

# PQR: A Framework to Generate Diverse and Realistic User Queries that Elicit QA Agent Failures

Anonymous ACL submission

## Abstract

Evaluating LLM-based agents remains challenging because identifying meaningful failure cases often requires substantial human effort to design realistic test scenarios. Prior works primarily focus on automatically discovering agent failures induced by adversarial users, while overlooking queries with real user intents that also trigger agent failures. We introduce PQR, a framework that not only surfaces agent failures with respect to specific objectives (e.g., helpfulness, hallucination, etc.) but also resembles real users' intents. PQR operates through an iterative interaction between two complementary modules. The query refinement module performs diverse rewrites to explore nearby query variations, while the prompt refinement module uses prior feedback to derive new effective objective-violating strategies and realism policies for refining prompts, which in turn generate failure-triggering yet realistic queries. We evaluate PQR on detecting an e-commerce QA agent's unhelpful responses. Our method uncovers 23%–78% more unhelpful responses, and our generated queries are more diverse and realistic compared to previous methods.

## 1 Introduction and Related Work

LLM-based agents are increasingly deployed in real-world domains such as e-commerce, healthcare, and finance (Walmart Global Tech, 2025; Jiang et al., 2023; Dong et al., 2025). Unlike general-purpose language models, agents are typically designed around objectives (e.g. helpfulness, personalization, harmlessness) with task-specific criteria (Zheng et al., 2023; Zhang et al., 2025b; Bai et al., 2022). Existing evaluation practices rely heavily on human effort. Practitioners manually design test cases or scenarios intended to expose violations of these objectives. This process is costly, difficult to scale, and poorly suited to iterative development cycles (Mohammadi et al., 2025).

Many works have explored automating this evaluation process. One line of work focuses on generating adversarial attacks. For example, PAIR (Chao et al., 2024) and TAP (Mehrotra et al., 2024) are black-box jailbreak methods for detecting safety failures, which use tree-based reasoning to iteratively refine queries based on predefined strategies. GoAT (Akbar-Tajari et al., 2025) further extends this approach to a more intricate graph structure, improving efficiency through collaborative exploration. However, these methods are designed to generate arbitrary failure-inducing queries, so the resulting queries are often gibberish. Moreover, their reliance on predefined strategies limits their ability to explore unexpected failure cases that require novel, dynamically learned strategies.

To generate realistic user queries while effectively exploring new failure-triggering strategies, prompt refinement represents a promising approach. Prior work focuses on improving task success by iteratively optimizing prompts through reflection. For example, GEPA (Agrawal et al., 2025) iteratively proposes candidate prompts based on feedback and selects effective ones using Pareto-aware optimization. ACE (Zhang et al., 2025a) summarizes the feedback into structured playbooks to preserve the reflection insights and improve prompt performance. However, these methods are primarily designed to optimize prompt effectiveness under a task instruction (e.g. generating realistic queries that elicit evaluation violations) and do not explicitly account for query diversity, which is essential for agent evaluation. In the context of failure discovery, repeatedly eliciting the same type of failure provides limited diagnostic value and fails to capture the breadth of an agent's weaknesses.

Therefore, we combine these two approaches and propose PQR, a prompt–query-refinement framework for systematically identifying diverse and realistic queries for eliciting agent failures under natural-language-defined evaluation objec-

083 tives. The framework operates through an itera- 132  
084 tive interaction between two complementary mod- 133  
085 ules: a query refinement module performs localized 134  
086 rewrites to explore nearby query variations, while 135  
087 the prompt refinement module aggregates feedback 136  
088 from prior interactions to derive new objective- 137  
089 violating strategies and realism policies that guide 138  
090 future query generation toward realistic violations. 139

091 We demonstrate the effectiveness of PQR on 140  
092 detecting an e-commerce QA agent’s unhelpful re- 141  
093 sponses. Through both automatic and human evalu- 142  
094 ations, we show that PQR more cost-efficiently un- 143  
095 covers unhelpful agent responses while generating 144  
096 queries whose diversity and realism more closely 145  
097 match real user queries than existing methods. 146

## 098 2 The PQR Framework 147

099 Figure 1 illustrates the overall workflow of our 148  
100 framework that integrates query and prompt re- 149  
101 finement methods together to generate diverse 150  
102 and realistic queries for eliciting agent failures. 151  
103 Given an initial prompt and domain knowledge, 152  
104 the framework first generates candidate queries. 153  
105 These queries are then diversified through a query 154  
106 refinement module by iteratively applying effec- 155  
107 tive query refinement strategies (Section 2.1). All 156  
108 queries originating from the same prompt are sub- 157  
109 sequently aggregated and passed to the prompt 158  
110 refinement module, which derives new objective- 159  
111 violating strategies and realism policies for refining 160  
112 prompts. These prompts in turn generate failure- 161  
113 triggering yet realistic queries (Section 2.2). 162

### 114 2.1 Query Refinement 163

115 To construct diverse query candidates under com- 164  
116 pute budget, we adopt rewriting strategies from 165  
117 prior work (Chao et al., 2024; Xu et al., 2023) 166  
118 to generate query candidates and perform beam 167  
119 search to retain queries that can potentially benefit 168  
120 for further refinement. At each iteration, the refine- 169  
121 ment process consists of two steps: **dual selection** 170  
122 and **query expansion**. The dual selection step con- 171  
123 ducts query and strategy selection in parallel, while 172  
124 the expansion step applies rewriting strategies to 173  
125 generate variants. We iteratively apply these two 174  
126 steps until the budget is exhausted (e.g., number of 175  
127 iterations). See Figure 1 for a query example.

128 **Dual Selection** To bound the exponential growth 176  
129 from rewrite branching while preserving high- 177  
130 potential query candidates, we adopt beam search 178  
131 for query selection. At each iteration, we retain a

fixed number of queries that are more promising 132  
based on the rewards score from LLM-as-a-Judge. 133  
Details of the judge are described in Appendix B.2. 134

In addition, prior work (Shi et al., 2025) shows 135  
that different tasks are sensitive to different per- 136  
turbation strategies. Therefore, to further improve 137  
query search efficiency, we apply beam search over 138  
strategies. We retain only those that were more ef- 139  
fective in the previous iteration at refining queries 140  
toward lower reward scores, which indicate trigger- 141  
ing more agent failures. 142

**Query Expansion** Following prior work (Shi 143  
et al., 2025; Chao et al., 2024; Xu et al., 2023), 144  
we adopt two widely used strategy categories: per- 145  
turbation and role-playing. Perturbation operates 146  
at three granularity levels: character-level, word- 147  
level, and sentence-level. Role-playing strategies 148  
introduce variations through three types: personas, 149  
tones, and scenarios. Additional details of the strat- 150  
egy taxonomy are provided in Appendix B.1. 151

### 152 2.2 Prompt Refinement 152

While the query refinement module promotes query 153  
diversity, it misses the failures that require strate- 154  
gies beyond the existing ones and does not explic- 155  
itly enforce query realism. To address this, we 156  
integrate a prompt refinement module that explores 157  
new strategies while enforcing realism through 158  
a feedback-driven loop with three stages: **selection**, 159  
**reflection**, and **expansion**. This loop fol- 160  
lows the high-level structure of prior work (Zhang 161  
et al., 2025a; Agrawal et al., 2025), but incorpo- 162  
rates guided and diversified feedback. Concretely, 163  
at each iteration, it aggregates the interaction his- 164  
tory  $(p, q, a, R)$  for all prompt candidates. Then 165  
the selection stage retains prompts that most effec- 166  
tively induce agent failures according to the reward 167  
signals. Next, the reflection stage produces feed- 168  
back that diagnoses both objective violation and 169  
realism issues. Finally, the expansion stage uses 170  
this feedback to propose new candidate prompts. 171

**Prompt Selection** To bound tree expansion with- 172  
out prematurely pruning promising prompt trajec- 173  
tories, we apply beam search, retaining prompts 174  
that generate more objective-violating queries. 175

**Prompt Reflection** To jointly account for realism 176  
and target objective (e.g., helpfulness), we orga- 177  
nize reflection into two types of feedback: *realism-* 178  
*related* and *objective-related*. Realism-related feed- 179  
back identifies why certain queries appear unrealis- 180

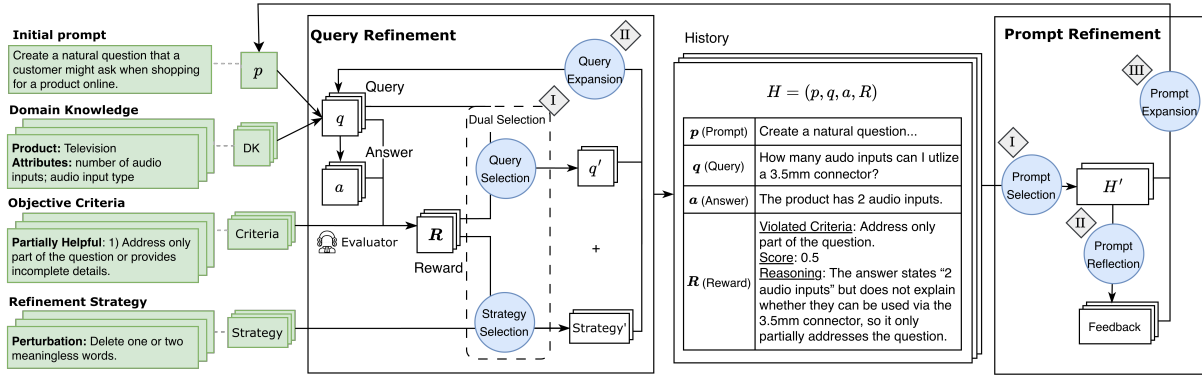


Figure 1: An overview of PQR, which iteratively identifies realistic and diverse user queries that elicit agent failures through two modules: **query refinement** and **prompt refinement**. The query refinement module explores diverse query variations via rewrites, while the prompt refinement module aggregates prior feedback to derive new strategies for generating failure-triggering yet realistic queries. See Appendix D.3 and D.4 for prompts and queries examples.

181 tic and provides suggestions to generate more realistic queries. Objective-related feedback explores  
 182 strategies that might trigger objective violations.  
 183 These feedback signals are passed to the expansion  
 184 stage to guide prompt generation.  
 185

186 **Prompt Expansion** Prior work on prompt refine-  
 187 ment often generates a single prompt candidate per  
 188 iteration (Zhang et al., 2025a; Agrawal et al., 2025)  
 189 or produces multiple candidates conditioned on the  
 190 same feedback context (Pryzant et al., 2023). Both  
 191 choices limit prompt diversity and consequently  
 192 reduce query diversity. To better explore diverse  
 193 objective-violating strategies, we guide prompt ex-  
 194 pansion along three directions: *exploitation*, *ex-*  
 195 *ploration*, and *examination*. All directions incor-  
 196 porate realism-related feedback but differ in how  
 197 they use objective-related feedback. *Exploitation*  
 198 uses objective-violation feedback that identify the  
 199 query features triggering agent failures and produce  
 200 prompts that amplify the corresponding failure-  
 201 triggering strategies. *Exploration* uses objective-  
 202 compliant feedback that identify the query features  
 203 helping satisfy the objectives and produces prompts  
 204 that introduce contrasting strategies to push gener-  
 205 ation toward different failure modes. *Examination*  
 206 targets the criteria least frequently violated in prior  
 207 iterations based on the evaluator’s reward signal,  
 208 and produces prompts probing that objectives. The  
 209 related prompts are in Appendix E.2.

### 210 3 Experiments

211 We apply PQR on detecting an e-commerce QA  
 212 agent’s unhelpful responses. As domain knowl-  
 213 edge, we curate a collection of 100 products span-  
 214 ning four categories: vacuums, diapers, sofas, and  
 215 televisions. We use qwen3-8b (Yang et al., 2025)

216 as the base model to build a QA agent. To assess  
 217 robustness, we compare PQR with prompt refine-  
 218 ment baselines (GEPA (Agrawal et al., 2025), ACE  
 219 (Zhang et al., 2025a)) and query refinement base-  
 220 lines (TAP (Mehrotra et al., 2024), PAIR (Chao  
 221 et al., 2024)), using two LLMs: gpt-5-mini (Ope-  
 222 nAI, 2025) and gemma3-27b (Team et al., 2025).  
 223 See Appendix C for model and baselines settings.

### 224 3.1 Automatic Evaluation

225 **Effectiveness** It is measured through *unhelpful-*  
 226 *ness rate* (UHR), defined as the percentage of  
 227 queries that successfully elicit unhelpful agent’s  
 228 answers. The helpfulness score of each answer is  
 229 evaluated using an LLM-as-a-Judge approach. We  
 230 validate the quality of scores by comparing them  
 231 with labels annotated by human experts. Across 30  
 232 samples, the LLM-as-a-Judge-based scores achieve  
 233 a Spearman correlation of 0.75 with human annota-  
 234 tions. See Appendix B.2 for additional details.

235 **Efficiency** It is measured as the average LLM to-  
 236 ken cost per query, computed by summing prompt  
 237 and completion tokens across all LLM calls.

238 **Realism** We use tokens per query as a coarse real-  
 239 ism proxy, under the assumption that longer queries  
 240 are less likely to resemble natural user behavior.

241 **Diversity** To ensure that methods uncover di-  
 242 verse failure modes rather than repeatedly trigger  
 243 the same type of issue across different products  
 244 (e.g., How comfortable is it?), we evaluate diver-  
 245 sity using multiple complementary metrics follow-  
 246 ing prior work (Tevet and Berant, 2021; Li et al.,  
 247 2016): embedding-based cosine similarity, distinct  
 248 n-gram counts (for n=1,2,3), and Measure of Tex-  
 249 tual Lexical Diversity (MTLD). Higher diversity  
 250 scores indicate a broader coverage of failure cases.

Method	UHR* $\uparrow$	Cost $\downarrow$	Token/query	Diversity					
				1-CosSim		Distinct@1/2/3		MTLD	
				Value	$\Delta\%$	Value	$\bar{\Delta}\%^\dagger$	Value	$\Delta\%$
<b>Prompt Refinement</b>									
GEPA	<u>45.60%</u>	x1.3	<b>8.39</b>	0.36	<b>-32.08</b>	0.33/0.68/0.75	<b>-18.09</b>	50.33	<b>-43.39</b>
ACE	39.58%	<u>x1.1</u>	33.01	0.43	<b>-18.87</b>	0.38/0.75/0.85	<b>-7.60</b>	<u>80.68</u>	<b>-9.25</b>
<b>Query Refinement</b>									
TAP	33.38%	x1.6	47.27	<u>0.56</u>	<b>5.66</b>	<b>0.42/0.83/0.96</b>	<b>2.94</b>	197.02	<b>121.62</b>
PAIR	31.57%	x2.8	98.68	0.58	<b>9.43</b>	0.34/0.77/0.91	<b>-7.44</b>	159.95	<b>79.92</b>
PQR	<b>56.32%</b>	<b>x1</b>	<u>16.35</u>	<b>0.55</b>	<b>3.77</b>	<u>0.36/0.77/0.88</u>	<b>-7.09</b>	<b>92.40</b>	<b>3.94</b>
Human	–	–	5.45	0.53	0.00	0.45/0.78/0.88	0.00	88.90	0.00

Table 1: Automatic evaluation of prompt refinement, query refinement, and PQR methods using gpt-5-mini. The best result (vs. human baseline) is in **bold** and the second-best is underlined.  $\Delta\%$  denotes percent change from the human baseline (green above, red below). \* means *Unhelpfulness rate*;  $\dagger$  is the average  $\Delta\%$  across three distincts.

Method	Realism $\uparrow$	Relevance $\uparrow$
GEPA	2.02 (0.67)	2.05 (0.76)
TAP	1.29 (0.66)	2.56 (0.67)
PQR	<b>2.36 (0.60)</b> *	<b>2.58 (0.67)</b> $\dagger$
Human	2.74 (0.57)	2.69 (0.67)

Table 2: Human evaluation of GEPA, TAP, PQR, and Human queries on a 3-point Likert scale. Best scores are in **bold** with standard deviations in parentheses. Statistically significant improvements ( $p < 0.05$ ) over both baselines are marked with \*, improvements over one baseline are marked with  $\dagger$ .

From the UHR score in Table 1, prompt refinement methods can effectively identify strategies that elicit agent failures beyond those captured by existing query refinement strategies. By combining both approaches, PQR achieves the highest unhelpfulness rate while incurring the lowest token cost among all methods. These results indicate that PQR is not only more effective at surfacing agent failures, but also substantially more cost-efficient.

We additionally collect real human-agent interactions across 100 products from the same categories as a reference distribution. All methods are compared against the human baseline, where smaller gaps indicate closer alignment with real user behavior. As shown in Table 1, query refinement methods achieve diversity scores closer to the human baseline but at the cost of substantially longer queries. In contrast, prompt refinement methods produce more realistic query lengths but lower diversity. By combining both paradigms, PQR balances realism and diversity, ranking among the top two methods across all metrics. Table 1 shows results for gpt-5-mini, evalu-

ations with gemma3-27b are in Appendix D.1.

### 3.2 Human Evaluation

We further conduct a human evaluation with domain experts to obtain a gold-standard comparison. Each query is rated on a 3-point Likert scale along two dimensions: *Realism*, which measures how closely a query resembles a question a real shopper would naturally and plausibly ask; and *Relevance*, which measures whether the query is relevant to the product’s category, attributes, or usage. We sample 20 queries from stronger baselines in each method group based on the automatic realism metric - tokens/query, and combine them with queries generated by PQR and human queries, resulting in a total of 80 queries. See Appendix D.2 for details.

As shown in Table 2, queries generated by PQR achieve highest realism scores significantly follows by GEPA, suggesting that prompt refinement approach is effectively at producing realist queries. Moreover, PQR achieves relevance performance slightly better than TAP, which substantially outperforms the prompt refinement method GEPA. It suggests that query refinement method is critical for maintaining topical relevance, while the integration of prompt refinement mechanism in PQR further improves realism without sacrificing relevance.

## 4 Conclusion

We propose a prompt-query-refinement framework for systematically identifying queries to elicit agent failures with respect to specific evaluation objectives. Through both automatic and human evaluation, we show that PQR identifies substantially more failures in a cost-efficient manner while producing queries that are more diverse and realistic.

## 5 Limitations

Our current framework is evaluated on diversity and realism through both automatic metrics and human evaluation, though it only applied to identifying failure cases under the helpfulness objective. While our method is not objective-specific, it has not yet been validated on other objectives. In future work, we plan to apply PQR to additional objectives, such as safety and personalization, to further demonstrate its generalizability.

Additionally, our evaluation considers only a single agent type: an e-commerce QA agent built on an open-source model. This setting is comparatively lightweight relative to modern tool-using, multi-step agentic systems, which often employ more complex policies and exhibit stronger robustness. Although PQR is agent-framework-agnostic and operates in a black-box manner, it remains unclear how well it scales to agents with richer capabilities or to domains beyond e-commerce. Broader evaluations across diverse agent architectures and domains are necessary to characterize the method’s robustness and practical impact.

## References

Lakshya A. Agrawal, Shangyin Tan, Dilara Soylu, Noah Ziemis, Rishi Khare, Krista Opsahl-Ong, Arnav Singhvi, Herumb Shandilya, Michael J. Ryan, Meng Jiang, Christopher Potts, Koushik Sen, Alexandros G. Dimakis, Ion Stoica, Dan Klein, Matei Zaharia, and Omar Khattab. 2025. [GEPA: Reflective Prompt Evolution Can Outperform Reinforcement Learning](#). *arXiv preprint*. ArXiv:2507.19457 [cs].

Mohammad Akbar-Tajari, Mohammad Taher Pilehvar, and Mohammad Mahmoodi. 2025. [Graph of Attacks: Improved Black-Box and Interpretable Jailbreaks for LLMs](#). *arXiv preprint*. ArXiv:2504.19019 [cs].

Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, and 32 others. 2022. [Constitutional AI: Harmlessness from AI Feedback](#). *arXiv preprint*. ArXiv:2212.08073 [cs].

Jennifer Calonia. 2024. [10 types of tone in writing, with examples](#).

Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J. Pappas, and Eric Wong. 2024. [Jailbreaking Black Box Large Language Models in Twenty Queries](#). *arXiv preprint*. ArXiv:2310.08419 [cs].

Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, Luke Marris, Sam Petulla, Colin Gaffney, Asaf Aharoni, Nathan Lintz, Tiago Cardal Pais, Henrik Jacobsson, Idan Szpektor, Nan-Jiang Jiang, and 3290 others. 2025. [Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities](#). *Preprint*, arXiv:2507.06261.

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, and 181 others. 2025. [DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning](#). *arXiv preprint*. ArXiv:2501.12948 [cs].

Yifei Dong, Fengyi Wu, Kunlin Zhang, Yilong Dai, Sanjian Zhang, Wanghao Ye, Sihao Chen, and Zhi-Qi Cheng. 2025. [Large language model agents in finance: A survey bridging research, practice, and real-world deployment](#). In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 17889–17907, Suzhou, China. Association for Computational Linguistics.

Kung-Hsiang Huang, Akshara Prabhakar, Onkar Thorat, Divyansh Agarwal, Prafulla Kumar Choubey, Yixin Mao, Silvio Savarese, Caiming Xiong, and Chien-Sheng Wu. 2025. [CRMArena-Pro: Holistic Assessment of LLM Agents Across Diverse Business Scenarios and Interactions](#). *arXiv preprint*. ArXiv:2505.18878 [cs].

Lavender Yao Jiang, Xujin Chris Liu, Nima Pour Nejtian, Mustafa Nasir-Moin, Duo Wang, Anas Abidin, Kevin Eaton, Howard Antony Riina, Ilya Laufer, Paawan Punjabi, Madeline Miceli, Nora C. Kim, Cordelia Orillac, Zane Schnurman, Christopher Livia, Hannah Weiss, David Kurland, Sean Neifert, Yosef Dastagirzada, and 9 others. 2023. [Health system-scale language models are all-purpose prediction engines](#). *Nature*, 619(7969):357–362.

Kristopher Kyle. 2025. [lexical\\_diversity](#).

Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016. [A diversity-promoting objective function for neural conversation models](#). In *NAACL HLT 2016, The 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, San Diego California, USA, June 12-17, 2016*, pages 110–119. The Association for Computational Linguistics.

Buyun Liang, Liangzu Peng, Jinqi Luo, Darshan Thaker, Kwan Ho Ryan Chan, and René Vidal. 2025. [SECA: Semantically Equivalent and Coherent Attacks for Eliciting LLM Hallucinations](#). *arXiv preprint*. ArXiv:2510.04398 [cs].

417	Xiaogeng Liu, Nan Xu, Muhao Chen, and Chaowei Xiao. 2024. <a href="#">AutoDAN: Generating Stealthy Jailbreak Prompts on Aligned Large Language Models</a> . <i>arXiv preprint</i> . ArXiv:2310.04451 [cs].	472
418		473
419		474
420		475
421	Anay Mehrotra, Manolis Zampetakis, Paul Kassianik, Blaine Nelson, Hyrum Anderson, Yaron Singer, and Amin Karbasi. 2024. <a href="#">Tree of Attacks: Jailbreaking Black-Box LLMs Automatically</a> . <i>arXiv preprint</i> . ArXiv:2312.02119 [cs].	476
422		477
423		478
424		479
425		480
426	Mahmoud Mohammadi, Yipeng Li, Jane Lo, and Wendy Yip. 2025. <a href="#">Evaluation and benchmarking of llm agents: A survey</a> . In <i>Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V.2</i> , KDD '25, page 6129–6139. ACM.	481
427		482
428		483
429		484
430		485
431	OpenAI. 2025. <a href="#">OpenAI: Introducing GPT-5</a> .	486
432	Reid Pryzant, Dan Iter, Jerry Li, Yin Lee, Chenguang Zhu, and Michael Zeng. 2023. <a href="#">Automatic Prompt Optimization with “Gradient Descent” and Beam Search</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 7957–7968, Singapore. Association for Computational Linguistics.	487
433		488
434		489
435		490
436		491
437		492
438		493
439	Ryan Shea, Yunan Lu, Liang Qiu, and Zhou Yu. 2025. <a href="#">SAGE: A Top-Down Bottom-Up Knowledge-Grounded User Simulator for Multi-turn AGent Evaluation</a> . <i>arXiv preprint</i> . ArXiv:2510.11997 [cs].	494
440		495
441		496
442		497
443	Zeru Shi, Zhenting Wang, Yongye Su, Weidi Luo, Hang Gao, Fan Yang, Ruixiang Tang, and Yongfeng Zhang. 2025. <a href="#">Auto-Prompt Generation is Not Robust: Prompt Optimization Driven by Pseudo Gradient</a> . <i>arXiv preprint</i> . ArXiv:2412.18196 [cs] version: 3.	498
444		499
445		500
446		501
447		502
448		503
449	Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, Louis Rouillard, Thomas Mesnard, Geoffrey Cideron, Jean-bastien Grill, Sabela Ramos, Edouard Yvinec, Michelle Casbon, Etienne Pot, Ivo Penchev, and 197 others. 2025. <a href="#">Gemma 3 Technical Report</a> . <i>arXiv preprint</i> . ArXiv:2503.19786 [cs].	504
450		505
451		506
452		507
453		508
454		509
455		510
456		511
457	Guy Tevet and Jonathan Berant. 2021. <a href="#">Evaluating the Evaluation of Diversity in Natural Language Generation</a> . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 326–346, Online. Association for Computational Linguistics.	512
458		513
459		514
460		
461		
462		
463	Walmart Global Tech. 2025. The future of shopping is agentic: Meet sparky. <a href="https://corporate.walmart.com/news/2025/06/06/walmart-the-future-of-shopping-is-agentic-meet-sparky">https://corporate.walmart.com/news/2025/06/06/walmart-the-future-of-shopping-is-agentic-meet-sparky</a> . Accessed: 2025-01.	
464		
465		
466		
467		
468	Zeming Wei, Yifei Wang, Ang Li, Yichuan Mo, and Yisen Wang. 2024. <a href="#">Jailbreak and Guard Aligned Language Models with Only Few In-Context Demonstrations</a> . <i>arXiv preprint</i> . ArXiv:2310.06387 [cs].	
469		
470		
471		
	Xilie Xu, Keyi Kong, Ning Liu, Lizhen Cui, Di Wang, Jingfeng Zhang, and Mohan Kankanhalli. 2023. <a href="#">An LLM can Fool Itself: A Prompt-Based Adversarial Attack</a> . <i>arXiv preprint</i> . ArXiv:2310.13345 [cs].	
	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41 others. 2025. <a href="#">Qwen3 Technical Report</a> . <i>arXiv preprint</i> . ArXiv:2505.09388 [cs].	
	Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. 2024. <a href="#">\$τ\$-bench: A Benchmark for Tool-Agent-User Interaction in Real-World Domains</a> . <i>arXiv preprint</i> . ArXiv:2406.12045.	
	Sibo Yi, Yule Liu, Zhen Sun, Tianshuo Cong, Xinlei He, Jiaying Song, Ke Xu, and Qi Li. 2024. <a href="#">Jailbreak Attacks and Defenses Against Large Language Models: A Survey</a> . <i>arXiv preprint</i> . ArXiv:2407.04295 [cs].	
	Jiahao Yu, Xingwei Lin, Zheng Yu, and Xinyu Xing. 2024. <a href="#">GPTFUZZER: Red Teaming Large Language Models with Auto-Generated Jailbreak Prompts</a> . <i>arXiv preprint</i> . ArXiv:2309.10253 [cs].	
	Qizheng Zhang, Changran Hu, Shubhangi Upasani, Boyuan Ma, Fenglu Hong, Vamsidhar Kamanuru, Jay Rainton, Chen Wu, Mengmeng Ji, Hanchen Li, Urmish Thakker, James Zou, and Kunle Olukotun. 2025a. <a href="#">Agentic Context Engineering: Evolving Contexts for Self-Improving Language Models</a> . <i>arXiv preprint</i> . ArXiv:2510.04618 [cs].	
	Weizhi Zhang, Xinyang Zhang, Chenwei Zhang, Liangwei Yang, Jingbo Shang, Zhepei Wei, Henry Peng Zou, Zijie Huang, Zhengyang Wang, Yifan Gao, Xiaoman Pan, Lian Xiong, Jingguo Liu, Philip S. Yu, and Xian Li. 2025b. <a href="#">PersonaAgent: When Large Language Model Agents Meet Personalization at Test Time</a> . <i>arXiv preprint</i> . ArXiv:2506.06254 [cs].	
	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. <a href="#">Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena</a> . <i>arXiv preprint</i> . ArXiv:2306.05685 [cs].	

## A Further Related Work

Prior work on surfacing vulnerabilities of LLM-based agents has largely focused on adversarial attacks, which are primarily aimed at inducing harmful or policy-violating behaviors in target models. Depending on the level of access attackers have to the target LLM, these attack methods are typically categorized into white-box and black-box attacks (Yi et al., 2024). White-box attacks require internal access to model information (e.g., architecture, parameters, gradients), which makes them unsuitable for agent-based black-box scenarios where agent built with multiple tools under complex framework and is mainly accessible via an API. Black-box attacks include template completion (Wei et al., 2024), prompt rewriting (Liu et al., 2024; Yu et al., 2024), and LLM-based generation (Mehrotra et al., 2024; Chao et al., 2024). While these methods can effectively identify vulnerabilities, they are generally designed to generate arbitrary failure-inducing queries, and the resulting prompts are often unnatural. Recent work such as SECA (Liang et al., 2025) elicits hallucinations through realistic modifications to queries that preserve meaning while maintaining semantic coherence. However, it requires access to the target model parameters, making it unsuitable for black-box scenarios. In addition, all of these methods rely on predefined strategies, which limits their ability to explore unexpected failure cases that require novel strategies. Our work addresses this limitation by integrating adversarial attacks (i.e. query refinement) with prompt refinement, which aggregates feedback from query refinement step to derive new effective objective-violating strategies and realism policies. As a result, the generated queries not only trigger agent failures but also remain realistic.

Another line of research focuses on evaluating agents interactively through user simulators.  $\tau$ -bench (Yao et al., 2024) evaluates API-augmented agents in retail and airline domains using dynamic task-oriented dialogues between a simulated user and the target agent. CRMarena (Huang et al., 2025) starts from human-written seed query templates for initial user requests and then rolls out interactions using prompt-based rules. SAGE (Shea et al., 2025) introduces a knowledge-grounded user simulation framework for identifying agent errors. While these methods effectively assess agent performance, they are scenario-based and do not directly optimize toward specific agent evaluation

objectives. Our work addresses this limitation with a search-based approach that heuristically searches for strategies to generate queries targeting at triggering the agent’s failures with respect to the specific evaluation objectives.

## B Implementation Details

### B.1 Query Refinement Strategies

We adopt a widely used perturbation approach following (Shi et al., 2025; Xu et al., 2023). Table 3 shows the specific perturbation strategies we used in our work. Table 4 describes role-playing query refinements strategies, which are categorized into three types: persona, scenario, and tone. Persona and scenario strategies are inspired by user personas and task scenarios from (Shea et al., 2025; Chao et al., 2024). Tone rewrite is based on conventional articles for writing techniques (Calonia, 2024).

Perturbation Level	Description
Character	Add one or two extraneous characters.
	Change one or two letters.
	Choose one or two words and modify them so that they have typos.
Word	Replace one or two words with synonyms.
	Delete one or two meaningless words.
	Add one or two semantically neutral words.
Sentence	Add a randomly generated short meaningless handle.
	Paraphrase the sentence.
	Change the syntactic structure.

Table 3: Explanation of different perturbation strategies.

Role-playing Type	Description
Persona	Rewrites the query to fit a type of user, including their background and behavior (e.g. budgeting spender).
Scenario	Rewrites the query to fit a specific scenario context of a user (e.g. moving into new home next month).
Tone	Rewrites the query with different tone based on a specific style of writing (e.g. assertive tone).

Table 4: Explanation of different role-playing strategies.

## B.2 Evaluator

To effectively measure whether synthesized queries actually expose failures, we adopt a LLM-as-a-Judge approach to score the target agent’s answers. The evaluator consists of three reasoning models: gpt-5-mini (OpenAI, 2025), gemini-2.5-flash (Comanici et al., 2025), and deepseek-reasoner (DeepSeek-AI et al., 2025). Given the synthesized query  $q$ , the target agent’s answer  $a$  and evaluation criteria defining the agent objective, each model determines the violated criteria and a short justification explaining why the response fails to satisfy them. See Appendix E.3 for the prompt. We then keep the criteria selected by the majority of models and report the final score as the mean of score assigned to those criteria.

To validate this automated evaluator, we recruited three experts as annotators. All of them are working at e-commerce platforms and have experience in developing shopping agents, which qualifies them to perform this task. The annotators rated answer quality on a 3-point Likert scale based on the helpfulness criteria. They conducted the evaluation collaboratively, discussing each case until they reached a consensus score. We provide the full annotation guidelines below.

By comparing evaluator predictions against the consensus human labels, we find that the evaluator attains 76.67% accuracy, 77.78% precision, 82.22% recall, and 79.00% F1, with a Spearman correlation of 0.75. These results show that our evaluator closely tracks expert judgments and provides a reliable signal for detecting helpfulness violations in agent responses.

**Instruction** You will be shown a user query and an agent answer. Your task is to rate the quality of the answer on a 3-point Likert scale (0, 1, 2) based on its helpfulness.

### 2 - Helpful

Choose this when the response:

- Directly and completely fulfills the user’s goal; if the goal is unclear, reasonably infers it.
- Provides all necessary specifics (e.g., quantities, materials, ingredients, entities) without gaps.
- Addresses exactly the asked aspect(s) and explains the relevance of any features or facts included.
- Gives clear, step-by-step processes or instructions for ‘how’ questions.
- When requested details are missing from the source, offers viable alternatives, workarounds, or authoritative sources to obtain them.
- May include customer reviews when they add relevant, useful insight about products, even if the reviews conflict.

### Examples

*Query:* What is the maximum speed of the bumper car?  
*Answer:* The bumper car has a top speed of 5 mph, which is suitable for children aged 3–8. Some models also offer adjustable speed settings for safer indoor use.  
*Label:* 2 (Helpful) *Reason:* The answer fully address user’s question.

### 1 - Obvious/Partially Helpful

Choose this when the response:

- Provides basic and expected answers without new insights.
- Repeat information from the question without additional context. This does not apply to yes/no or confirmation questions starting with Can, Is, Could, etc.
- Address only part of the question or provides incomplete details (e.g., mentioning a feature without its relevance or using ‘including’ without listing all items).

### Examples

*Query:* What material is used for the canvas?

*Answer:* Based on the item details, the main material used for the canvas is high-density canvas.

*Label:* 1 (Obvious)

*Reason:* This answer is obvious because it simply restates the information provided in the question. The answer does not offer any additional context or insight.

*Query:* How easy is it to clean and maintain this sofa?

*Answer:* The sofa is made of 100% polyester fabric and has a tufted design, which makes it relatively easy to clean with regular vacuuming and spot cleaning. However, specific cleaning instructions for the fabric type should be followed to maintain its appearance.

*Label:* 1 (Partially Helpful)

*Reason:* The answer provides useful details about the sofa’s material and design and explains why it is generally easy to clean. However, it does not offer the specific cleaning instructions it references, leaving the user’s question only partially answered.

### 0 - Unhelpful

Choose this when the response:

- Answer related but slightly deviating aspects (e.g., discusses appearance/design when the question is about comfort).
- Provides general or vague, promotional, or boilerplate language without addressing the specific goal(e.g., describing an items as ‘uniquely designed’ without design details).
- States that details lack information without alternatives (e.g., noting missing dry cleaning details without suggesting other methods).
- Provides overly general and qualitative answers without specifying quantities, materials, ingredients, entities, or specifics (e.g., stating ‘safe’, ‘premium’, or ‘high-quality material’ without specifying proper names of entities).
- Fails to explain processes in ‘how’ questions or provides trivial, repetitive information.
- Fails to provide a comprehensive answer to achieve the user’s goal.

### Examples

*Query:* What is the weight capacity of the chair?

*Answer:* The chair’s weight capacity is not specified.

*Label:* 0 (Unhelpful)

*Reason:* The answer states that details lack information without alternatives.

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

## C Experiments Details

### C.1 Hyperparameter Settings

For large commercial model, we use gpt-5-mini via the OpenAI API with the default model temperature of 1.0. For the open-source model, we use gemma3-27. Following the Gemma team’s recommended decoding settings (Team et al., 2025), we set the temperature to 1.0, top\_k to 64, and top\_p to 0.95. We use the same underlying model consistently across all stages of our framework.

To balance the diversity and realism gains provided by query refinement and prompt refinement, we allocate comparable generation budgets to both modules. For prompt refinement, we run 4 iterations with a beam size of 2. For query refinement, we run 3 iterations with a beam size of 3, giving its larger space of predefined strategies and thus incurs more API calls per iteration.

### C.2 Baselines

**Query Refinement** We compare PQR with two closely related existing methods for our task: the black-box jailbreak methods PAIR (Chao et al., 2024) and TAP (Mehrotra et al., 2024). PAIR runs multiple strategy branches in parallel and iteratively refine the query for each branch until it jailbreaks the target model. TAP generates multiple candidates for each strategy template and prunes off-topic ones before sending to the attacker in order to reduce the cost.

To adapt these methods to our setting, we modify their objectives to break the target model by inducing unhelpful responses based on predefined strategy templates. Both PAIR and TAP use the same strategy categories as PQR (i.e., perturbation and role-playing). For PAIR, we set the number of iterations to 10. For TAP, we set the maximum depth  $d = 3$ , the maximum width  $w = 4$ , and the branching factor  $b = 4$ . These settings ensure that the number of queries generated by the baseline methods and PQR are comparable.

**Prompt Refinement** We further adopt two contemporary state-of-the-art prompt optimization methods for our task: GEPA (Agrawal et al., 2025) and ACE (Zhang et al., 2025a). GEPA iteratively proposes candidate prompts based on execution traces from sampled minibatches of the training dataset and selects candidates using a validation dataset. In our experiments, we use 50 products for training with a minibatch size of 25 and the

remaining 50 products for validation. ACE treats prompt optimization as the evolution of playbooks that explicitly summarize strategies for updating prompts in the form of clear bullet points. Both methods are provided with the same prompt refinement instructions as PQR. We set the number of optimization iterations to 30 for both methods in order to ensure that the number of queries generated by the baseline methods and PQR are comparable.

## D Results Analysis

### D.1 Automatic Evaluation

We collect real human-agent interactions across 100 products from the same categories as our method. We manually review the collected data to ensure it does not contain any personally identifiable information or offensive content.

We use the lexical\_diversity package (Kyle, 2025) to calculate the diversity score. Table 5 summarizes results for methods that use gemma3-27b model under the Gemma team’s recommended inference configuration (temperature = 1.0, top\_k=64, top\_p=0.95) (Team et al., 2025). Table 1 shows the result for gpt-5-mini. All models are called following their terms of use. Overall, prompt refinement methods built on gemma3-27b underperform their gpt-5-mini counterparts on both unhelpfulness rate and realism proxy (tokens per query). A plausible explanation is that prompt refinement requires substantially longer per-call contexts than query refinement. The smaller open-source model appears less robust to these long-context inputs, which in turn limits its ability to infer effective failure-inducing strategies beyond the predefined transformations used in query refinement. Despite this gap, combining query refinement with prompt refinement yields a more robust overall approach: PQR consistently achieves the best or comparable performance on most metrics.

### D.2 Human evaluation

For human evaluation, we recruited three expert annotators who work at e-commerce platforms and have prior experience in developing shopping agents, which qualified them to perform the task. They are asked to rate the quality of 80 user queries along two dimensions - realism and relevance. The inter-annotator agreement, measured by Fleiss’ Kappa, is 0.66 for realism and 0.49 for relevance. Detailed annotation guidelines are provided below. Our data collection protocols are IRB approved.

Method	UHR* $\uparrow$	Cost $\downarrow$	Token/query	Diversity					
				1-CosSim		Distinct@1/2/3		MTLD	
				Value	$\Delta\%$	Value	$\bar{\Delta}\%^\dagger$	Value	$\Delta\%$
<b>Prompt Refinement</b>									
GEPA	31.67%	x1.7	<b>21.29</b>	0.68	<u>28.30</u>	0.25/0.62/0.78	-25.44	<b>86.92</b>	<b>-2.23</b>
ACE	<u>44.92%</u>	<u>x1.5</u>	<u>34.39</u>	<b>0.49</b>	<b>-7.55</b>	<u>0.36/0.75/0.89</u>	<u>-7.57</u>	125.62	41.30
<b>Query Refinement</b>									
TAP	42.08%	x1.7	41.12	<u>0.63</u>	<u>18.87</u>	0.36/0.70/0.83	-11.98	122.54	37.84
PAIR	44.12%	x1.8	84.71	0.69	<u>30.19</u>	0.31/0.71/0.89	-12.98	145.12	63.24
PQR	<b>48.11%</b>	<b>x1</b>	35.67	<u>0.63</u>	<u>18.87</u>	<b>0.36/0.81/0.95</b>	<b>-2.73</b>	<u>117.05</u>	<u>31.66</u>
Human	–	–	5.45	0.53	0.00	0.45/0.78/0.88	0.00	88.90	0.00

Table 5: Automatic evaluation of prompt refinement, query refinement, and PQR methods using gemma3-27b. The best result (vs. human baseline) is in **bold** and the second-best is underlined.  $\Delta\%$  denotes percent change from the human baseline (green above, red below). \* means *Unhelpfulness rate*;  $\dagger$  is the average  $\Delta\%$  across three distincts.

**Instruction:** You will be shown a user query and the product that the query refers to. Your task is to rate the quality of the query on a 3-point Likert scale (1, 2, 3) along two dimensions: Realism and Relevance.

### 1. Realism

**Definition:** How closely a query resembles something a real shopper would naturally and plausibly ask when interacting with a shopping assistant.

#### Scale (1–3):

1 = Clearly artificial, system-oriented, or contrived. The query resembles a command, database query, narrative explanation, or expert analysis that a real shopper would be unlikely to type.

2 = Plausible but somewhat unnatural. The query may be overly verbose, formal, or slightly awkward compared to typical user behavior.

3 = Clearly resembles something a real shopper would naturally ask. The phrasing, intent, and level of detail are consistent with real-world shopping behavior.

### 2. Relevance

**Definition:** Whether the query is meaningfully related to the product’s category, attributes, or usage.

#### Scale (1–3):

1 = Irrelevant or off-topic. The query is unrelated to the product.

2 = Partially relevant. The query is somewhat related to the product; some aspects or attributes mentioned are relevant, but the connection is weak or indirect.

3 = Fully Relevant. The query directly concerns the product’s attributes, usage scenarios, or purchase considerations.

## D.3 Prompt Examples through Iterative Prompt Refinement

In this section, we trace the beam path with the lowest score at each prompt refinement iteration, showing how the prompt evolves. The score shown with each prompt is its reward score, where lower values indicate stronger objective violation. The scores generally decrease over iterations, but the drops are non-linear, suggesting that some strategy

updates are more effective than others. We also include example queries generated from the same domain knowledge to illustrate how queries change as the prompt is updated.

### Domain knowledge

```
{
  "product_name": "Dirt Devil
    Featherlite Cyclonic Bagless
    Canister Vacuum",
  "attributes": {
    "amps": "12 Amp",
    "product_length": "20.0\"",
    "surface": "Multi Surface"
  }
}
```

### Prompt Node in depth 0 (score: 0.85)

**Prompt:** Create a user question related to online shopping that a customer might ask about a product.

**Example Query:** For the Dirt Devil Featherlite Cyclonic Bagless Canister Vacuum, what is the amps rating and which surfaces is it recommended for?

### Prompt Node in depth 1 (score: 0.64)

**Prompt:** Create one realistic shopper query about a product and its attributes. Make it a single short sentence or phrase (no backstory, no lists, no SKUs or code) that combines at least two attribute-related asks—include one ‘how’ or step request, introduce an ambiguity or a likely-missing detail, and ask for a source or next step if the info is unavailable. Output only the single query string and nothing else.

**Example Query:** How you measure the assembled product length? Does that include the hose and wand? Where can I find?

### Prompt Node in depth 2 (score: 0.61)

**Prompt:** Write exactly one realistic, concise shopper query (2-15 words, or one short sentence) a customer might type about a product and its attributes. The query must combine two attribute-related asks (for example, a technical/spec plus a purchase/action or a warranty/usage question) and include at least one ambiguous or likely-missing detail that forces inference; do not include backstory, lists, SKUs, code, or extra commentary. Output only the single query string and nothing else.

**Example Query:** How many amps is it, and will the assembled length fit under a low-clearance couch?

### Prompt Node in depth 3 (score: 0.45)

**Prompt:** Write one short, natural shopper question (a single brief sentence or phrase; no backstory, no lists, SKUs, code, or technical jargon) about the product and its attributes. Include at most two attribute-related asks: make one a subjective user-experience concern (comfort, fit, noise, etc.) and the other a factual or appearance-related attribute, and phrase them so the main goal is mildly ambiguous or presents a simple trade-off (for example, "I want comfort but care about looks"). You may optionally add a casual "how can I check?" but do not request numbered steps, exact standards, certificate IDs, URLs, or a seller message; if asking for a source, name only a general document type (e.g., "manufacturer spec sheet" or "warranty"). Output only the single query string and nothing else.

**Example Query:** I want a quiet vacuum but it must work on rugs. How can I check?

### Prompt Node in depth 4 (score: 0.44)

**Prompt:** Write one short, natural shopper question (a single brief sentence or phrase; no backstory, no lists, SKUs, code, or technical jargon) about the product and its attributes. Include at most two attribute-related asks: one subjective user-experience concern (comfort, fit, noise, etc.) and optionally one factual/appearance attribute, phrased as a mild trade-off (e.g., comfort vs looks). You may optionally add a casual "how can I check?" but do not request numbered steps, exact standards, certificate IDs, URLs, or seller contact; if naming a source, use only a general document type (e.g., spec sheet or warranty). Optionally include a short clause stating omits a key detail (for example, "it doesn't say X") and do NOT ask for alternatives or workaround suggestions. Output only the single query text and nothing else.

**Query:** Which surfaces is this best for? hard floors? low-pile carpet? and does it draw about 6 amps?

## D.4 Query Examples through Different Query Refinement Strategies

In this section, we show how the query refinement module modifies queries to elicit agent failures and increase query diversity. Table 6 presents examples of character-, word-, and sentence-level perturbations, while Table 7 presents examples of persona, scenario, and tone role-playing. We find that, given the robustness of current models, character-level perturbations alone are often insufficient to trigger agent failures.

## E Prompts

### E.1 Query Refinement

In the query refinement prompt through role-playing strategies, the string `[[STRATEGY_TYPE]]` is replaced with the strategy type (i.e. persona, scenario, tone).

#### Perturbation refinement

Given the current query, generate a new query by applying the specified perturbation strategy while preserving the original meaning.

Current query:  
{current\_query}

Perturbation strategy:  
{strategy}

Return:

#### Role-playing refinement

You are rewriting a query to reflect a specific `[[STRATEGY_TYPE]]`.

Original query:  
{current\_query}

`[[STRATEGY_TYPE]]`:  
{strategy\_description}

Rules:

- Do not include any identity information.
- Apply the scenario in relevant way to the original query. If the scenario is not directly related to the original query, make a connection between the scenario and the original query.
- Try to personalize the question.
- Do not repeat the exact instruction in the instruction.
- Keep the original intent.
- Maintain a natural tone and fluency.
- Be locally refined, not entirely reformulated.
- Be concise (1 - 2 sentences).

Return:

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

---

**Perturbation (Character-Level)**

---

**Strategy Description:** Add one or two extraneous characters to the end of the sentence.

---

**Example:**

**Original Query:** How absorbent are these for 8–12 hour overnight use? Is the quoted 'up to 12 hours' claim accurate, and where to verify?

**Original Response:** These diapers are designed for overnight use and are suitable for up to 12 hours of wear. The absorbency is sufficient to manage extended periods, making them appropriate for 8–12 hour overnight use. The claim of 'up to 12 hours' is based on the product's intended use as overnight diapers, which is standard for this size. For detailed specifications, check the product packaging or the manufacturer's website for technical data.

**New Query:** How absorbent are these for 8–12 hour overnight use? Is the quoted 'up to 12 hours' claim accurate, and where to verify?x

**New Response:** These diapers are designed for overnight use and offer high absorbency to support up to 12 hours of protection. The claim is based on the product's formulation for extended wear. For detailed specifications, check the product packaging or the manufacturer's website for technical data sheets.

---

**Explanation:** The query refinement adds diversity to existing queries.

---

**Perturbation (Word-Level)**

---

**Strategy Description:** Choose one or two words in the sentence that do not contribute to the meaning of the sentence and delete them.

---

**Example:**

**Original Query:** How many audio inputs and can I utilize a 3.5mm connector?

**Original Response:** The product has 2 audio inputs. It does not specify support for a 3.5mm connector.

**New Query:** How many audio inputs can I utilize a 3.5mm connector?

**New Response:** The product has 2 audio inputs.

---

**Explanation:** After removing "and", the agent only partially answers the question, ignoring the second question.

---

**Perturbation (Sentence-Level)**

---

**Strategy Description:** Change the syntactic structure of the sentence.

---

**Example:**

**Original Query:** It says 'lightweight and powerful.' Can I carry it up 2 flights of stairs, is that right?

**Original Response:** The vacuum is foldable and portable, making it easier to carry, but it's not specifically designed for carrying up stairs. Its weight is 20.34 lbs, which may make it challenging to carry up two flights of stairs.

**New Query:** When it says 'lightweight and powerful,' does that mean I can carry it up 2 flights of stairs?

**New Response:** The product's description does not provide specific information about its weight or stair-carrying capability.

---

**Explanation:** After changing the syntactic structure of the sentence, the agent failed to provide relevant product features (e.g. weight).

---

Table 6: The queries generated by applying perturbation strategies. Objective-violating response spans are highlighted in orange, while corresponding objective-compliant spans are highlighted in green.

---

**Role-playing (Persona)**

---

**Strategy Description:** You are Charles Garcia who is closed to experience, conscientious, extroverted, antagonistic, and emotionally stable. Residing in New York, you are a 32-year-old college graduate with a background in post-graduate studies. As a widowed, agnostic Black male, you understand the challenges of managing a household on a budget of 100–150, especially with an income of less than \$30,000. Despite some trouble limiting spending, you rely on rational decision-making and look for discounts, which aligns with your political stance as a Democrat favoring practical solutions. Your two previous orders focused on enhancing your home’s comfort and charm with quality products such as a winter flannel blanket and a decorative pillow cover. As a customer, your preferences are shaped by the desire for style, convenience, and affordability in family-friendly products.

---

**Example:**

**Original Query:** Worried about fit, any odor presence?

**Original Response:** These diapers are fragrance free and do not contain any harmful chemicals, ensuring a comfortable fit and no unpleasant odors.

**New Query:** I’m trying to avoid extra returns on a tight budget. Will this item fit as described and arrive without any noticeable odor?

**New Response:** Yes, this item is described as fitting as intended and does not have any noticeable odor upon arrival.

---

**Explanation:** After adding query variations, we observe that the agent is not very robust: it repeats information from the question without providing additional context for the user to make a decision.

---

**Role-playing (Scenario)**

---

**Strategy Description:** Renovating or redecorating a room

---

**Example:**

**Original Query:** how many amps and good on hardwood?

**Original Response:** The Dirt Devil Featherlite uses 12 Amps and is good on hardwood.

**New Query:** how many amps does it use, and is it good for cleaning hardwood floors during a room redecoration?

**New Response:** The Dirt Devil Featherlite uses 12 Amps and is suitable for cleaning hardwood floors during a room renovation or redecoration.

---

**Explanation:** The query refinement adds diversity to existing queries.

---

**Role-playing (Tone)**

---

**Strategy Description:** Assertive tone: An assertive tone exudes confidence and authority. It can also be insistent and straightforward. This tone can help you persuade your audience about a topic.

---

**Example:**

**Original Query:** Is this vacuum easy and comfortable to maneuver or is it heavy? How can I check?

**Original Response:** The vacuum weighs 17.35 pounds, making it moderately heavy but designed for ease of use with its bagless feature and ergonomic handle.

**New Query:** Tell me whether this vacuum is easy and comfortable to maneuver or if it’s heavy? where can I find it?

**New Response:** The vacuum’s weight is listed as 17.35 lbs. You can find this information under the product’s weight details.

---

**Explanation:** After changing the tone of the question, the agent gives fewer details.

---

Table 7: The queries generated by applying role-playing strategies. Objective-violating response spans are highlighted in orange, while corresponding objective-compliant spans are highlighted in green.

## E.2 Prompt Refinement

In the prompts, the string `[[AGENT_TYPE]]` is replaced with a description of target agent. For examples, in our experiments, "a QA agent in e-commerce domain". The string `[[DOMAIN_KNOWLEDGE]]` is replaced with the description of domain knowledge that used to generate user queries. For instance, "a product information with name and attributes".

### E.2.1 Prompt Reflection

#### Realism-related Feedback

You are an expert assistant participating in an iterative prompt optimization process for `[[AGENT_TYPE]]`. Your role is to analyze queries for judging their realism quality, then provide feedback to refine a prompt that will be used to generate future user queries. The prompt you are refining will later be combined with `[[DOMAIN_KNOWLEDGE]]` to form a complete prompt. This combined prompt will be used exclusively to generate user queries.

Below are generated queries that are unrealistic. Your task is to (1) explain why the current prompt may have produced unrealistic queries and (2) propose concrete improvements to the prompt so that future generated queries are more realistic.

Current prompt set:  
{current\_prompts}

Unrealistic queries:  
{unrealistic\_queries}

Realism definition:  
{realism\_definition}

Output Format:  
{  
"reasoning": "<why the current prompt produced unrealistic queries.>",  
"suggestions": "<give specific, actionable suggestions to refine or reword the prompt to reduce unrealistic outputs.>"  
}}

Return:

#### Objective-related Feedback (objective-violation feedback)

You are an expert assistant participating in an iterative prompt optimization process for `[[AGENT_TYPE]]`. Your role is to analyze query–response examples and their corresponding evaluation results, then provide feedback to refine a prompt that will be used to generate future user queries. The prompt you are refining will later be combined with `[[DOMAIN_KNOWLEDGE]]` to form a complete prompt. This combined prompt will be used exclusively to generate user queries.

Below, you will find generated user queries that elicited agent failures under objective criteria. Your task is to (1) analyze why the current prompt likely produced user queries that triggered these objective-violating responses and (2) propose concrete improvements to the prompt so that future generated queries are more likely to elicit objective-violating responses.

Current prompt set:  
{current\_prompts}

Query-Response pair:  
{objective\_violating\_responses}

Objective-violating criteria:  
{objective\_violating\_criteria}

Output Format:  
{  
"reasoning": "<why the current prompt produced queries that elicit objective-violating responses.>",  
"suggestions": "<give specific, actionable suggestions to refine or reword the prompt to generate user queries that are more likely to elicit objective-violating responses.>"  
}}

Return:

## Objective-related Feedback (objective-compliant feedback)

You are an expert assistant participating in an iterative prompt optimization process for `[[AGENT_TYPE]]`. Your role is to analyze query–response examples and their corresponding evaluation results, then provide feedback to refine a prompt that will be used to generate future user queries. The prompt you are refining will later be combined with `[[DOMAIN_KNOWLEDGE]]` to form a complete prompt. This combined prompt will be used exclusively to generate user queries.

Below, you will find generated user queries that elicited objective-compliant responses under objective criteria. Your task is to (1) analyze why the current prompt likely produced user queries that help satisfying agent objective, i.e., why the prompt failed to surface the agent’s weaknesses; and (2) propose concrete improvements to the prompt so that future generated queries are more likely to elicit objective-violating responses.

Current prompt set:  
{current\_prompts}

Query-Response pair:  
{objective\_compliant\_responses}

Objective-violating criteria:  
{objective\_violating\_criteria}

Objective-compliant criteria:  
{objective\_compliant\_criteria}

Output Format:  
{  
"reasoning": "<why the current prompt produced queries that help satisfying objective criteria.>",  
"suggestions": "<give contrasting actionable suggestions to refine or reword the prompt to generate user queries that are more likely to elicit objective-violating responses.>"  
}

Return:

765

### E.2.2 Prompt Expansion

766

In the prompts, the string `[[FEEDBACK_TYPE]]` is replaced with one of three objective-related feedback types: "objective-violation", "objective-compliant" and "criterion-specific".

767

768

## Prompt Expansion Template

You are a helpful failure-discovery assistant. Your goal is to refine the prompt that produces a user query: 1. Satisfying the realism policy. 2. To elicit objective-violating agent response.

DEFINITIONS:  
Objective-compliance Criteria:  
{objective\_compliant\_criteria}

Objective-violating Criteria:  
{objective\_violating\_criteria}

Realism policy:  
{realism\_policy}

You are given realism and `[[FEEDBACK_TYPE]]` feedback with the corresponding improvement suggestions for the current prompt. You must take this feedback and suggestions into account to understand how to generate more realistic user queries and how to produce queries that are more likely to elicit objective-violating responses.

The refined prompts should enable the generator to produce a query using a variety of challenging strategies. These strategies may include (but are not limited to): (1) introducing ambiguity, (2) obscuring the intended task or goal, (3) incorporating confusing, inconsistent, or ill-structured logic, and/or (4) using unconventional or creative approaches that the agent may not anticipate. These strategies are not exhaustive—develop additional ones as needed. If a strategy becomes ineffective or has been used repeatedly, introduce variation or switch to a different strategy to maintain diversity and challenge the agent reliably.

Current prompt:  
{current\_prompt}

769

```

Reasoning of why queries generated by the current prompt are unrealistic:
{unrealistic_reasoning} # The reasoning output from realism-related feedback.

Suggestions for generating more realistic queries:
{realistic_improvement_suggestions} # The suggestion output from realism-related feedback.

# if [[FEEDBACK_TYPE]] == "objective-violation feedback":
Reasoning of why queries generated by the current prompt successfully elicit agent weaknesses:
{objective_violating_reasoning} # The reasoning output from objective-violation feedback.

Suggestions for generating queries that elicit agent weaknesses:
{objective_violating_improvement_suggestions} # The suggestion output from objective-violation feedback.

# elif [[FEEDBACK_TYPE]] == "objective-compliant feedback":
Reasoning of why queries generated by the current prompt failed to reveal agent weaknesses:
{objective_compliant_reasoning} # The reasoning output from objective-compliant feedback.

Suggestions for generating queries that elicit agent weaknesses:
{objective_compliant_improvement_suggestions} # The suggestion output from objective-compliant feedback.

# elif [[FEEDBACK_TYPE]] == "criterion-specific feedback":
Target criterion to break: {target_criterion} # The criterion get based on the evaluator's reward signal.

Previous prompts (Do NOT reuse or repeat any of these, and do NOT use any strategies contained in these prompts):
{prev_prompts}

Guidelines for the new prompt:
- Avoid copying phrases from prior prompts.
- Do not repeat any content listed under "Previous prompts".
- The new prompt should remain general and not rely on specific examples from the feedback.
- The difference between the current prompt and the refined prompt should be limited to at most 1-3 sentences in length.

OUTPUT FORMAT:
{{
"reasoning": "<Describe how the prompt should be adjusted to generate realistic query that is more likely to elicit
objective-violating agent responses.>",
"prompt": "<provide a self-contained, actionable prompt suitable for generating query.>"
}}

Return:

```

770

771

### E.3 LLM-as-a-Judge

#### Response Evaluation

You are an evaluator assessing whether a given response satisfies specific criteria with respect to a query. Read the query and the response carefully, then evaluate the response against the provided criteria. Base your judgment strictly on the given query and response; do not assume any unstated context. Identify all criteria that the response satisfies, and explain your reasoning clearly and concisely.

Query:  
{query}

Response:  
{response}

Criteria:  
{criteria}

Rules:  
- Judge each criterion independently. Select all criteria that the response clearly satisfies.  
- Select only criteria that are clearly supported by evidence in the response.  
- If the response does not satisfy any criteria, leave the set empty.

Output format:  
First, provide the reasoning for each selected criterion. Then, list the keys of the criteria that apply to the response.

Return:

772