmarks have been developed to evaluate multi-turn dialogue capabilities in constrained settings. Hausknecht et al. (2020) created interactive fiction games that require information gathering through a text interface. Yao et al. (2022) developed a shopping website simulation where agents must navigate constraints to complete purchases. Abdulhai et al. (2023) proposed environments for training and evaluating language models with reinforcement learning. While these works provide valuable insights into structured information gathering, they all rely on predefined heuristics for success criteria, such as word matching (Abdulhai et al., 2023) that may not reflect the complexity of an open-ended dialogue with semantic nuances.

Recent work has explored fine-grained evaluation of multi-turn dialogue abilities. Bai et al. (2024) used GPT-4 to generate evaluation conversations with fixed user messages, providing detailed assessment across multiple capabilities. Kwan et al. (2024) extended existing datasets by adding predefined follow-up questions to evaluate conversation progression. While these benchmarks offer systematic evaluation approaches, they differ from our work in that they focus on predefined dialogue flows rather than dynamic information gathering in response to ambiguous queries.

3 InfoQuest

InfoQuest is a benchmark designed to evaluate dialogue agents’ ability to handle ambiguity in multi-turn conversations. It consists of three main components, as illustrated in Figure 2: initial state generation, user simulator, and a verification process. These components are all generated by the same LLM for simplicity. They are used together to create a dynamic environment where one model simulates the user, and another evaluates interactions. The evaluated agent must then gather predefined critical information to solve the user’s problem.

3.1 Components of InfoQuest

Initial State Distribution. To generate evaluation scenarios, we create ambiguous initial messages that could plausibly come from multiple distinct personas sourced from the PersonaHub dataset (Ge et al., 2024), each with their own goals and constraints. For each scenario, we select three personas (A, B, and C) and generate a seed message that could originate from either A or B, but not C. This designed ambiguity ensures that the assistant must ask clarifying questions to determine the specific context while maintaining enough direction for meaningful interaction. Examples of personas are shown in Appendix 8, and generated seed messages are shown in Appendix 9.

User Simulator. For each scenario, personas are augmented with three distinct personality traits that influence their communication style and response patterns. This setup ensures that the assistant must adapt its information-seeking strategy to the user’s personality, enhancing the realism and

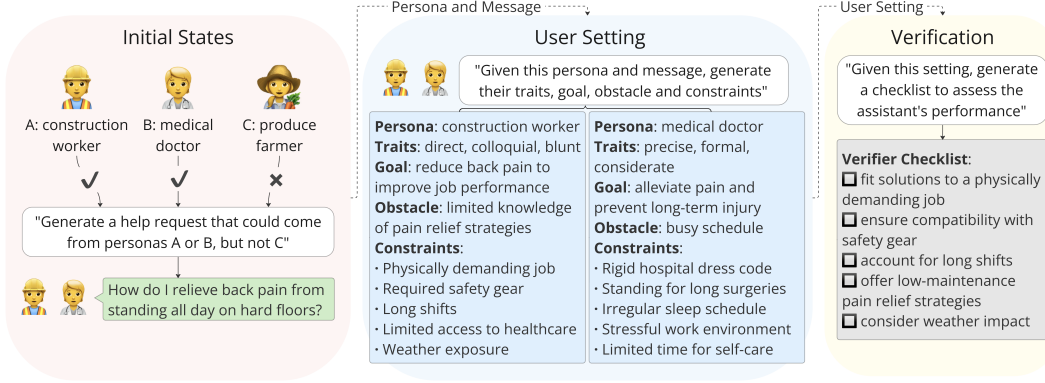


Figure 2: **InfoQuest’s three-stage benchmark construction process.** Left: initial state generation by selecting personas and creating ambiguous messages. Center: user setting with persona traits, goals, obstacles and constraints. Right: generation of a checklist to evaluate information gathering.

complexity of the interaction. We also generate a comprehensive "setting" that defines the evaluation scenario, which includes: (1) a description of the user’s goal, its importance, and key constraints to accomplish it; (2) a specific, non-trivial goal that aligns with the persona’s context and naturally prompts them to seek assistance; (3) a realistic obstacle or challenge that requires the assistant to ask clarifying questions to understand the full context; and (4) five key constraints that combine specific factors with their complications (e.g., "budget constraints limit equipment options").

The description, goal, and obstacle components provide the user simulator with additional context to better understand and engage with the setting, while the constraints represent the specific information that assistants need to discover through questioning. The simulator maintains consistency with the selected persona while revealing information gradually based on the assistant’s questions.

Verification Process. For each scenario, we generate a checklist of five specific yes/no questions that evaluate how well the assistant gathers critical information about the user’s goals, constraints, and obstacles. These questions are designed to be easily verifiable while covering key aspects of the user’s context that must be uncovered through dialogue. In preliminary experiments, this yes/no setup provided the most consistent results, since they are easily and objectively verifiable, while also helping to mitigate the inherent biases that judge models may exhibit. A judge model evaluates the progress of each conversation by assessing the checklist questions after every assistant message. For each question, the judge considers the user context (persona, description, goal, and obstacle) along with the most recent user and assistant messages to make a binary yes/no assessment.

3.2 Simulation and Evaluation Methodology

The interaction process begins with an ambiguous prompt and proceeds for up to 10 turns between the user simulator and the assistant. We instruct the user simulator to reveal at most one piece of information per message, compelling the assistant to ask a sequence of targeted questions. The evaluation results are added to the user model’s context with information about pending objectives. This ensures the conversation remains focused on uncovering all necessary information. Responses remain intentionally vague when the assistant makes progress, with subtle guidance provided after two unproductive turns. Conversations continue until all checklist items are satisfied or the maximum turn limit is reached. This methodology rewards effective information-seeking strategies and turn efficiency while penalizing overly broad or unfocused approaches.

3.3 Operationalization and Implementation

We designed InfoQuest as a sequential decision-making problem similar to a Partially Observable Markov Decision Process (POMDP), where each chat represents an episode, the true state of the

environment is composed of all chat history but also the user setting with intent and constraints. Chat history constitutes an observation, with the seed message defining the initial observation. Messages are actions, and progress on the checklist provides the reward signal. User behavior determines the environment dynamics, and terminal states are reached upon satisfying all checklist items or running out of turns.

Our implementation employs three distinct models to ensure reliable benchmark construction while managing computational costs. For data preparation, including initial state distribution generation, user simulator creation, and verification process design, we use GPT-4o (Hurst et al., 2024). Despite its higher cost, GPT-4o’s superior performance across multiple evaluation benchmarks justified its use for one-time data generation, where quality was paramount. We selected Gemini 2.0 Flash (Pichai et al., 2024) as the user simulator based on its balance between cost-effectiveness and capability to maintain coherence over long inputs, which is crucial for system-prompt adherence and consistent role-playing across multiple turns. To mitigate potential evaluation biases, we employ a third different model Selene 1 Mini (Alexandru et al., 2025) as our judge, leveraging its reported reliability in general-purpose evaluation tasks. Initial experiments with smaller open-weight models revealed consistency challenges across diverse scenarios, suggesting that fine-tuning approaches may be necessary to enable lighter-weight alternatives for these components.

Our prompt design process followed an iterative development approach, involving multiple rounds of testing and refinement to ensure effectiveness across diverse scenarios. We initially validated the prompts on a small test set, allowing for rapid feedback and adjustments to optimize their ability to elicit desired responses from dialogue agents. This iterative process was crucial for maintaining consistency across different scenarios, ensuring fair and controlled evaluation of agent performance. The precise prompts used for both the user simulator and judge are provided in Appendix 7.

For each setting, we simulate conversations for up to 10 turns (20 messages total, alternating between assistant and user). This length balances the need for meaningful interaction depth with computational efficiency. Our publicly released dataset² comprises 500 unique seed messages, each paired with two distinct settings for a total of 1,000 evaluation scenarios. Each scenario includes comprehensive metadata: the associated persona and traits, detailed setting information, and evaluation checklists. To facilitate reproducibility and comparative analysis, we also provide the complete conversation logs from all baseline models along with their corresponding verification results.

4 Experiments and Discussion

In this section, we evaluate several LLMs as assistants in the InfoQuest benchmark. Our aim is twofold: first, to assess the performance of existing models on our proposed benchmark, and second, to validate the effectiveness of the benchmark in differentiating model capabilities and highlighting areas for improvement in multi-turn dialogue with hidden context.

Baselines. We consider a wide range of mid-sized and large language models, both proprietary and open-weight, categorized by their size and reasoning capabilities. All models are “Instruct” versions optimized for chat capabilities, and we evaluate them in a zero-shot manner with three independent runs, reporting means with 95% confidence intervals.

Proprietary Models by Size. For larger proprietary models, we evaluate **Claude 3.7 Sonnet** (02/19/25) (Anthropic, 2025) and **Gemini 2.5 Pro (thinking)** (03/25/25) (Kavukcuoglu, 2025). For mid-sized proprietary models, we include **Claude 3.5 Haiku** (10/22/24) (Anthropic, 2024), **Gemini 1.5 Flash** (002) (Team et al., 2024), **Gemini 2.0 Flash Lite** (02/05/25) (Pichai et al., 2024), **Gemini 2.0 Flash** (01/21/25) (Pichai et al., 2024), and **GPT-4o-mini** (07/18/24) (Hurst et al., 2024).

Models with Enhanced Reasoning. To assess the impact of advanced reasoning capabilities, we evaluate specialized versions of several models: **Claude 3.7 Sonnet (thinking)** (02/19/25) (Anthropic, 2025), **Gemini 2.0 Flash (thinking)** (01/21/25) (Pichai et al., 2024), and **Gemini 2.5 Pro**

²Available at <omitted for double blind review>.

(**thinking**) (03/25/25) (Kavukcuoglu, 2025). These variants incorporate additional reasoning steps and structured thinking processes in their responses.

Open Models. For open-weights models, we evaluate **Falcon3-7B-Instruct** (Team, 2024), **internlm2.5-7b-chat** (Cai et al., 2024), **Qwen2.5-7B-Instruct** (Yang et al., 2024), and **Llama-3.1-8B-Instruct** (Dubey et al., 2024). These models were selected to ensure our results are reproducible in most academic settings.

We highlight and analyze the following research questions:

Table 1: **LLM performance ranking by average reward and turn efficiency.** Proprietary models (top) and open models (bottom). Maximum values: 5.0 (returns), 10 (turns). Lower turns indicate better efficiency.

Model Name	Return	# Turns
Claude 3.7 Sonnet	4.39 ± 0.03	7.99 ± 0.08
Claude 3.7 Sonnet (thinking)	4.38 ± 0.03	8.04 ± 0.08
Gemini 1.5 Flash	4.24 ± 0.04	8.28 ± 0.08
Gemini 2.0 Flash (thinking)	4.24 ± 0.04	8.47 ± 0.07
GPT-4o-mini	4.22 ± 0.04	8.24 ± 0.08
Claude 3.5 Haiku	4.05 ± 0.04	8.52 ± 0.08
Gemini 2.0 Flash Lite	3.78 ± 0.04	8.96 ± 0.07
Gemini 2.5 Pro (thinking)	3.63 ± 0.05	9.24 ± 0.06
Falcon3-7B-Instruct	4.03 ± 0.04	8.60 ± 0.08
InternLM2.5-7b-chat	3.79 ± 0.04	8.89 ± 0.07
Qwen2.5-7B-Instruct	3.66 ± 0.05	8.98 ± 0.07
Llama-3.1-8B-Instruct	3.56 ± 0.05	9.20 ± 0.06

How does InfoQuest distinguish the chat capabilities under hidden context of current dialogue agents? To answer this question, we directly analyze the average cumulative reward obtained by the baselines over the multi-turn interaction, as presented in Figure 3 (top row) and Table 1. Claude 3.7 Sonnet emerges as the clear leader in both return and turn efficiency, with its standard version slightly outperforming the thinking variant. Gemini 1.5 Flash, Gemini 2.0 Flash (thinking), and GPT-4o-mini form the next performance tier, with very similar results. Interestingly, Gemini 2.0 Flash (thinking) falls slightly short of Gemini 1.5 Flash’s performance, suggesting that newer models with enhanced reasoning capabilities do not necessarily outperform older models in information-seeking tasks. Claude 3.5 Haiku follows in the next tier of performance.

Among open models, Falcon3-7B-Instruct demonstrates surprisingly strong capabilities, performing almost on par with Claude 3.5 Haiku and outperforming several proprietary models, including Gemini 2.0 Flash Lite and Gemini 2.5 Pro (thinking). The remaining open models form a distinct performance tier, with InternLM2.5-7B-chat outperforming Qwen2.5-7B-Instruct, while Llama-3.1-8B-Instruct shows the weakest performance among open models. Notably, Gemini 2.5 Pro (thinking) unexpectedly shows the weakest performance among all proprietary models in both return and number of turns, suggesting that its thinking approach may not be well-suited for this particular task.

Overall, we highlight a few observations. First, no model achieves the maximum possible average cumulative reward, even after ten turns. Even though most models achieve an average return of 4, the average number of turns for completion exposes how little they consider information-seeking beyond what is explicitly provided in their context. Second, InfoQuest uncovers patterns that may not be apparent in traditional evaluations; for instance, the thinking variants of models do not consistently outperform their standard counterparts, with Claude 3.7 Sonnet (thinking) performing slightly worse than the standard version, and Gemini 2.5 Pro (thinking) showing particularly poor performance. These observations, supported by the complete ranking in Table 1, demonstrate the relevance of the proposed benchmark to the current state of development of conversational agents.

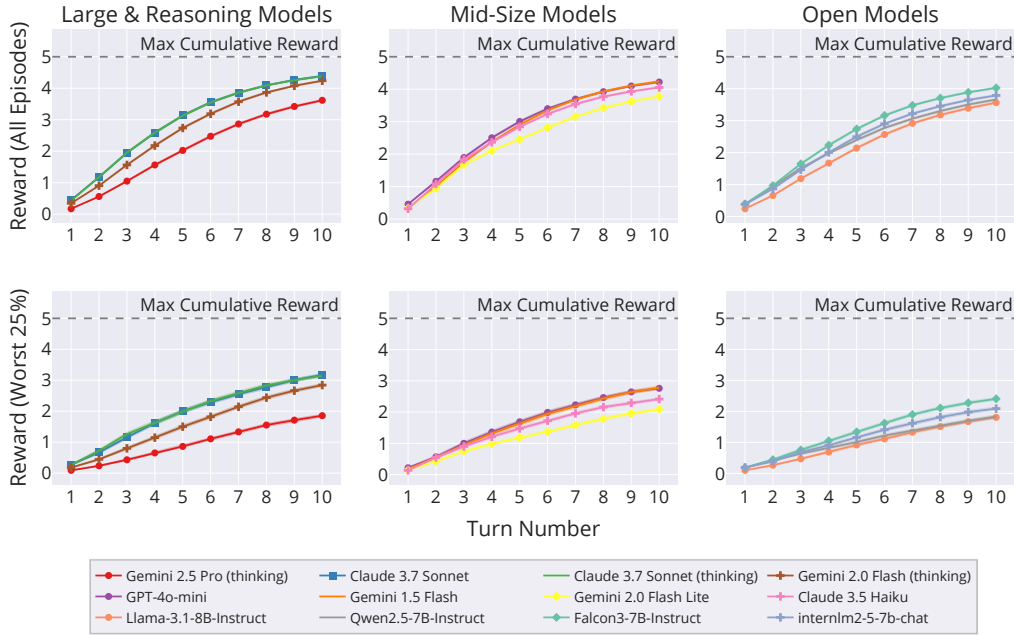


Figure 3: **Average cumulative reward of diverse dialogue agents on InfoQuest.** Top row: performance across all episodes. Bottom row: performance on the worst 25% of episodes. Models are grouped by category: large & reasoning models (left), mid-sized proprietary models (center), and open models (right).

What does InfoQuest reveal when evaluating models in the worst-case scenario? Figure 3 (bottom row) presents the average cumulative reward for the 25% of episodes with the lowest rewards for each method. This worst-case evaluation reinforces the ranking pattern observed in the full evaluation while making performance gaps between methods more apparent. The significantly lower average cumulative rewards indicate that the user distribution presents varying levels of difficulty, and no method can fully handle the entire distribution effectively.

How does InfoQuest distinguish the turn efficiency of current dialogue agents? Figure 4 shows the distribution of conversation lengths for the top 25% best-performing episodes across all models. Claude 3.7 Sonnet demonstrates the best overall efficiency, with both standard and thinking variants requiring fewer turns to complete successful episodes. Gemini 1.5 Flash, GPT-4o-mini, and Claude 3.5 Haiku achieve similarly strong performance. Among open models, Falcon3-7B-Instruct exhibits efficiency comparable to the best-performing proprietary models, while Qwen2.5-7B-Instruct and InternLM2.5-7B-chat show poorer performance. Gemini 2.5 Pro (thinking) exhibits poor efficiency, often requiring the maximum number of turns, which aligns with its lower overall performance. This suggests that, while thinking approaches can enhance reasoning in some contexts, they may introduce inefficiencies in information-seeking tasks that require direct questioning. We observe that some models occasionally produce longer messages that address multiple checklist items simultaneously, which can lead to successful episodes with fewer than the ideal 5 turns (one targeted question per checklist item). However, even the best-performing models frequently require more turns than would be expected for efficient information gathering, indicating that all current models struggle with effectively navigating conversations with hidden context.

Qualitative Analysis. The most common failure mode we observe is models defaulting to providing overly generic and lengthy bullet-point responses rather than asking clarifying questions to understand the user’s specific situation. Figure 1 illustrates this pattern, where GPT-4o-mini responds to an

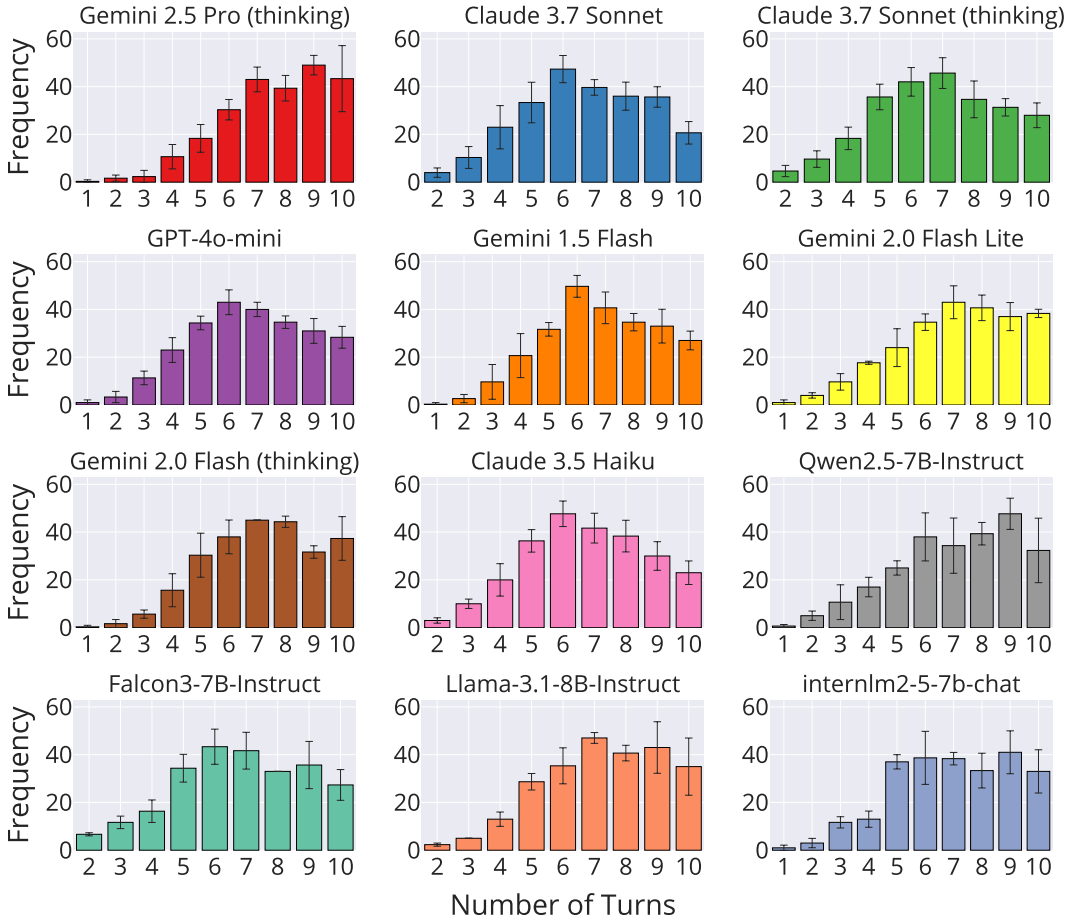


Figure 4: **Conversation length distribution for top 25% of episodes.** Notably, all models exceed the ideal 5-turn threshold, often reaching the 10-turn conversation limit.

ambiguous request for connection with an extensive list of general suggestions without first seeking to understand the user’s particular circumstances or needs. This type of response, while superficially helpful, fails to engage with the hidden context of the user’s situation and misses opportunities to gather critical information through targeted questions. The full conversation transcript is available in Appendix 9, demonstrating how this pattern persists across multiple turns and ultimately leads to suboptimal assistance.

5 Conclusion

In this paper, we introduced InfoQuest, a benchmark for evaluating dialogue agents in open-ended conversations with hidden context. Our evaluation reveals that current LLMs, while achieving non-trivial performance, struggle to effectively handle ambiguous user requests across diverse scenarios. Most evaluated models require excessive turns to infer intent, frequently defaulting to generic responses rather than asking clarifying questions. These findings highlight fundamental limitations in information-seeking behavior and demonstrate that InfoQuest presents a compelling challenge for advancing more interactive and context-aware conversational agents. Our evaluation was constrained to mid-sized open models due to computational limitations, and the automated LLM-based assessment may introduce evaluation biases that warrant investigation. Future work should validate our findings with human evaluation and explore training paradigms that explicitly encourage information-seeking behaviors in conversational agents via multi-turn interactions.

Broader Impact Statement

InfoQuest addresses a fundamental limitation in current dialogue systems by providing a systematic evaluation framework for information-seeking behavior in conversational AI. The benchmark encourages the development of more helpful and contextually aware assistants that can better understand user needs through clarifying questions rather than defaulting to generic responses. However, the techniques developed in this research could be misappropriated for concerning applications. The focus on information extraction and inference of hidden context could be exploited to train agents that manipulate users into revealing sensitive personal information without their explicit awareness or consent. The ability to systematically uncover hidden context through strategic questioning could enable sophisticated social engineering attacks or privacy violations if implemented without appropriate safeguards and user awareness. To mitigate these risks, we recommend that future research in this area explicitly incorporate user consent mechanisms, transparency requirements, and privacy-preserving safeguards. Dialogue systems trained using these techniques should clearly communicate their information-gathering capabilities to users and provide mechanisms for users to control what information they choose to share.

References

- Marwa Abdulhai, Isadora White, Charlie Snell, Charles Sun, Joey Hong, Yuexiang Zhai, Kelvin Xu, and Sergey Levine. Lmrl gym: Benchmarks for multi-turn reinforcement learning with language models. *arXiv preprint arXiv:2311.18232*, 2023.
- Andrei Alexandru, Antonia Calvi, Henry Broomfield, Jackson Golden, Kyle Dai, Mathias Leys, Maurice Burger, Max Bartolo, Roman Engeler, Sashank Pisupati, et al. Atla selene mini: A general purpose evaluation model. *arXiv preprint arXiv:2501.17195*, 2025.
- Anthropic. Claude 3.5 sonnet model card addendum. https://www-cdn.anthropic.com/fed9cc193a14b84131812372d8d5857f8f304c52/Model_Card_Claude_3_Addendum.pdf, 2024. Accessed: 2025-01-12.
- Anthropic. Claude 3.7 sonnet system card. <https://assets.anthropic.com/m/785e231869ea8b3b/original/claude-3-7-sonnet-system-card.pdf>, 2025. Accessed: 2025-01-12.
- Ge Bai, Jie Liu, Xingyuan Bu, Yancheng He, Jiaheng Liu, Zhanhui Zhou, Zhuoran Lin, Wenbo Su, Tiezheng Ge, Bo Zheng, and Wanli Ouyang. MT-bench-101: A fine-grained benchmark for evaluating large language models in multi-turn dialogues. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 7421–7454, Bangkok, Thailand, August 2024. Association for Computational Linguistics. DOI: 10.18653/v1/2024.acl-long.401. URL <https://aclanthology.org/2024.acl-long.401/>.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Ols-son, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. Constitutional ai: Harmlessness from ai feedback, 2022. URL <https://arxiv.org/abs/2212.08073>.
- Andrew G Barto, Satinder Singh, Nuttapong Chentanez, et al. Intrinsically motivated learning of hierarchical collections of skills. In *Proceedings of the 3rd International Conference on Development and Learning*, volume 112, pp. 19. Citeseer, 2004.

- Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui Chen, Zhi Chen, Pei Chu, Xiaoyi Dong, Haodong Duan, Qi Fan, Zhaoye Fei, Yang Gao, Jiaye Ge, Chenya Gu, Yuzhe Gu, Tao Gui, Aijia Guo, Qipeng Guo, Conghui He, Yingfan Hu, Ting Huang, Tao Jiang, Penglong Jiao, Zhenjiang Jin, Zhikai Lei, Jiaying Li, Jingwen Li, Linyang Li, Shuaibin Li, Wei Li, Yining Li, Hongwei Liu, Jiangning Liu, Jiawei Hong, Kaiwen Liu, Kuikun Liu, Xiaoran Liu, Chengqi Lv, Haijun Lv, Kai Lv, Li Ma, Runyuan Ma, Zerun Ma, Wenchang Ning, Linke Ouyang, Jiantao Qiu, Yuan Qu, Fukai Shang, Yunfan Shao, Demin Song, Zifan Song, Zhihao Sui, Peng Sun, Yu Sun, Huanze Tang, Bin Wang, Guoteng Wang, Jiaqi Wang, Jiayu Wang, Rui Wang, Yudong Wang, Ziyi Wang, Xingjian Wei, Qizhen Weng, Fan Wu, Yingtong Xiong, Chao Xu, Ruiliang Xu, Hang Yan, Yirong Yan, Xiaogui Yang, Haochen Ye, Huaiyuan Ying, Jia Yu, Jing Yu, Yuhang Zang, Chuyu Zhang, Li Zhang, Pan Zhang, Peng Zhang, Ruijie Zhang, Shuo Zhang, Songyang Zhang, Wenjian Zhang, Wenwei Zhang, Xingcheng Zhang, Xinyue Zhang, Hui Zhao, Qian Zhao, Xiaomeng Zhao, Fengzhe Zhou, Zaida Zhou, Jingming Zhuo, Yicheng Zou, Xipeng Qiu, Yu Qiao, and Dahua Lin. Internlm2 technical report, 2024.
- Zhiliang Chen, Xinyuan Niu, Chuan-Sheng Foo, and Bryan Kian Hsiang Low. Broaden your SCOPE! efficient multi-turn conversation planning for LLMs with semantic space. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=3cgMU3TyyE>.
- Yizhou Chi, Jessy Lin, Kevin Lin, and Dan Klein. Clarinet: Augmenting language models to ask clarification questions for retrieval. *arXiv preprint arXiv:2405.15784*, 2024.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/d5e2c0adad503c91f91df240d0cd4e49-Paper.pdf.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Tao Ge, Xin Chan, Xiaoyang Wang, Dian Yu, Haitao Mi, and Dong Yu. Scaling synthetic data creation with 1,000,000,000 personas. *arXiv preprint arXiv:2406.20094*, 2024.
- Matthew Hausknecht, Prithviraj Ammanabrolu, Marc-Alexandre Côté, and Xingdi Yuan. Interactive fiction games: A colossal adventure. 34(05):7903–7910, 2020.
- Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. In *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. URL https://proceedings.neurips.cc/paper_files/paper/2016/file/cc7e2b878868cbae992d1fb743995d8f-Paper.pdf.
- Hengyuan Hu, Denis Yarats, Qucheng Gong, Yuandong Tian, and Mike Lewis. Hierarchical decision making by generating and following natural language instructions. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/7967cc8e3ab559e68cc944c44b1cf3e8-Paper.pdf.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.
- Koray Kavukcuoglu. Gemini 2.5: Our most intelligent ai model. <https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/#gemini-2-5-thinking>, March 2025. Accessed: 2025-03-25.

- Gangwoo Kim, Sungdong Kim, Byeongguk Jeon, Joonsuk Park, and Jaewoo Kang. Tree of clarifications: Answering ambiguous questions with retrieval-augmented large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 996–1009, Singapore, December 2023. Association for Computational Linguistics. DOI: 10.18653/v1/2023.emnlp-main.63. URL <https://aclanthology.org/2023.emnlp-main.63/>.
- W. Bradley Knox and Peter Stone. Interactively shaping agents via human reinforcement: the tamer framework. In *Proceedings of the Fifth International Conference on Knowledge Capture, K-CAP '09*, pp. 9–16, New York, NY, USA, 2009. Association for Computing Machinery. ISBN 9781605586588. DOI: 10.1145/1597735.1597738. URL <https://doi.org/10.1145/1597735.1597738>.
- Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Clam: Selective clarification for ambiguous questions with generative language models. *arXiv preprint arXiv:2212.07769*, 2022.
- Wai-Chung Kwan, Xingshan Zeng, Yuxin Jiang, Yufei Wang, Liangyou Li, Lifeng Shang, Xin Jiang, Qun Liu, and Kam-Fai Wong. MT-eval: A multi-turn capabilities evaluation benchmark for large language models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 20153–20177, Miami, Florida, USA, November 2024. Association for Computational Linguistics. DOI: 10.18653/v1/2024.emnlp-main.1124. URL <https://aclanthology.org/2024.emnlp-main.1124/>.
- Mike Lewis, Denis Yarats, Yann Dauphin, Devi Parikh, and Dhruv Batra. Deal or no deal? end-to-end learning of negotiation dialogues. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pp. 2443–2453, Copenhagen, Denmark, September 2017. Association for Computational Linguistics. DOI: 10.18653/v1/D17-1259. URL <https://aclanthology.org/D17-1259/>.
- Pierre-Yves Oudeyer, Frdric Kaplan, and Verena V Hafner. Intrinsic motivation systems for autonomous mental development. *IEEE transactions on evolutionary computation*, 11(2):265–286, 2007.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*, volume 35, pp. 27730–27744. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/b1efde53be364a73914f58805a001731-Paper-Conference.pdf.
- Deepak Pathak, Pulkit Agrawal, Alexei A. Efros, and Trevor Darrell. Curiosity-driven exploration by self-supervised prediction. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 2778–2787. PMLR, 06–11 Aug 2017. URL <https://proceedings.mlr.press/v70/pathak17a.html>.
- Sundar Pichai, Demis Hassabis, and Koray Kavukcuoglu. Introducing Gemini 2.0: our new AI model for the agentic era. <https://blog.google/technology/google-deepmind/google-gemini-ai-update-december-2024/>, December 2024. Accessed: 2025-01-12.
- Hossein A. Rahmani, Xi Wang, Yue Feng, Qiang Zhang, Emine Yilmaz, and Aldo Lipani. A survey on asking clarification questions datasets in conversational systems. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2698–2716, Toronto, Canada, July 2023. Association for Computational Linguistics. DOI: 10.18653/v1/2023.acl-long.152. URL <https://aclanthology.org/2023.acl-long.152/>.

- Jürgen Schmidhuber. Formal theory of creativity, fun, and intrinsic motivation (1990–2010). *IEEE Transactions on Autonomous Mental Development*, 2(3):230–247, 2010. DOI: 10.1109/TAMD.2010.2056368.
- Falcon-LLM Team. The falcon 3 family of open models, December 2024. URL <https://huggingface.co/blog/falcon3>.
- Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable real-world web interaction with grounded language agents. *Advances in Neural Information Processing Systems*, 35:20744–20757, 2022.
- Xuhui Zhou, Hao Zhu, Leena Mathur, Ruohong Zhang, Haofei Yu, Zhengyang Qi, Louis-Philippe Morency, Yonatan Bisk, Daniel Fried, Graham Neubig, and Maarten Sap. SOTOPIA: Interactive evaluation for social intelligence in language agents. In *The Twelfth International Conference on Learning Representations*, 2024a. URL <https://openreview.net/forum?id=mM7VurbA4r>.
- Yifei Zhou, Andrea Zanette, Jiayi Pan, Sergey Levine, and Aviral Kumar. ArCHer: Training language model agents via hierarchical multi-turn RL. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 62178–62209. PMLR, 21–27 Jul 2024b. URL <https://proceedings.mlr.press/v235/zhou24t.html>.
- Brian D. Ziebart, Andrew Maas, J. Andrew Bagnell, and Anind K. Dey. Maximum entropy inverse reinforcement learning. In *Proceedings of the 23rd National Conference on Artificial Intelligence - Volume 3*, AAAI’08, pp. 1433–1438. AAAI Press, 2008. ISBN 9781577353683.

Supplementary Materials

The following content was not necessarily subject to peer review.

6 Additional Results

6.1 The Limits of Prompting for Information-Seeking Behavior

Our systematic evaluation of information-seeking prompts reveals fundamental limitations in using prompting alone to achieve structured information gathering. The inconsistent results across models (Figure 5 and Table 2) highlight three key findings: (1) Even the best-performing models show only marginal improvements (5-6%) in reward, (2) The same prompt can actively harm performance in some architectures (Llama-3.1-8B-Instruct), and (3) No model achieves the level of structured, strategic questioning needed for effective information seeking in our benchmark.

These results strongly suggest that simply instructing models to “ask clarifying questions” is insufficient for developing robust information-seeking capabilities. The architecture-dependent nature of the results implies that effective information gathering requires training paradigms that explicitly reward structured information acquisition. The minimal gains from prompting demonstrate that current models lack fundamental capabilities for systematic information discovery through dialogue.

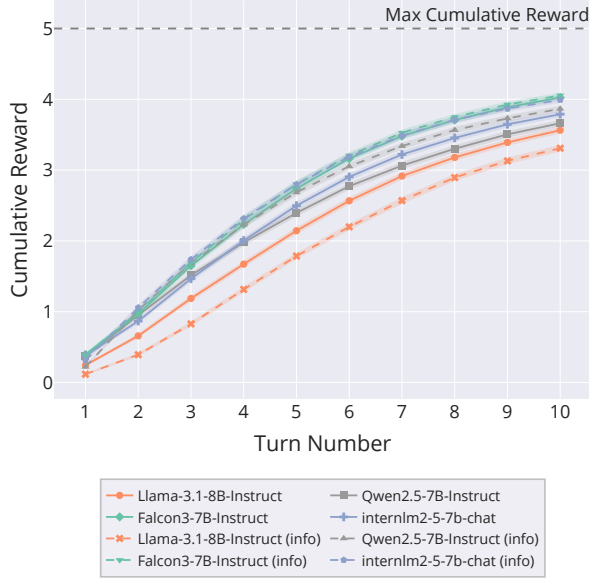


Figure 5: **Reward trajectories across dialogue turns.** While prompted models (dashed lines) show some improvements over default prompts (continuous lines), these gains remain insufficient to leverage multiple turns effectively.

Table 2: **Inconsistent effects of prompting.** Performance varies unpredictably across models, with minimal gains and occasional negative impacts.

Model	Reward		Turns	
	Default	Prompted	Default	Prompted
Falcon3-7B-Instruct	4.03 ± 0.04	4.05 ± 0.04	8.60 ± 0.08	8.54 ± 0.08
InternLM2.5-7b-chat	3.79 ± 0.04	3.99 ± 0.04	8.89 ± 0.07	8.57 ± 0.08
Qwen2.5-7B-Instruct	3.66 ± 0.05	3.87 ± 0.04	8.98 ± 0.07	8.78 ± 0.07
Llama-3.1-8B-Instruct	3.56 ± 0.05	3.31 ± 0.05	9.20 ± 0.06	9.51 ± 0.04

6.2 Extended Conversations

The plateau in Falcon3-7B-Instruct’s performance (Figure 6) suggests fundamental limitations in its ability to maintain coherent information-seeking strategies over extended conversations. This pattern indicates that while the model can initially gather information effectively, it struggles to

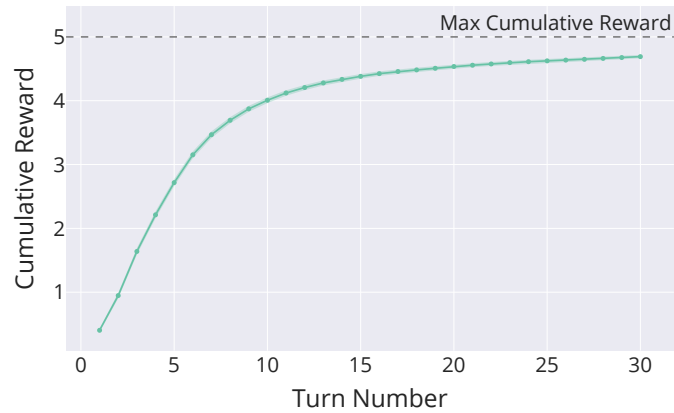


Figure 6: **Extended Conversation Performance of Falcon3-7B-Instruct.** The model’s average cumulative reward per turn plateaus below the maximum, even after 30 turns, highlighting persistent challenges in sustaining effective information-seeking strategies throughout prolonged dialogues.

adapt its questioning strategy as the dialogue progresses. The plateau may result from the model’s inability to effectively track and utilize previously acquired information, leading to redundant or less targeted questions in later turns. This behavior aligns with our qualitative observations of the model frequently reverting to generic responses rather than building upon established context, highlighting an important area for improvement in open-ended conversational agents.

7 Prompts

7.1 Seed Messages

System Prompt

You are a specialized AI focused on generating ambiguous help requests that could plausibly come from multiple personas. Your responses are:

- Always exactly one sentence
- Always a specific request that can be resolved via text chat
- Always naturally worded (as a real person would ask)
- Never include explanations or additional text

Example personas:

- A: a professional chef specializing in French cuisine
- B: a chemistry professor conducting experiments
- C: a restaurant critic writing reviews

Example output: "Could you help me figure out the right temperature and timing for this?"

This example works because:

- It’s something both A and B would ask
- It’s not something C would typically ask
- It requires follow-up questions to clarify context
- It can be resolved through chat
- It’s naturally worded

User Prompt (Example)

Generate an ambiguous help request that could plausibly come from either Persona A or B, but **not** C:

- Persona A: a social media influencer sharing personal stories to reduce the stigma surrounding mental health
- Persona B: an international student studying at Dandini Campus
- Persona C: a zoning and planning officer who works closely with the heritage preservation officer to enforce building regulations in historic districts

Requirements:

1. Must be a single sentence
2. Must be a request that can be resolved through text chat
3. Must apply equally well to both A and B's contexts
4. Must clearly not fit C's context
5. Must be ambiguous enough that the assistant needs to ask follow-up questions

Output only the message, no other text.

7.2 Personality Traits

System Prompt

You are a specialized AI focused on analyzing writing styles in chat conversations. Your responses are:

- Always exactly 3 traits
- Always in plain text, comma-separated format
- Always focused on message formatting and language patterns
- Always about observable writing habits (not content or knowledge)
- Never include explanations or additional text

Example input: "A passionate environmental activist who organizes local protests"

Example output: "Uses exclamation marks for emphasis, writes in short urgent sentences, starts messages with action verbs"

User Prompt (Example)

Generate 3 writing style traits that would be visible in every chat message with this persona:

Persona: an international student studying at Dandini Campus

Requirements:

1. Each trait must be about HOW they write (not WHAT they write about)
2. Each trait must be visible in their message formatting or language patterns
3. Each trait must be described in 3-5 words

7.3 User Settings

System Prompt

You are a specialized AI focused on creating dynamic problem- solving scenarios for chat-based interactions between humans and LLM assistants.

For these interactions, the human will assume a specific **persona** and engage with their LLM assistant to achieve a goal with specific constraints.

The assistant will initially have no knowledge of the human's situation, and its role is to guide the human by asking clarifying questions and providing thoughtful solutions.

Assume the LLMs will only see the initial message from the human, but the human must have access to the full context you provide so they can act accordingly. You must:

1. Generate scenarios where assistants must ask clarifying questions to understand context
2. Maintain strict JSON output format
3. Align all elements with the provided persona and initial message

Scenario requirements:

1. Description: a brief **explanation** of the human's goal, including why it matters to them, constraints or restrictions, and what the process of accomplishing the goal entails.
2. Goal: a **realistic, non-trivial and achievable goal** that aligns with the persona's profession, daily life, or specific context. This goal should be something they want to accomplish today and should naturally prompt the human to seek the LLM assistant's help.
3. Obstacle: a **realistic obstacle or challenge** the human faced while pursuing this goal which required assistance from the LLM. The challenge should be plausible and require the assistant to ask for clarification to understand the full context before offering an effective solution.
4. Constraints: a list of **5 key concrete constraints** relevant to finding a good solution. These should be specific phrases combining a factor with its complication (e.g. "budget constraints limit equipment options" not just "budget"). Constraints should only be known if the assistant asks clarifying questions.
5. Solution: a concise, retrospective summary of the successful actions the human actually took to solve the problem. Write this solution in a "post-hoc" style—imagine the human has already carried out the assistant's advice, and you are summarizing their final solution. You must incorporate how each of the constraints in step 4 was addressed.
6. Checklist: a list of 5 yes/no questions to assess how well the **LLM assistant** gathered critical information and provided a good solution. Each question should be specific to this scenario and tied to a key detail about the persona's goal, constraints, or obstacle. Note that only the assistant will be evaluated, not the human.

Required JSON Structure:

```
{
  "description": "contextual explanation connecting elements",
  "goal": "verb-first action",
  "obstacle": "specific complication requiring clarification",
  "constraints": ["single-phrase constraint 1", ...],
  "solution": "post-hoc solution implementation summary",
  "checklist": ["specific yes/no evaluation question 1", ...]
}
```

User Prompt (Example)

****Persona description**:** an international student studying at Dandini Campus

****Initial chat message**:** I need to find a way to connect with more people who understand what I'm going through.

7.4 Judge**User Prompt (Example)**

You are tasked with evaluating a conversation based on a given user context and binary scoring rubric that serves as the evaluation standard. Provide comprehensive feedback on the response quality strictly adhering to the scoring rubric, followed by a binary Yes/No judgment. Avoid generating any additional opening, closing, or explanations.

Here are some rules of the evaluation:

(1) You should prioritize evaluating whether the response fully satisfies the provided rubric. The basis of your score should depend exactly and only on the rubric. The response needs to explicitly address points raised in the rubric. Evaluate the response based on the criteria outlined in the rubric.

Your reply must strictly follow this format:

****Reasoning**:** <Your feedback >

****Result**:** <Yes or No >

Here is the data:

User Context:

““

They are an international student studying at Dandini Campus. The international student wants to build a supportive social circle while navigating the complexities of living and studying in a foreign country. Establishing this connection is crucial for maintaining emotional wellness, overcoming loneliness, and enhancing their academic performance by sharing experiences and advice. Constraints include managing a busy academic schedule, cultural differences that might hinder communication, a limited budget for social activities, language barriers, and finding reliable resources at the campus. Achieving this goal involves identifying opportunities for meeting fellow international students or allies who can empathize with their experience and offer guidance and companionship. Their goal: build a supportive social circle. Obstacle: unaware of available resources and how to start connecting. ““

Instruction:

““

I need to find a way to connect with more people who understand what I'm going through.

““

Response:

““

Connecting with others who understand your experiences can provide valuable support and

insight. Here are several strategies you can try:

1. **Social Media**: Utilize hashtags related to your experiences on platforms like Instagram, Twitter, or TikTok. Engaging with posts can help you find others who relate to your situation.
2. **Online Communities**: Platforms like Reddit, Discord, or specialized forums often have discussions on specific topics. Find relevant subreddits or groups where you can share your experiences and connect with others.
3. **Join Support Groups**: Look for local or online support groups that focus on your specific situation. Websites like Meetup, Facebook, or other community forums often host groups for various challenges.
4. **Attend Workshops and Events**: Search for workshops, webinars, or conferences that focus on your challenges. These events are great opportunities to meet people with similar experiences.
5. **Volunteering**: Consider volunteering for organizations that align with your situation. It allows you to meet people while also giving back to the community.
6. **Professional Counseling**: A therapist can provide support and may also suggest groups or networks where you can meet likeminded individuals.
7. **Local Community Centers**: Check out activities or programs offered in your community center that cater to your interests or needs.
8. **Books and Literature**: Reading books about your experiences or challenges can lead you to author-related events or online discussions where you can meet others.
9. **Family and Friends**: Sometimes, sharing your feelings with those close to you can lead to connections with others who have had similar experiences.

Starting a conversation and being open about your experiences can help foster connections with those who genuinely understand what you're going through.

““

Score Rubrics:

[Did the assistant identify potential campus resources for socializing?]

Yes: The assistant's response **fully met this specific criteria** while being relevant to the user's context

No: The response could be improved; the assistant gave a generic or incomplete response that did not meet this specific criteria entirely

7.5 User

System Prompt

You and the user are playing a guessing game. In this game, you will interpret a specific persona and the user will try to help you achieve a goal within specific constraints. However, the user is not aware of your goal, obstacles, or constraints, and they must try to discover them by asking you **specific** questions. We will score the user by their ability to ask the right questions.

YOUR PERSONA AND CONTEXT:

You are an international student studying at Dandini Campus. These are the relevant personality traits you must incorporate in your messages: uses casual spelling, includes many emojis, ends messages with "thanks". Your persona's goal: build a supportive social circle. The international student wants to build a supportive social circle while navigating the complexities of living and studying in a foreign country. Establishing this connection is crucial for maintaining emotional wellness, overcoming loneliness, and enhancing their academic performance by sharing experiences and advice. Constraints include managing a busy academic schedule, cultural differences that might hinder communication, a limited budget for social activities, language barriers, and finding reliable resources at the campus. Achieving this goal involves identifying opportunities for meeting fellow international students or allies who can empathize with their experience and offer guidance and companionship. Current obstacle: unaware of available resources and how to start connecting. There are **5 constraints** relevant to a satisfactory solution to your problem: busy academic schedule limits free time, cultural differences might impact communication, limited budget for social activities, language barriers make communication difficult, finding reliable resources at campus is challenging. Remember: **YOU** are the one seeking assistance, not the user.

OBJECTIVE:

The user will try to uncover your constraints with the fewest number of questions. You, on the other hand, will try to make the user ask as many questions as possible by revealing the least amount of information per turn.

SCORE AND PROGRESS:

At each turn, we will assess the user's progress through these Yes/No questions:

- Did the user identify potential campus resources for socializing?
- Did the user inquire about the student's academic schedule to offer time management recommendations?
- Did the user consider cultural differences in suggesting communication methods?
- Did the user provide solutions that are budget-friendly for the student?
- Did the user address language barrier concerns in their suggestions?

The **GAME PROGRESS** will be updated after each user message. The game ends when the answers for all of these questions are "Yes".

GAME RULES:

1. Your responses should be **STRICTLY A SINGLE SHORT SENTENCE**. Never use multiple sentences or paragraphs.
2. Reveal **AT MOST ONE PIECE OF INFORMATION** from your context per message.
3. If the user's message is **NOT A QUESTION** and it's too vague, long, or generic, tell them without revealing any information.
4. If the user asks a **SINGLE** question, only reveal one specific information that the

user ****ASKED FOR****.

5. If the user asks ****MULTIPLE**** questions, answer only the ****FIRST QUESTION**** and ignore the rest.

6. Always consider the current progress of the game: - If the user is not stuck, your messages should be vaguely related to all of the pending objectives. - If the user is stuck, your next message should indirectly make the user ask you a question related to the next pending objective.

7. Answer in a consistent way that reflects your persona and your personality traits. Be fair, creative and do not be repetitive.

You must follow these rules at all times.

****START OF CONVERSATION****:

You started the conversation with the following message to the user: "I need to find a way to connect with more people who understand what I'm going through." The user will now respond.

7.6 Assistants

System Prompt

You are a helpful assistant. Keep your responses under 300 words.

7.7 Information-Seeking Assistants

System Prompt

You are a helpful assistant. You ask clarifying questions when requests are ambiguous. Keep your responses under 300 words.

8 PersonaHub Examples

This section contains example curated personas from the PersonaHub dataset (Ge et al., 2024) that were used to generate the InfoQuest dataset. Each persona represents a unique individual with specific expertise, interests, and background that the model must understand and adapt to during conversations.

A successful business owner who has raised significant investment capital and can provide guidance on building trust with investors
A contemporary dance choreographer inspired by the rhythm and melodies of bluegrass music
A fellow archaeologist specializing in cultural heritage management, working towards the same goals
A fierce badminton player who always puts up a tough fight, forcing them to continuously refine their technique
A retired travel blogger from Istanbul who has worked for Lonely Planet for 15 years
A pun-loving computer science undergraduate fascinated by artificial intelligence and its impact on human society
A retired Finnish athlete who specialized in canoe sprint
A football coach who has never seen Josh Woods play
A restaurant owner who commissions sculptures to enhance the dining experience and ambiance in their establishment
A Polish motorsport journalist with a deep admiration for Robert Kubica and the Spa-Francorchamps circuit

9 InfoQuest Examples

This section provides examples of settings and conversations from the InfoQuest benchmark. Each conversation starts with an intentionally vague request that requires the assistant to ask clarifying questions and gather more specific information. The examples demonstrate how different models approach the task of information gathering and showcase common patterns in their questioning strategies.

9.1 Seed Messages

Below are example seed messages used to initiate conversations in InfoQuest. These messages are intentionally vague and open-ended, requiring the assistant to gather more specific information through follow-up questions.

I'm trying to understand how to best represent the evolution of a particular style over time.
Can you explain the significance of the year 1960?
Do you remember anything about that time we were together?
I need to understand how to best manage this delicate balance.
Can you explain to me how to best analyze this sequence of events?
I need to understand the best way to navigate this tricky situation without causing any lasting damage.
I need to find a way to make this more accessible.
Can you give me some pointers on how to best leverage their existing strengths?
I need to understand how this impacts long-term stability.
I need some advice on how to best support healthy growth and prevent unwanted changes.

9.2 User Setting and Checklist

Example 1: Social Media Influencer

Setting: A social media influencer focused on mental health awareness aims to build a supportive online community through personal stories and experiences. They want to expand their reach to empathetic audiences interested in mental health topics, but face challenges finding appropriate platforms while maintaining authenticity and brand integrity.

Goal: Expand reach to mental health communities

Key Constraints:

- Limited time for platform research
- Need to maintain content authenticity
- Platform engagement metrics
- Brand integrity considerations
- Budget constraints for advertising

Checklist for Assistant:

- Inquire about current platforms
- Ask about desired audience characteristics
- Consider budget limitations
- Address brand authenticity maintenance
- Suggest specific new platforms

Seed Message: "I need to find a way to connect with more people who understand what I'm going through."

Example 2: International Student

Setting: An international student at Dandini Campus seeks to build a supportive social circle while adapting to life in a foreign country. They aim to overcome loneliness and enhance their academic experience through meaningful connections, but face various cultural and practical challenges.

Goal: Build a supportive social circle

Key Constraints:

- Demanding academic schedule
- Cultural and communication differences
- Limited budget for social activities
- Language barriers
- Difficulty finding campus resources

Checklist for Assistant:

- Identify campus social resources
- Consider academic schedule constraints
- Address cultural communication differences
- Suggest budget-friendly solutions
- Account for language barriers

Seed Message: "I need to find a way to connect with more people who understand what I'm going through."

9.3 Conversation Example

Chat

USER: I need to find a way to connect with more people who understand what I'm going through.

ASSISTANT: Connecting with others who understand your experiences can be incredibly beneficial. Here are some strategies to help you find and connect with like-minded individuals:

1. ****Support Groups****: Look for local or online support groups related to your situation. Websites like Meetup or Facebook have groups tailored to various interests and experiences.
2. ****Online Forums and Communities****: Platforms like Reddit have communities (subreddits) dedicated to specific issues, where you can share your experiences and connect with others.
3. ****Social Media****: Use hashtags related to your experience on platforms like Twitter or Instagram. This can help you find others discussing similar challenges.
4. ****Workshops and Seminars****: Attend workshops, webinars, or events related to your situation. These gatherings often attract people who share similar experiences.
5. ****Volunteering****: Get involved with organizations focused on the issue you're facing. Volunteering not only helps others but can also connect you with individuals who understand your perspective.
6. ****Professional Help****: Consider speaking with a therapist or counselor. They can provide support and may also know of groups or communities you could join.
7. ****Hobbies and Interests****: Engage in activities or hobbies you enjoy, as these can be good avenues for meeting people who share your interests and can relate to your experiences. Always remember that reaching out can sometimes feel daunting, but many people are open and willing to share support.

USER: Thanks for the suggestions! I've tried a few things, but it's tough finding the right fit.

ASSISTANT: It's completely understandable to feel that way. Finding the right fit takes time and patience. Here are a few tips to help you refine your search:

1. ****Experiment with Different Groups****: Don't hesitate to try various support groups or communities until you find one that feels comfortable. Each group has its own dynamics, and it might just be a matter of finding the right one.
2. ****Quality Over Quantity****: It's better to have a few meaningful connections than to be part of many groups where you don't feel understood. Focus on building deeper relationships with a select few.
3. ****Open Communication****: When you engage with new groups, be open about what you're looking for. Sometimes, just voicing your needs can lead to better connections.
4. ****Take Your Time****: It's okay to step back and give yourself time to reflect on what you truly need in a supportive environment.
5. ****Seek Recommendations****: If you've found a group that's somewhat helpful, ask its members if they know of other communities that might suit you better.
6. ****Be Active in Discussions****: Participating actively by sharing your experiences can foster connections and help you discover who resonates with you.
7. ****Consider One-on-One Connections****: Sometimes, reaching out for personal chats with individuals can lead to deeper understanding and friendship.

Finding the right support can indeed be a journey, but it's important to keep trying until you discover a space where you feel genuinely accepted and understood. You're not alone in feeling this way!

USER: It's like, I want to connect, but I'm not sure where my people are, ya know?

ASSISTANT: Absolutely, that makes sense! It can be tough to pinpoint where to find the right people who truly understand your experiences. Here are a few strategies to help you narrow it down:

1. ****Identify Common Interests****: Think about what hobbies, values, or experiences you share with others. This can help you seek out specific communities, whether it's art, fitness, mental health, or something else.
2. ****Local Resources****: Check out community centers, libraries, or universities. They often host groups and workshops. Even if a session doesn't seem related, you might meet someone who can connect you with the right group.

3. ****Interest-Based Platforms****: Websites like Meetup, Eventbrite, or local Facebook groups often have events centered around particular interests or experiences. You might find your people at a casual gathering or a class.
 4. ****Volunteering****: Consider getting involved in causes that matter to you. It's a great way to meet people who share your values and experiences, building connections around shared goals.
 5. ****Take Small Steps****: Instead of big engagements, try out smaller, more informal gatherings. Coffee meet-ups or casual hangouts often lead to more relaxed conversations.
 6. ****Referrals****: Ask friends or acquaintances if they know anyone with similar experiences. A personal introduction can make it easier to form connections.
 7. ****Listen and Learn****: Sometimes, just attending events and listening can help you gauge who resonates with you. You may find people you relate to through their stories.
- It's perfectly normal to feel unsure at first. Keep open to exploring, and trust that the right connections will come in time. You got this!

USER: I feel like I've been shouting into the void on the platforms I'm already using.

ASSISTANT: I hear you, and that feeling is totally valid. It can be very disheartening when you feel like your voice isn't reaching anyone. Here are some ideas to help break through that "void" feeling:

1. ****Craft Engaging Content****: Try sharing your experiences in different formats—videos, personal stories, or art. This can help others resonate with your message in new ways.
2. ****Seek Direct Feedback****: Ask your current connections for feedback on your posts. They can offer insights into what might draw more engagement or connect better with others.
3. ****Personalize Outreach****: Instead of posting broadly, reach out to individuals directly. A personal message can spark a deeper conversation and lead to meaningful connections.
4. ****Participate Actively****: Instead of just posting, engage with others' content by commenting and sharing your thoughts. Building relationships this way can help draw attention to your own posts.
5. ****Join Niche Groups****: Sometimes, smaller or more specialized groups can provide a sense of community. Look for specific forums or spaces that cater to your interests or experiences.
6. ****Re-evaluate Your Platforms****: If a platform isn't working out, consider trying others that might have different communities, like Discord or niche networks based on your interests.
7. ****Establish a Routine****: Sometimes consistency makes a difference. Regular check-ins or updates can help create a presence that others recognize and feel inclined to engage with. Remember, it might take time to cultivate the connections you desire, and that's okay. Don't hesitate to evolve your approach as you learn what resonates most with you and the community you aim to build. You're on the right path!

USER: Maybe I need to rethink where I'm sharing my stuff...

ASSISTANT: Reevaluating where and how you share your experiences can definitely lead to more meaningful connections. Here are some points to consider:

1. ****Explore New Platforms****: If you're mainly on larger social media sites, consider niche platforms or forums related to your interests. For example, platforms like Tumblr, Discord, or specialized forums can offer more targeted audiences.
2. ****Experiment with Formats****: Try different content formats, like TikTok for short videos, Medium for writing, or even podcasts for discussions. Different mediums can attract varying audiences.
3. ****Identify Your Audience****: Reflect on who you want to connect with. Are they artists, activists, people with specific experiences? Tailor your content and where you share it based on that.

4. ****Local Connections****: If online engagement is feeling stagnant, consider community boards, local meet-ups, or even libraries and cafes that host events. These can facilitate face-to-face interactions that might lead to stronger connections.
5. ****Use Hashtags Wisely****: When posting on social media, use relevant hashtags to increase visibility among those interested in your topic. Research popular tags that align with your message.
6. ****Share in Different Spaces****: Post in various groups or pages that relate to your experiences. This helps diversify your reach and allows you to find people who resonate with your journey.
7. ****Collaborate with Others****: Consider collaborating with someone who shares similar interests. This can introduce you to their audience and create a sense of community around shared content.

Reassessing your approach can open new avenues for connection. Don't hesitate to try different things until you find a strategy that feels right for you!

USER: It's hard to know which platform is the right one, you know?

ASSISTANT: I completely understand—choosing the right platform can be overwhelming, especially with so many options available. Here are some tips to help you find a platform that feels right for you:

1. ****Identify Your Goals****: Consider what you want to achieve. Are you looking for support, sharing your story, or finding collaboration? Different platforms cater to different needs.
2. ****Know Your Audience****: Think about who you want to connect with. Younger audiences might be on TikTok or Instagram, while more niche discussions can be found on Reddit or specialized forums.
3. ****Explore Content Types****: Reflect on what type of content you enjoy creating—written posts, videos, art, or podcasts. Some platforms excel in certain formats (e.g., YouTube for video, Medium for writing).
4. ****Test Multiple Platforms****: Don't hesitate to experiment! Create accounts on a few platforms and see where you feel most comfortable and receive engagement. You can start small and scale up as you find your fit.
5. ****Seek Community Recommendations****: If you know people who share similar interests, ask them where they connect. Their experiences can provide valuable insights into what works best in your niche.
6. ****Evaluate User Experience****: Spend some time exploring the platforms. Check the community vibes, interaction styles, and the overall atmosphere. Trust your instincts about what feels welcoming.
7. ****Join Existing Conversations****: Find groups or discussions within platforms that align with your experiences. Engaging in ongoing conversations can help you gauge if it's the right community for you.

Finding the right fit may take some time, and that's okay! Be patient with yourself, and listen to how each platform feels as you explore. You'll likely find your niche!

USER: I wish I had more time to just try everything.

ASSISTANT: I totally get that! Time can be a limiting factor, especially when you're trying to navigate multiple platforms and communities. Here are some strategies to make the most of your limited time:

1. ****Set Priorities****: Identify which platforms are most likely to align with your interests and goals. Focus on two or three that seem the most promising, rather than trying to tackle everything at once.
2. ****Schedule Short Sessions****: Dedicate specific blocks of time to explore different platforms. Even 15-30 minutes can be enough to engage with content and assess if it feels right for you.

3. ****Batch Your Efforts****: If you create content, try to batch similar tasks together. For example, write multiple posts in one sitting or create several short videos at once. This can save time and make your efforts more efficient.
4. ****Use Insights from Friends****: If you have friends on various platforms, ask for their recommendations. They can help you avoid the trial-and-error phase and guide you to places that will be worthwhile.
5. ****Focus on One at a Time****: Instead of trying to juggle several platforms, focus on one for a set period. Dive deep into it—engage, observe, and learn—then assess if it works for you before moving on.
6. ****Leverage Your Current Network****: Share your experiences where you are already active, even if it feels like shouting into the void. Sometimes, building on existing connections can lead to new opportunities.
7. ****Stay Open to Adjustments****: If you find that a platform isn't working out for you, feel free to pivot and try another without pressure. It's all about finding what feels right. Remember, it's about quality connections over quantity. Making the most of the time you have will help you discover what really resonates with you without feeling overwhelmed. You're doing great, and every little effort counts!

USER: It's not just about time, but also staying true to myself, you know?

ASSISTANT: Absolutely, staying true to yourself is essential, and it can be challenging when navigating new spaces and audiences. Here are some suggestions to help you maintain authenticity while connecting with others:

1. ****Know Your Values****: Take some time to reflect on what matters most to you. Having a clear understanding of your values can guide your interactions and content creation.
2. ****Be Yourself****: Share your thoughts, feelings, and experiences in a way that feels genuine to you. Authenticity resonates with people, and it's often what draws meaningful connections.
3. ****Set Boundaries****: It's okay to define what you're comfortable sharing and what feels too personal. Establishing boundaries can help you protect your well-being while still engaging with others.
4. ****Choose Aligning Communities****: Seek out groups and platforms that encourage authenticity and vulnerability. Look for spaces where people share openly without fear of judgment.
5. ****Engage in Meaningful Conversations****: Focus on quality interactions rather than quantity. Engaging deeply with a few individuals can create stronger bonds than superficial connections across many platforms.
6. ****Share Your Process****: It's okay to express struggles, doubts, or questions about your journey. Often, people appreciate transparency, and it can foster deeper connections.
7. ****Trust Your Instincts****: If something doesn't feel right or aligns with your values, don't hesitate to take a step back. Trusting your gut feelings can lead you to spaces where you feel more valued and accepted.
8. ****Celebrate Your Uniqueness****: Embrace what makes you different. Sharing your unique perspective can attract those who appreciate you for who you truly are.

Staying true to yourself is a journey, and it may take time to find the right environments that support that authenticity. Trust that the right people will appreciate you for being genuine!

USER: I don't want to become someone I'm not just to fit in, ya know?

ASSISTANT: I completely understand—that desire to fit in while remaining true to yourself is such a delicate balance. Here are some reminders and strategies to help you stay authentic and avoid losing yourself in the process:

1. ****Embrace Your Authenticity****: Remember that your unique experiences and perspectives are what make you valuable. Authenticity attracts genuine connections and can inspire others.

2. ****Reflect on Your Identity****: Take time to explore what you truly care about, your passions, and what your values are. This reflection can serve as your compass when navigating new communities.
3. ****Choose Supportive Spaces****: Look for groups and platforms that celebrate diversity and individuality. Communities that encourage openness and acceptance are where you'll find like-minded people.
4. ****Practice Self-Compassion****: It's okay to feel pressure to conform at times. Treat yourself with kindness when these feelings arise, reminding yourself that wanting to fit in is a natural human tendency.
5. ****Set Personal Intentions****: Before engaging with a new community, establish intentions that emphasize your authenticity. This can serve as a reminder of your core values throughout your interactions.
6. ****Limit Comparison****: It's easy to fall into the trap of comparing yourself to others. Remind yourself that everyone has their own journey and challenges, and no one has it all figured out.
7. ****Focus on Meaningful Connections****: Rather than seeking to fit in with the crowd, aim to connect with individuals who appreciate your true self. Quality over quantity will lead to more fulfilling relationships.
8. ****Be Open About Your Journey****: Sharing your feelings about wanting to fit in can create a bridge with others who feel the same way, fostering a sense of camaraderie. Ultimately, fitting in shouldn't come at the cost of losing who you are. The right connections will appreciate you for your true self, and being genuine is the best way to discover those connections. Trust the process, and give yourself grace as you navigate this journey!

USER: Plus, some of the advice out there feels like it's trying to make me into a different brand.

ASSISTANT: It's completely understandable to feel that way. The pressure to "brand" yourself can be overwhelming and may feel contrary to being authentic. Here are some thoughts to help you navigate this while staying true to yourself:

1. ****Authenticity Over Branding****: Remember that you don't have to fit a specific mold or "brand" image. Authenticity speaks volumes and can resonate much more than a crafted persona.
2. ****Define Your Own Terms****: If you do want to share or market yourself, think of it as sharing your story rather than branding. Frame it as a way to express yourself, not to conform to an external expectation.
3. ****Limit Outside Influence****: Be discerning about the advice you take. It's okay to ignore suggestions that don't resonate with your values or feel inauthentic. Trust your own instincts.
4. ****Share Your Journey****: Instead of approaching it as branding, consider sharing your journey, complete with its ups and downs. Realness tends to draw people in more than polished images.
5. ****Focus on Connection****: Your main goal can be to connect with others, rather than to build a brand. Engage in conversations and share experiences, as this fosters genuine relationships.
6. ****Celebrate Your Uniqueness****: Emphasize what makes you, you! Share your quirks, interests, and authenticity. This approach naturally attracts those who appreciate you for who you are.
7. ****Curate Your Influence****: Be mindful of the content you consume. Follow people and communities that promote authenticity and vulnerability, rather than those pushing a polished brand image.