## The Road Not Taken: Hindsight Exploration for LLMs in Multi-Turn RL

Anonymous Author(s) Affiliation Address email

#### Abstract

Multi-turn reinforcement learning provides a principled framework for training 1 2 LLM agents, but exploration remains a key bottleneck. Classical exploration 3 strategies such as  $\epsilon$ -greedy and upper confidence bounds select random actions, failing to efficiently explore the combinatorial space of multi-turn token sequences. 4 Our key insight is that LLMs can use hindsight to guide exploration: by analyzing 5 completed trajectories and proposing counterfactual actions that could have led 6 to higher returns. We propose HOPE (Hindsight Off-Policy Exploration), which 7 integrates hindsight-guided exploration into both the actor and critic stages of multi-8 turn RL. HOPE improves the critic's state-action coverage by generating rollouts 9 from counterfactual actions, and steers the actor's exploration in RL by using 10 a learned counterfactual generator to propose alternative actions. Experimental 11 results show that HOPE outperforms strong multi-turn RL baselines in task-oriented 12 dialogue tasks, TwentyQuestions (success:  $0.82 \rightarrow 0.97$ ), GuessMyCity (success: 13  $0.68 \rightarrow 0.75$ ), and tool-use dialogue task CarDealer (success:  $0.72 \rightarrow 0.77$ ). 14

#### **15 1** Introduction

Reinforcement learning (RL) has proven effective for aligning large language models (LLMs)
with human preferences [31, 36, 41] and enhancing their reasoning capabilities [13, 24, 45, 48].
Recent works apply RL to train LLM agents for multi-turn decision-making tasks, such as code
generation [18, 23], web navigation [4, 32, 33], and task-oriented dialogues [1, 46, 55], where each
action by the agent influences future states in the environment.

However, exploration remains a key challenge in RL [6, 42, 44]. Classical exploration strategies that sample random actions, such as  $\epsilon$ -greedy and upper confidence bound [3], fail to efficiently explore the expansive action space of LLM agents, where actions are token sequences spanning multiple turns. Random perturbations rarely produce coherent alternative actions, making it difficult to explore the space of multi-turn sequences in a purposeful way.

Our key insight is that LLMs can guide exploration through hindsight reasoning—by analyzing
 completed trajectories and proposing counterfactual actions that might have led to higher returns.
 Prior works empirically show that LLMs exhibit this form of hindsight reasoning [5, 19, 29, 39],
 often revising their responses based on feedback. However, naively imitating these counterfactual
 actions does not guarantee policy improvement—they may be suboptimal or unrealizable.
 We present HOPE (Hindsight Off-Policy Exploration), which leverages environment rewards and a

learned critic to identify the counterfactual actions that have resulted in improved outcomes. HOPE
 integrates hindsight-guided exploration into both stages of actor-critic RL. In critic training, it

<sup>34</sup> augments the replay buffer by splicing in high-value counterfactual actions at intermediate timesteps

<sup>35</sup> of past trajectories and then rolling out the current policy, thereby broadening state–action coverage.

- 36 In actor updates, it biases exploration toward promising regions by sampling candidate actions from
- a learned counterfactual generator and evaluating them via the critic. We show that HOPE is both
- easy to incorporate into existing multi-turn RL pipelines and closely aligned with the principles of
- <sup>39</sup> posterior sampling, as it explores by generating actions conditioned on past outcomes.
- 40 Our key contributions are:
- A novel framework, HOPE, that injects hindsight-guided exploration into both critic and actor
   stages of multi-turn RL.
- 43 2. A counterfactual action generator that steers the actor's exploration during RL by proposing44 alternative actions.
- 45 3. Experiments validation on diverse multi-turn decision-making domains [1], showing that HOPE
- 46 outperforms leading multi-turn RL approaches on task-oriented dialogue and tool-use tasks.



Figure 1: **Overview of** HOPE that integrates hindsight-guided exploration in actor-critic RL training. The training loop is iterative. Stage 1: Collect diverse rollouts to compute Q targets by augmenting past rollouts with counterfactual actions and rolling out the remaining steps with the current policy  $\pi_i$ . Stage 2: Train the critic  $Q_i(s, a)$  via supervised learning. Stage 3: Train the policy via RL and hindsight guided action exploration, where it samples counterfactual actions with probability  $\beta$ .

#### 47 **2 Preliminaries**

48 **Problem Formulation.** We model the multi-turn decision-making problem as a Markov Decision 49 Process (MDP) for the LLM agent. At each turn t, the agent receives state  $s_t = \{o_0, a_0, \dots, o_t\}$ , the 50 history of observations and actions so far. It then takes an action  $a_t$ , a token sequence, and transitions 51 to a new state  $s_{t+1}$  according to the environment's transition dynamics. We assume access to a sparse 52 reward signal  $r(s_t, a_t)$  only at the end of the trajectory. The goal is to learn a policy  $\pi(a_t \mid s_t)$  that 53 maximizes the expected discounted return:  $\mathbb{E}_{\pi} \left[ \sum_{t=0}^{T-1} \gamma^t r(s_t, a_t) \right]$ .

Actor Critic for Multi-Turn RL. Multi-turn RL can be viewed as a hierarchical RL problem consisting of two nested RL loops: (1) turn-level RL and (2) token-level RL. We solve the turn-level RL problem by rolling out the agent in the environment, collecting state-action trajectories and training a Process Reward Model (PRM) Q(s, a), akin to a critic in actor-critic RL. The LLM policy is then solved by optimizing the PRM using a token-level RL algorithm. While this recipe has been used to great success in math and reasoning settings [24, 38, 45, 48], this framework is underexplored in the LLM agent setting.

<sup>61</sup> The training loop is iterative: at each iteration *i*, the current policy  $\pi_{i-1}$  is rolled out to generate <sup>62</sup> trajectories. For each encountered state-action pair (s, a), all trajectories passing through (s, a) are <sup>63</sup> stored in a replay buffer  $\mathcal{G}(s, a)$ . The PRM target dataset  $\mathcal{D} = \{(s, a, \hat{Q})\}$  is constructed using Monte <sup>64</sup> Carlo estimates:

$$\hat{Q}(s,a) = \frac{1}{|\mathcal{G}(s,a)|} \sum_{(s_t,a_t)\in\mathcal{G}(s,a)} \sum_{k=t}^{T-1} \gamma^{k-t} r(s_k,a_k).$$
(1)

- The PRM  $Q_i$  is trained on  $\mathcal{D}$ , and the policy  $\pi_i$  is then trained to maximize  $Q_i$  while staying close to 65  $\pi_{i-1}$ , using algorithms such as PPO [37], rejection sampling [15], or Online DPO [16]. 66
- The effectiveness of this framework hinges on the quality and diversity of trajectories in the replay 67
- buffer. When exploration is limited, the buffer contains few high-reward state-action pairs, leading to 68
- poor value estimates and stagnated policy improvement. 69 70
  - Algorithm 1 Actor-Critic RL with HOPE (Hindsight Off-Policy Exploration)
  - 1: Initialize with agent policy  $\pi_0$
  - 2: for iteration  $i = 1, \ldots, K$  do
  - // Stage 1: Collect rollouts (Explore with HOPE) 3:
  - 4: Collect on-policy rollouts  $\{(\ldots, s_t, a_t, r_t, \ldots)\}$  using  $\pi_{i-1}$  and store in replay buffer  $\mathcal{G}(s, a)$
  - 5: Collect hindsight-guided rollouts Algorithm 2 and maintain data ratio  $\alpha$  in  $\mathcal{G}(s, a)$
  - 6: // Stage 2: Train Process Reward Model
  - Compute PRM targets via (1) and aggregate into dataset  $\mathcal{D} = \{(s, a, \hat{Q})\}$ 7:
  - Train PRM  $Q_i$  to minimize soft binary cross-entropy loss: 8:

$$Q_i = \arg\min_{Q} -\mathbb{E}_{(s,a,\hat{Q})\sim\mathcal{D}}\left[\hat{Q}\log Q(s,a) + (1-\hat{Q})\log(1-Q(s,a))\right]$$
(2)

- 9: // Stage 3: Train Policy via RL (Explore with HOPE)
- Train policy  $\pi_i$  to maximize  $Q_i$ , exploring with probability  $\beta$  via hindsight proposer  $\pi^H$ 10:

$$\pi_{i} = \arg\max_{\sigma} \mathbb{E}_{s \sim \mathcal{D}, a \sim \text{Sample}(s, \pi_{i}, \pi^{H}, \beta)} \left[ Q_{i}(s, a) \right] - \beta \mathcal{D}_{\text{KL}} \left[ \pi(a \mid s) \| \pi_{i-1}(a \mid s) \right]$$
(3)

11: end for

12: **return** Best  $\pi \in {\pi_1, \ldots, \pi_K}$  on validation dataset

#### Approach 3 71

We introduce HOPE (Hindsight Off-Policy 72 Exploration) in Algorithm 1, a framework that 73 leverages LLMs' hindsight reasoning to guide 74 exploration in multi-turn RL. Unlike conven-75 tional approaches that rely on random action 76 noise, HOPE proposes counterfactual actions by 77 analyzing completed trajectories. A hindsight 78 proposer generates outcome-conditioned alter-79 natives from trajectory summaries (Section 3.1), 80 which are then used to improve state-action 81 82 coverage during critic training (Section 3.2) and steer exploration during actor training (Sec-83 tion 3.3). 84

#### Algorithm 2 HOPE Data Collection For Critic

- 1: **Input:** Dataset of onpolicy rollouts  $\{\tau\}$ ; current policy  $\pi$ ; rollout summarizer  $\pi^{\text{sum}}(\phi \mid \tau)$ ; hindsight proposer  $\pi^H(a^H \mid s_t, \phi)$ .
- 2: for each onpolicy rollout  $\tau$  do
- Get summary  $\phi \sim \pi^{\text{sum}}(\cdot \mid \tau)$ 3:
- 4:
- 5:
- for each randomly sampled timestep t do Get counterfactual  $a_t^H \sim \pi^H(\cdot \mid s_t, \phi)$ Create partial rollout  $\tau^H = (\dots, s_t, a_t^H)$ 6:
- Complete  $\tau^H$  with  $\pi$  and store in dictio-7: nary  $\mathcal{G}^H(s,a)$

#### end for 8: 9: end for

10: **return** State-action dictionary  $\mathcal{G}^{H}(s, a)$ 

#### 3.1 Hindsight Proposer 85

We introduce a *hindsight proposer*,  $\pi^{H}$ , a LLM 86

- policy that generates a counterfactual action given a state  $s_t$  and a completed trajectory au = 87  $(s_0, a_0, r_0, \ldots)$ . The goal is to produce an alternative action that might have led to a better out-88
- 89 come.
- A key challenge is that completed trajectories are long and entangled: they contain many state-action 90 pairs whose contributions to final reward are hard to disentangle. Prompting an LLM directly on  $\tau$ 91 often yields shallow or noisy revisions. To address this, we introduce a two-step decomposition. First, 92 we use a summarizer LLM  $\pi^{\text{sum}}$  to generate a compact, natural language summary  $\phi \sim \pi^{\text{sum}}(\cdot \mid \tau)$ 93 that reflects the policy's high-level strategy and assesses the effectiveness of individual actions. This 94 summary typically includes both behavioral intent (e.g., "the agent fails to consider a common 95 category...") and retrospective feedback ("this prevents the agent from searching effectively..."). Then, 96

conditioned on the state  $s_t$  and summary  $\phi$ , the hindsight proposer  $\pi^H$  generates a counterfactual action  $a_t^H \sim \pi^H(\cdot | s_t, \phi)$ .

Prior work has demonstrated that LLMs are capable of hindsight reasoning—either by refining their responses based on feedback [5, 19, 29, 39] or generating post-hoc explanations conditioned on ground-truth answers [47, 50, 52]. However, these capabilities have primarily been applied to improve generation quality or interpretability. In contrast, we use hindsight reasoning to guide exploration: the LLM summarizes completed trajectories and proposes counterfactual actions that could have led to better outcomes. This allows the agent to incorporate structured exploratory data during training, rather than relying solely on stochastic or undirected exploration strategies.

#### 106 3.2 Hindsight-Guided Critic Training

Algorithm 2 presents hindsight-guided data collection to expand state-action coverage in the critic training data. We first collect a small set of on-policy trajectories  $\{\tau\}$  with the current policy  $\pi_{i-1}$ . We summarize each trajectory  $\tau$  as  $\phi$  and randomly select a state  $s_t$  to query the hindsight proposer for a counterfactual action  $a_t^H \sim \pi^H(\cdot \mid s_t, \phi)$ . To evaluate whether  $a_t^H$  can improve outcomes, we augment the original trajectory with  $a_t^H$  before completing the rollout with  $\pi_{i-1}$  and obtaining  $\tau^H = (\dots, s_t, a_t^H, r_t^H, s_{t+1}, a_{t+1}, \dots)$ .

Then, we store all collected trajectories that pass through each encountered state-action pair (s, a) in a replay buffer  $\mathcal{G}(s, a)$  before computing target Q-values via Monte Carlo Estimation. To account for off-policy data introduced by the hindsight proposer, we construct  $\mathcal{G}(s, a)$  using a mixture that consists  $\alpha$  percentage of hindsight data and  $(1 - \alpha)$  data from on-policy trajectories  $\{\tau\}$ . We empirically investigate the effect of ratio  $\alpha$  in Section D.1.

#### 118 3.3 Hindsight-Guided Actor Training

To guide action exploration during RL, HOPE samples counterfactual actions from the hindsight proposer  $a^H \sim \pi^H(s, \phi)$  with probability  $\beta$  given a state s and a corresponding summary  $\phi$ . Setting  $\beta = 0$  recovers the original actor objective where actions are only sampled from the LLM agent. The critic Q(s, a) provides supervision on the counterfactual actions, allowing the actor to learn from actions that lead to higher Q-estimates. Because counterfactual actions are explicitly generated to differ from collected experiences, they help overcome premature convergence when actions sampled from the policy are clustered around local optima [8, 49, 54].

#### 126 3.4 Connection to Posterior Sampling

127 Our method draws inspiration from posterior sampling (PS) which is a well-studied exploration strategy in multi-armed bandit [9, 43] and RL [17, 42]. In the context of RL, the algorithm maintains 128 and samples from a posterior over MDPs M given the history of all trajectories collected so far: 129  $P(M|\mathcal{H})$ . In every episode i, it samples a new MDP  $M_i \sim P(M \mid \mathcal{H})$ , computes the optimal 130 policy  $\pi_i$ , rolls it out to collect interaction data and adds it to  $\mathcal{H}$ , and updates the posterior. PS is 131 provably exploration-efficient in terms of Bayesian regret, achieving near-optimal bounds in tabular 132 settings [17]. The key intuition is that posterior sampling balances exploration and exploitation 133 by implicitly encouraging optimism through sampling: sampled MDPs often favor under-explored 134 regions of the environment due to posterior uncertainty. While elegant in theory, extending PS to 135 large-scale or function-approximation settings introduces practical challenges, such as maintaining a 136 posterior over high-dimensional model classes and solving for optimal policies in sampled MDPs at 137 each episode. 138

In contrast to classical PS, we make a number of practical approximations. Rather than explicitly 139 maintaining and updating a full posterior over MDPs, we instead generate counterfactual actions 140 from a hindsight policy  $\pi^{H}(s, \phi)$  conditioned on a trajectory summary  $\phi$ . This hindsight generator 141 implicitly captures posterior over plausible MDPs in its reasoning, and then samples alternative 142 actions that would be optimal under the MDP. Unlike PS, which updates its posterior using the 143 full interaction history  $\mathcal{H}$ , our approach uses only a compact trajectory-level summary  $\phi$ . This 144 simplification sacrifices the formal guarantees of posterior sampling, since the hindsight policy may 145 not reflect a true posterior update. However, when the LLM prior is well-calibrated with the MDP 146 distribution, this approximation can yield effective exploration without the computational burden of 147 maintaining and updating an explicit posterior. 148



		TwentyQuestions		Guess	MyCity	CarDealer (Tool)		
		Return ↑	Success ↑	Return ↑	Success ↑	Return ↑	Success ↑	
	$\frac{\text{gpt4o}}{3B}$ $\pi_0$	1.00 (0.11) 0.09 (0.04) 1.08 (0.11)	0.61 (0.06) 0.28 (0.05) 0.67 (0.05)	<b>1.00 (0.12)</b> 0.01 (0.01) 0.80 (0.09)	0.56 (0.05) 0.01 (0.01) 0.70 (0.05)	1.00 (0.09) 0.40 (0.07) 1.03 (0.08)	0.59 (0.05) 0.51 (0.05) 0.66 (0.05)	
LEAP	$\pi_1 \\ \pi_2$	1.16 (0.12) 0.52 (0.09)	0.68 (0.06) 0.45 (0.06)	0.35 (0.04) 0.29 (0.03)	0.59 (0.05) 0.62 (0.05)	0.77 (0.08) 0.86 (0.08)	0.52 (0.05) 0.57 (0.05)	
Reject-S	$\pi_1 \ \pi_2$	0.63 (0.11) 0.90 (0.13)	0.55 (0.06) 0.63 (0.06)	0.46 (0.02) 0.31 (0.03)	0.69 (0.02) 0.62 (0.05)	1.07 (0.08) 1.00 (0.08)	0.68 (0.05) 0.62 (0.05)	
PRM+RL	$\pi_1 \ \pi_2$	1.19 (0.16) 1.05 (0.09)	0.55 (0.06) 0.82 (0.05)	0.35 (0.05) 0.54 (0.03)	0.66 (0.05) 0.66 (0.02)	1.00 (0.08) 1.14 (0.08)	0.67 (0.05) 0.72 (0.04)	
HOPE	$rac{\pi_1}{\pi_2}$	1.15 (0.13) <b>1.91 (0.09)</b>	0.77 (0.05) <b>0.97 (0.02</b> )	0.45 (0.05) 0.91 (0.10)	0.74 (0.05) <b>0.77 (0.04)</b>	1.07 (0.08) <b>1.20 (0.08)</b>	0.69 (0.05) <b>0.77 (0.04</b> )	

Table 1: **Policy performance on test sets of three task-oriented conversational environments.** We report the average normalized return and success rate with standard error, and we highlight the top-performing approaches. The return is normalized with respect to GPT-4o's performance.

### 149 **4** Experiments

Our experiments aim to evaluate the effectiveness of HOPE in multi-turn conversational domains with sparse environment rewards. Below, we summarize the key research questions we investigated along with our main findings:

153 1. Does HOPE outperform state-of-the-art IL and RL algorithms? (Section 4.2) As shown in 154 Table 1, HOPE achieves the highest success rates, improving over baselines across all domains, 155 with gains in TwentyQuestions ( $0.82 \rightarrow 0.97$ ), GuessMyCity ( $0.68 \rightarrow 0.75$ ), and CarDealer 156 ( $0.72 \rightarrow 0.77$ ).

How effective is hindsight-guided data collection for critic training? (Section 4.3) Figure 2
 shows that HOPE achieves a success rate of 0.77 on TwentyQuestions, outperforming both the
 unguided H-Temp (0.55) and Explore-Prompt (0.72) exploration strategies.

3. Does hindsight-guided exploration for RL training improve performance? (Section 4.4) Table 2 shows that hindsight-guided exploration improves success rates in the TwentyQuestions environment, especially when the critic data collection is unguided  $(0.55 \rightarrow 0.62)$ .

#### 163 4.1 Setup

Multi-Turn Environments. We evaluate our approach across 3 multi-turn conversational environments proposed in the LMRL Gym benchmark [1]. Each environment models a conversation, where the LLM agent must talk with a simulated user to achieve some goal. Below, we describe the conversation structure, reward function, success criteria, and dataset splits for each environment:

• **TwentyQuestions.** In this environment, the LLM agent must identify a secret object selected by the simulator. To do so, the agent has up to 20 turns, each consisting of a single yes-or-no question aimed at narrowing down the possibilities. The answers are simulated by prompting a Llama-3.2-3B-Instruct model [15] with the secret object name and the question. The LLM agent receives a -1 reward on each turn, and the episode ends with a reward of 0 when the secret object is
correctly guessed. The train and validation set share the same object categories but contain different
objects. The test set contains new objects from categories not seen in the train/validation sets.
GuessMyCity. As in the previous environment, the simulator selects a secret city that the LLM
agent must identify. However, instead of asking yes-or-no questions, the agent can ask up to 10
open-ended questions, each eliciting a free-form response. We use Llama-3.2-3B-Instruct [15] to
simulate the environment. The episode ends successfully with reward of 0 when the agent correctly

- guesses the secret city, accruing a reward of -1 for all previous turns. The training and validation sets include cities from the same set of countries, while the test set uses an unseen set of countries.
  CarDealer (Tool). We extend the original Car Dealer environment to incorporate tool use. In this environment, the LLM agent plays the role of a car dealer aiming to sell a vehicle to a buyer within
- 10 dialogue turns. At each turn, the agent may first issue a database tool call to search the inventory and then respond to the buyer. The interaction ends when the buyer agrees to a purchase, with a reward of (purchase\_cost)<sup>2</sup>/(budget × market\_price); all prior turns yield zero reward. We simulate buyers with diverse personalities and constraints using Qwen2.5-14B-Instruct [35]. The validation set models the same buyer profiles as the training set, while varying the available cars.
- 188 The test set has both unseen buyer profiles and cars.

**Baselines.** We evaluate the effectiveness of HOPE against a range of baseline approaches. We begin by 189 benchmarking the zero-shot performance of GPT-40 and the base model, Llama-3.2-3B-Instruct [15] 190 (denoted as **3B**). Following RLHF [31], we supervised finetune the base model for 3 epochs using 191 gpt4o's rollouts on the train set and denote this policy as  $\pi_0$ . Using  $\pi_0$  as the starting model for 192 193 training, we compare various imitation learning (IL) and reinforcement learning (RL) methods. LEAP [11] is an IL approach that directly finetunes the policy to imitate corrective actions from a 194 teacher policy, which is the hindsight proposer for our setting (Sec. 3.1). For RL, we compare with 195 Rejection Sampling (Reject-S for short) [15] that skips training a critic. It iteratively finetunes the 196 policy on successful trajectories as indicated by the environment sparse reward. Meanwhile, PRM+RL 197 is an actor-critic method that explicitly trains a critic similar to HOPE, but it does not leverage the 198 hindsight proposer for its critic's data collection (i.e,  $\alpha = 0$ ) and simply samples actions from the 199 current policy with high temperature 1.0. 200

**Training Details**. Appendix includes detailed information about the ratio  $\alpha$  of hindsight data used to train HOPE's critic in each domain, prompts, and training hyperparameters,

Metrics. We report the success rate and normalized returns (i.e., cumulative rewards normalized with respect to GPT-40's performance) on the test set. We evaluate all intermediate training checkpoints on the validation set and select the best-performing one for final evaluation on the test set.

#### 206 4.2 Does HOPE outperform state-of-the-art IL and RL algorithms?

Table 1 shows that HOPE consistently outperforms 207 all baselines, achieving the highest rewards and suc-208 cess rates across all domains. Among the RL base-209 lines, Reject-S often fails to improve over the ini-210 211 tial policy  $\pi_0$ , particularly in the TwentyQuestions and CarDealer environments. This is because it only 212 updates the policy using state-action pairs from suc-213 cessful trajectories, ignoring potentially informative 214 rollouts that result in failure. In contrast, PRM+RL is 215 more competitive, as its actor is trained to maximize 216 Q-values estimated by a critic, enabling it to learn 217 from both successful and unsuccessful trajectories. 218 However, since its critic is trained only on on-policy 219 data, it suffers from limited state-action coverage. 220 This limits the actor's ability to discover higher-value 221 behaviors. HOPE overcomes this limitation through 222 hindsight-guided exploration, which augments the 223 critic training data with diverse, high-value counter-224 factual actions and improves critic supervision. Fi-225



Figure 2: Policy performance given different exploration strategies. We use the test set of TwentyQuestions to evaluate  $\pi_1$  trained with different exploration approaches and report the average success with standard error.

nally, the IL baseline LEAP often underperforms, as it lacks access to reward signals and simply imitates the hindsight proposer, which may generate suboptimal actions in certain contexts.

#### **4.3** How effective is hindsight-guided data collection for critic training?

We evaluate the effectiveness of hindsight-guided exploration against both unguided and guided exploration strategies. For each method, we train a policy  $\pi_1$  using actor-critic RL and report its test-time performance on the TwentyQuestions domain.

As an unguided baseline, H-Temp (PRM+RL) explores by sampling actions from  $\pi_0$  with a high 232 temperature of 1.0, collecting 30 rollouts per task to train the critic. We also compare against two 233 guided strategies, Explore-Prompt and Oracle-Explore, both of which follow the HOPE protocol 234 in Algorithm 2 by collecting 14 on-policy rollouts with  $\pi_0$ , followed by 16 guided rollouts. Instead of 235 using a hindsight proposer, Explore-Prompt samples 5 actions from  $\pi_0$  (at temperature 1.0) and then 236 augments the prompt to explicitly request actions that differ from those samples. Oracle-Explore, 237 in contrast, uses an oracle proposer—specifically, the final HOPE policy  $\pi_2$ —to generate actions, 238 serving as an upper bound for guided exploration. 239

#### **4.4 Does hindsight-guided exploration for RL training improve performance?**

We investigate how hindsight-guided exploration influences actor training by varying the hyperparameter  $\beta$  that controls the probability of sampling counterfactual actions from the hindsight generator.

**Setup.** We compare four configurations: ( $\alpha = 0.0, \beta = 0.0$ ), which does not utilize any hindsight exploration; ( $\alpha = 0.0, \beta = 0.5$ ), which only uses hindsight-guided exploration during actor training; ( $\alpha =$  $0.4, \beta = 0.0$ ), which only uses

Hindsight	Critic	Х	Х	$\checkmark$	$\checkmark$	
Explore in	Actor	Х	$\checkmark$	Х	$\checkmark$	
success		0.55 (0.06)	0.62 (0.06)	0.77 (0.05)	0.77 (0.05)	

Table 2: Effect of hindsight-guided exploration during actor training.  $\checkmark$  in Critic represents  $\alpha = 0.4$  where the critic is trained with 40% hindsight data.  $\checkmark$  in Actor represents  $\beta = 0.5$ . We report the average success rate and standard error in TwentyQuestions test set.

hindsight-guided critic trained on 40% hindsight data; ( $\alpha = 0.4, \beta = 0.5$ ), which uses hindsight exploration both in critic and actor training. To cost-efficiently generate counterfactual actions during RL training, we distill the GPT-40 hindsight proposer  $\pi^H$  into a Llama-3.2-3B-Instruct model via supervised fine-tuning for three epochs on data  $(s_t, \phi, a_t^H) \sim \pi^H(\cdot | s_t, \phi)$ .

**Results.** Table 2 reports the average success rate of the policy trained with these configurations on the TwentyQuestions test set. We observe that hindsight-guided exploration during critic training contributes to the most significant performance gain  $(0.55 \rightarrow 0.77)$ . When a critic is ineffective  $(\alpha = 0)$ , guided exploration during actor training can boost performance  $(0.55 \rightarrow 0.62)$ . When the critic is trained with a sufficient state-action coverage, additional exploration in actor training does not yield additional gains.

#### 264 **5** Discussion

We present HOPE, a framework that leverages hindsight reasoning in LLMs to guide exploration in multi-turn reinforcement learning. Instead of directly imitating counterfactual actions, HOPE incorporates them into both actor and critic training, enabling iterative policy improvement through a simple and scalable approach. By combining environment rewards with a learned critic, the method identifies beneficial counterfactuals and uses them to shape future behavior.

270 Experiments show that HOPE consistently outperforms imitation learning and standard RL exploration

strategies across domains. Nonetheless, two limitations remain: (1) HOPE inserts a single counterfactual per trajectory before reverting to on-policy rollouts, as full counterfactual rollouts require

expensive environment simulation; and (2) evaluations were limited to LLMs with up to three billion

274 parameters, leaving scalability to larger models an open question.

#### 275 **References**

- [1] 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, 2023.
- [2] Dilip Arumugam and Thomas L. Griffiths. Toward efficient exploration by large language model agents, 2025.
- [3] Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Mach. Learn.*, 47(2–3):235–256, May 2002.
- [4] Hao Bai, Yifei Zhou, Mert Cemri, Jiayi Pan, Alane Suhr, Sergey Levine, and Aviral Kumar.
   Digirl: Training in-the-wild device-control agents with autonomous reinforcement learning. In
   A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors,
   *Advances in Neural Information Processing Systems*, volume 37, pages 12461–12495. Curran
   Associates, Inc., 2024.
- [5] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, 288 Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine 289 Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli 290 Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal 291 Ndousse, Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, 292 Noemi Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, 293 Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, 294 Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben 295 Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. 296 Constitutional ai: Harmlessness from ai feedback, 2022. 297
- [6] Ronen I Brafman and Moshe Tennenholtz. R-max-a general polynomial time algorithm for
   near-optimal reinforcement learning. *Journal of Machine Learning Research*, 3(Oct):213–231,
   2002.
- [7] Thomas Carta, Clément Romac, Thomas Wolf, Sylvain Lamprier, Olivier Sigaud, and Pierre Yves Oudeyer. Grounding large language models in interactive environments with online
   reinforcement learning. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara En gelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*,
   pages 3676–3713. PMLR, 23–29 Jul 2023.
- [8] Shicong Cen, Jincheng Mei, Katayoon Goshvadi, Hanjun Dai, Tong Yang, Sherry Yang, Dale
   Schuurmans, Yuejie Chi, and Bo Dai. Value-incentivized preference optimization: A unified
   approach to online and offline RLHF. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [9] Olivier Chapelle and Lihong Li. An empirical evaluation of thompson sampling. In J. Shawe Taylor, R. Zemel, P. Bartlett, F. Pereira, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc., 2011.
- [10] Baian Chen, Chang Shu, Ehsan Shareghi, Nigel Collier, Karthik Narasimhan, and Shunyu Yao.
   Fireact: Toward language agent fine-tuning, 2023.
- [11] Sanjiban Choudhury and Paloma Sodhi. Better than your teacher: LLM agents that learn from
   privileged AI feedback. In *The Thirteenth International Conference on Learning Representa- tions*, 2025.
- [12] Julian Coda-Forno, Marcel Binz, Zeynep Akata, Matt Botvinick, Jane Wang, and Eric Schulz.
   Meta-in-context learning in large language models. In A. Oh, T. Naumann, A. Globerson,
   K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 65189–65201. Curran Associates, Inc., 2023.

[13] DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin 323 Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, 324 Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan 325 Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, 326 Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli 327 Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng 328 Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, 329 Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian 330 Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean 331 Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan 332 Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, 333 Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong 334 Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan 335 Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting 336 Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, 337 T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, 338 Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao 339 Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, 340 Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang 341 Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. 342 Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao 343 Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang 344 Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, 345 Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong 346 Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, 347 Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan 348 Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, 349 Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, 350 and Zhen Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement 351 learning, 2025. 352

- [14] Vikranth Dwaracherla, Seyed Mohammad Asghari, Botao Hao, and Benjamin Van Roy. Efficient exploration for LLMs. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp, editors, *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 12215–12227. PMLR, 21–27 Jul 2024.
- [15] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ah-358 mad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela 359 Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem 360 Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, 361 Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, 362 Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, 363 Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, 364 365 Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab 366 AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco 367 Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind 368 Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah 369 Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan 370 Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason 371 Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya 372 Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, 373 Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Va-374 suden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, 375 Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal 376 Lakhotia, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz 377 Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke 378 de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin 379 Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kam-380

badur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, 381 Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, 382 Onur Celebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal 383 Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao 384 Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert 385 Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, 386 Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hos-387 seini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, 388 Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, 389 Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane 390 Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, 391 Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal 392 Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, 393 Virginie Do, Vish Vogeti, Vítor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin 394 Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, 395 Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine 396 Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, 397 Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, 398 Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay 399 Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit 400 Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, 401 Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, 402 Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, 403 Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, 404 Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, 405 Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, 406 Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester 407 Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon 408 Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, 409 Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin 410 Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, 411 Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, 412 Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank 413 Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, 414 Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan 415 Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison 416 Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, 417 418 Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff 419 Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, 420 Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh 421 Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun 422 Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, 423 Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro 424 Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, 425 Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew 426 Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao 427 Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel 428 Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, 429 Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, 430 Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich 431 Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem 432 Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, 433 Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, 434 Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, 435 Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ 436 Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, 437 Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, 438 Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao 439

Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, 440 Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen 441 Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, 442 Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, 443 Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim 444 Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, 445 Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu 446 Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Con-447 stable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, 448 Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin 449 Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary 450 DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The llama 3 451 herd of models, 2024. 452

[16] Shangmin Guo, Biao Zhang, Tianlin Liu, Tianqi Liu, Misha Khalman, Felipe Llinares, Alexan dre Rame, Thomas Mesnard, Yao Zhao, Bilal Piot, Johan Ferret, and Mathieu Blondel. Direct
 language model alignment from online ai feedback, 2024.

[17] Osband Ian, Van Roy Benjamin, and Russo Daniel. (more) efficient reinforcement learning via
 posterior sampling. In *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2*, NIPS'13, page 3003–3011, Red Hook, NY, USA, 2013. Curran
 Associates Inc.

[18] Arnav Kumar Jain, Gonzalo Gonzalez-Pumariega, Wayne Chen, Alexander M Rush, Wenting
 Zhao, and Sanjiban Choudhury. Multi-turn code generation through single-step rewards. In
 *Workshop on Reasoning and Planning for Large Language Models*, 2025.

<sup>463</sup> [19] Geunwoo Kim, Pierre Baldi, and Stephen Marcus McAleer. Language models can solve <sup>464</sup> computer tasks. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.

[20] Akshay Krishnamurthy, Keegan Harris, Dylan J. Foster, Cyril Zhang, and Aleksandrs Slivkins.
 Can large language models explore in-context?, 2024.

[21] Nicholas Kroeger, Dan Ley, Satyapriya Krishna, Chirag Agarwal, and Himabindu Lakkaraju.
 Are large language models post hoc explainers? In *R0-FoMo:Robustness of Few-shot and Zero-shot Learning in Large Foundation Models*, 2023.

470 [22] Aviral Kumar, Vincent Zhuang, Rishabh Agarwal, Yi Su, John D Co-Reyes, Avi Singh, Kate
471 Baumli, Shariq Iqbal, Colton Bishop, Rebecca Roelofs, Lei M Zhang, Kay McKinney, Disha
472 Shrivastava, Cosmin Paduraru, George Tucker, Doina Precup, Feryal Behbahani, and Aleksandra
473 Faust. Training language models to self-correct via reinforcement learning, 2024.

Hung Le, Yue Wang, Akhilesh Deepak Gotmare, Silvio Savarese, and Steven Chu Hong Hoi.
Coderl: Mastering code generation through pretrained models and deep reinforcement learning.
In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 21314–21328. Curran Associates,
Inc., 2022.

[24] Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee,
 Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step. In *The Twelfth International Conference on Learning Representations*, 2024.

[25] Zhihan Liu, Hao Hu, Shenao Zhang, Hongyi Guo, Shuqi Ke, Boyi Liu, and Zhaoran Wang.
 Reason for future, act for now: A principled framework for autonomous llm agents with provable
 sample efficiency, 2024.

[26] Jianqiao Lu, Zhiyang Dou, Hongru Wang, Zeyu Cao, Jianbo Dai, Yingjia Wan, and Zhijiang
 Guo. Autopsv: Automated process-supervised verifier, 2024.

[27] Liangchen Luo, Yinxiao Liu, Rosanne Liu, Samrat Phatale, Meiqi Guo, Harsh Lara, Yunxuan Li,
 Lei Shu, Yun Zhu, Lei Meng, Jiao Sun, and Abhinav Rastogi. Improve mathematical reasoning
 in language models by automated process supervision, 2024.

- [28] Kaixin Ma, Hongming Zhang, Hongwei Wang, Xiaoman Pan, and Dong Yu. LASER: LLM
   agent with state-space exploration for web navigation. In *NeurIPS 2023 Foundation Models for Decision Making Workshop*, 2023.
- [29] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri
  Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad
  Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Selfrefine: Iterative refinement with self-feedback. In A. Oh, T. Naumann, A. Globerson, K. Saenko,
  M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 46534–46594. Curran Associates, Inc., 2023.
- [30] Allen Nie, Yi Su, Bo Chang, Jonathan N. Lee, Ed H. Chi, Quoc V. Le, and Minmin Chen.
   Evolve: Evaluating and optimizing llms for exploration, 2024.
- [31] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. 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 Christiano,
   Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feed back. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA, 2022. Curran Associates Inc.
- [32] Pranav Putta, Edmund Mills, Naman Garg, Sumeet Motwani, Chelsea Finn, Divyansh Garg, and
   Rafael Rafailov. Agent q: Advanced reasoning and learning for autonomous ai agents, 2024.
- [33] Zehan Qi, Xiao Liu, Iat Long Iong, Hanyu Lai, Xueqiao Sun, Jiadai Sun, Xinyue Yang, Yu Yang,
   Shuntian Yao, Wei Xu, Jie Tang, and Yuxiao Dong. WebRL: Training LLM web agents via self evolving online curriculum reinforcement learning. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [34] Yuxiao Qu, Tianjun Zhang, Naman Garg, and Aviral Kumar. Recursive introspection: Teaching
  language model agents how to self-improve. In A. Globerson, L. Mackey, D. Belgrave, A. Fan,
  U. Paquet, J. Tomczak, and C. Zhang, editors, *Advances in Neural Information Processing Systems*, volume 37, pages 55249–55285. Curran Associates, Inc., 2024.
- [35] Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2025.
- [36] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and
   Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model.
   In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [37] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal
   policy optimization algorithms. *CoRR*, abs/1707.06347, 2017.
- [38] Amrith Setlur, Chirag Nagpal, Adam Fisch, Xinyang Geng, Jacob Eisenstein, Rishabh Agarwal,
   Alekh Agarwal, Jonathan Berant, and Aviral Kumar. Rewarding progress: Scaling automated
   process verifiers for LLM reasoning. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [39] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflex ion: language agents with verbal reinforcement learning. In A. Oh, T. Naumann, A. Globerson,
   K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 8634–8652. Curran Associates, Inc., 2023.
- [40] Yifan Song, Da Yin, Xiang Yue, Jie Huang, Sujian Li, and Bill Yuchen Lin. Trial and error:
   Exploration-based trajectory optimization of LLM agents. In Lun-Wei Ku, Andre Martins,
   and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7584–7600, Bangkok, Thailand,
   August 2024. Association for Computational Linguistics.

- [41] Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec
   Radford, Dario Amodei, and Paul Christiano. Learning to summarize from human feedback. In
   *Proceedings of the 34th International Conference on Neural Information Processing Systems*,
   NIPS '20, Red Hook, NY, USA, 2020. Curran Associates Inc.
- [42] Malcolm J. A. Strens. A bayesian framework for reinforcement learning. In *Proceedings of the Seventeenth International Conference on Machine Learning*, ICML '00, page 943–950, San
   Francisco, CA, USA, 2000. Morgan Kaufmann Publishers Inc.
- [43] William R. Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3/4):285–294, 1933.
- [44] S. Thrun. The role of exploration in learning control. In D.A. White and D.A. Sofge, editors,
   *Handbook for Intelligent Control: Neural, Fuzzy and Adaptive Approaches*. Van Nostrand
   Reinhold, Florence, Kentucky 41022, 1992.
- Ionathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, Lisa Wang,
   Antonia Creswell, Geoffrey Irving, and Irina Higgins. Solving math word problems with
   process- and outcome-based feedback, 2022.
- [46] Siddharth Verma, Justin Fu, Sherry Yang, and Sergey Levine. CHAI: A CHatbot AI for taskoriented dialogue with offline reinforcement learning. In Marine Carpuat, Marie-Catherine
  de Marneffe, and Ivan Vladimir Meza Ruiz, editors, *Proceedings of the 2022 Conference*of the North American Chapter of the Association for Computational Linguistics: Human
  Language Technologies, pages 4471–4491, Seattle, United States, July 2022. Association for
  Computational Linguistics.
- [47] Somin Wadhwa, Silvio Amir, and Byron C Wallace. Investigating mysteries of CoT-augmented
   distillation. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 6071–6086,
   Miami, Florida, USA, November 2024. Association for Computational Linguistics.
- [48] Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhifang Sui. Math-shepherd: Verify and reinforce LLMs step-by-step without human annotations. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9426–9439, Bangkok, Thailand, August 2024. Association for Computational Linguistics.
- [49] Tengyang Xie, Dylan J. Foster, Akshay Krishnamurthy, Corby Rosset, Ahmed Awadallah, and
   Alexander Rakhlin. Exploratory preference optimization: Harnessing implicit q\*-approximation
   for sample-efficient rlhf, 2024.
- [50] Rongwu Xu, Zehan Qi, and Wei Xu. Preemptive answer "attacks" on chain-of-thought reasoning.
   In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Findings of the Association for Computational Linguistics: ACL 2024*, pages 14708–14726, Bangkok, Thailand, August 2024.
   Association for Computational Linguistics.
- [51] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and
   Yuan Cao. React: Synergizing reasoning and acting in language models. In *The Eleventh International Conference on Learning Representations*, 2023.
- [52] Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. STar: Bootstrapping reasoning with
   reasoning. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors,
   *Advances in Neural Information Processing Systems*, 2022.
- [53] Yuexiang Zhai, Hao Bai, Zipeng Lin, Jiayi Pan, Shengbang Tong, Yifei Zhou, Alane Suhr,
  Saining Xie, Yann LeCun, Yi Ma, and Sergey Levine. Fine-tuning large vision-language
  models as decision-making agents via reinforcement learning. In A. Globerson, L. Mackey,
  D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors, *Advances in Neural Information Processing Systems*, volume 37, pages 110935–110971. Curran Associates, Inc.,
  2024.

- [54] Shenao Zhang, Donghan Yu, Hiteshi Sharma, Han Zhong, Zhihan Liu, Ziyi Yang, Shuohang
   Wang, Hany Hassan Awadalla, and Zhaoran Wang. Self-exploring language models: Active
- <sup>594</sup> preference elicitation for online alignment. *Transactions on Machine Learning Research*, 2025.
- 595 [55] Yifei Zhou, Andrea Zanette, Jiayi Pan, Sergey Levine, and Aviral Kumar. ArCHer: Training
- language model agents via hierarchical multi-turn RL. In *Forty-first International Conference on Machine Learning*, 2024.

# **Appendix**

## **Table of Contents**

602	А	Related Work	15
603		A.1 Exploration with LLMs	15
604		A.2 Training LLMs for Multi-Turn Tasks	16
605	В	Broader Impacts	16
606	С	Approach Details	16
607		C.1 Supervised Finetuning (SFT)	16
608		C.2 Collecting Trajectories	18
609		C.3 Training Critic via Supervised Learning	18
610		C.4 Training Actor via OnlineDPO	18
611	D	Additional Experimental Results	19
612		D.1 What is the optimal ratio between on-policy and hindsight data in the replay buffer?	19
613		D.2 Full Results	19
614		D.3 Qualitative: How does hindsight help exploration?	20
615	E	Domain: TwentyQuestions Details	21
616		E.1 Environment Setup	21
617		E.2 Agent Prompt	21
618		E.3 Summary Prompt	22
619		E.4 Hindsight Proposer Prompt	23
620	F	Domain: GuessMyCity Details	25
621		F.1 Environment Setup	25
622		F.2 Agent Prompt	25
623		F.3 Summary Prompt	27
624		F.4 Hindsight Proposer Prompt	29
625	G	Domain: CarDealer Details	30
626		G.1 Environment Setup	30
627		G.2 Agent Prompt	32
628		G.3 Summary Prompt	36
629		G.4 Hindsight Proposer Prompt	38
630 631 -			

## 633 A Related Work

#### 634 A.1 Exploration with LLMs

Early works investigate LLMs' ability to explore under in-context learning, where only the input to the model is updated instead of the model parameters. Some study this ability in multi-armed bandit problems [12, 20, 25, 30], while others investigate whether LLMs can self-refine and explore actions given feedback [5, 19, 29, 39]. However, prompting-based approaches require significant engineering efforts (e.g., by explicitly modeling the structure of the state space [28]).

RLHF [31] provides a practical framework for training and aligning LLMs. To improve the sample
efficiency, recent works propose adding an optimistic exploration bonus to the RLHF loss [8, 14,
49, 54], but they are limited to single-turn and are difficult to apply to the multi-turn setting that
requires exploration both at the turn level and the token level. Closer to our work are ones that utilize

posterior sampling to guide exploration [2, 14]. Dwaracherla et al. [14] estimate uncertainty via an 644 ensemble of reward models, but it does not directly train the LLM and is also limited to single-turn 645 tasks. Concurrent work [2] uses LLMs to explicitly sample from the posterior over possible MDPs 646 and show positive results for multi-turn deterministic environments. However, they discuss how their 647 approach fails to scale in stochastic domains, which include task-oriented dialogue tasks explored in 648 our work. Instead of explicitly estimating and sampling from the posterior, HOPE takes inspiration 649 from works that utilize LLMs' hindsight reasoning ability [21, 47, 50, 52]. Given the summary of a 650 completed trajectory in hindsight, the LLM implicitly estimates a posterior over plausible MDPs and 651 samples possible actions that would be optimal under the MDPs in its reasoning. 652

#### 653 A.2 Training LLMs for Multi-Turn Tasks

Reinforcement learning has been leveraged to train LLM agents on reasoning [13, 24, 48], learning 654 self-correction [22, 34], code generation [18, 23], and interactive environments [4, 7, 32, 33, 53]. A 655 class of approaches utilizes rejection sampling to only fine-tune the LLM on successful trajectories [10, 656 15, 52]. Some works also utilize failed trajectory by learning from expert correction [11] or a 657 contrastive loss [40]. Meanwhile, ARCHER [55] frames the multi-turn problem as a hierarchical 658 MDP where the lower-level MDP considers each token as actions and the higher-level MDP optimizes 659 rewards over turns. In addition to only optimizing sparse reward signals, recent works adopt an 660 actor-critic framework, where they train a critic, or a process reward model, that provides dense 661 supervision [24, 26, 27, 38, 45, 48]. While our work is agnostic to a specific RL algorithm as we 662 focus on improving exploration in multi-turn tasks, we also practically show a simple, scalable 663 implementation of actor-critic RL that trains a critic and an actor over iterations. 664

#### **B Broader Impacts**

Enabling LLMs to learn through online interaction unlocks significant societal benefits. In domains
 like customer service and software engineering, LLM assistants can automate routine tasks, allowing
 professionals to focus on higher-order, intellectually demanding problems.

However, self-improving agents introduce risks if not properly constrained. Without well-defined
 rewards and safeguards, such agents may learn behaviors misaligned with human values even if it is
 optimizing reward signals. Moreover, these capabilities can be weaponized—malicious actors could
 make LLM automatically optimize for harmful tasks such as spreading misinformation. Mitigating
 these risks requires rigorous safety mechanisms, ethical oversight, and robust evaluation protocols to
 ensure alignment with human interests and societal good.

### 675 C Approach Details

Table 3 reports the hyperparameters, the number of gpu/cpu/memory, and estimated training time used during training. We train our models on NVIDIA RTX 6000/NVIDIA RTX 6000 Ada, and we build upon the training code in OpenInstruct<sup>1</sup>. Below, we describe each stage of training: SFT to get  $\pi_0$  (Section C.1), collecting critic data (Section C.2), iteratively training critic (Section C.3), and iteratively training actor (Section C.4).

#### 681 C.1 Supervised Finetuning (SFT)

We train Llama-3.2-3B-Instruct on gpt4o rollouts on the training set (and there are 3 trajectories per game). We format the dataset in standard SFT format, where the input includes instruction and the state (the history of observations and actions so far)  $s_t = \{o_0, a_0, \dots, o_t\}$ , and the output is the gpt4o's action  $a_t$ . Note that we follow ReAct [51] and make  $a_t$  contain both the actual action to take and the reasoning for this action.

For CarDealer, because the agent has to first make API calls to the dealership database before talking to the user, the model is trained on equal amount of data with corresponding prompts for both tasks.

<sup>&</sup>lt;sup>1</sup>https://github.com/allenai/open-instruct

Dataset	TwentyQuestion	GuessMyCity	CarDealer
SFT			
batch size	4	4	2
gradient accumulation steps	16	16	12
train epochs		3	
learning rate		3.00e-05	
lr schedular		cosine	
# gpus	2	2	4
# cpus	2	2	4
Mem (GB)	80	80	200
Estimated Time (hrs)	1.25	1.25	1.50
Collecting Trajectory			
# onpolicy traj		14	
policy sample temp		1.0	
# offpolicy traj	16	16	32
# timesteps to sample		2	
sampling range $(T = len(traj))$	[0.37, 0.67)	[0.37, 0.87)	[2, T - 1)
# counterfactuals at a timestep		2	
hindsight proposer $\pi^{II}$ temp		0.3	
Critic Training			
$\alpha$ (% of hindsight data)	0.4	0.5	0.5
batch size		4	
gradient accumulation steps		16	
train epochs		1	
learning rate		5.00e-06	
lr schedular	2	linear	4
# gpus	2	2	4
# cpus	2	2	4
Mem (GB)	80	80	200
Estimated Time (nrs)	2.00	2.00	2.50
Actor Training - OnlineDPO		2	
batch size		2	
gradient accumulation steps	6	6	1
train epochs		1	
learning rate		8.00e-08	
Ir schedular		linear	
generation temp	4	0.7	(
# gpus (1 used for generator)	4	4	0
# cpus	200	200	4
Ivicili (UD) Estimated Time (hrs)	200	200	230
Esumated Time (firs)	2.23	5.00	5.40
Shared parameters	20.49	2000	2500
max seq length	2048	5000 A damW	3300
opumizer		Auantw	

Table 3: Training hyperparameters and estimated training time.

#### 689 C.2 Collecting Trajectories

<sup>690</sup> We leverage a fast inference library, SG-Lang<sup>2</sup> to serve both the environment simulator and the <sup>691</sup> current policy efficiently.

For all domains, we first collect 14 trajectories per game on the training set via the current policy 692 with a temperature of 1.0. Then, we randomly sample 4 trajectories to augment. For each of 693 these trajectories, we randomly sample 2 timesteps based on a domain-specific range with respect 694 to the trajectory length. On a first principle, for domains that require search (TwentyQuestions, 695 GuessMyCity), the counterfactual action is not useful when the timestep is too early (when the policy 696 can already eliminate common categories) or when the timestep is too late (when the proposer can 697 directly end the game by guess the word). At each sampled timestep, we use the hindsight proposer 698 to generate 2 distinct counterfactual actions. Then, we augment the original trajectory by splicing 699 in the counterfactuals before completing the trajectory with the original policy. Thus, with all the 700 parameters, we generate  $4 \times 2 \times 2 = 16$  hindsight trajectories. 701

Note that, because CarDealer requires generating two-part counterfactuals (first 2 counterfactual API calls, then 2 counterfactual responses to the user), we generate  $4 \times 2 \times (2+2) = 32$  hindsight trajectories in total.

#### 705 C.3 Training Critic via Supervised Learning

We fix the dataset to 10k datapoints for TwentyQuestions/GuessMyCity and 20k datapoints for CarDealer (because we have 10k datapoints for making API call and 10k datapoints for responding to buyer).

After calculating the Q-values via MC estimate, we normalize the Q-targets to be between [0, 1]. To maintain a balance dataset, we ensure that 50% of datapoints have low values [0, 0.5) and 50% of datapoints have high values [0.5, 1].

For low-value data, we prioritize using datapoints from on-policy trajectories (i.e., trajectories that are generated by only sampling from the current policy). If there isn't enough datapoints, we add hindsight trajectories until sufficient.

• For high-value data, we use  $\alpha$  to control the percentage of datapoints from hindsight trajectories.

Critic is initialized with weights from  $\pi_0$ . It has an additional randomly initialized linear layer of dimension [hidden\_dim, 1] because it predicts a scalar value. To select the best intermediate checkpoint for actor training, we evaluate the Best-of-N performance of using the current policy  $\pi_{i-1}$ as the generator and the checkpoint as the critic. Specifically, at each step,  $\pi_{i-1}$  generates (N = 15) actions at temperature 0.7, and we execute the action with the highest score from the critic. The best critic is one that has the highest Best-of-N success rate on the validation set.

#### 722 C.4 Training Actor via OnlineDPO

Similar to the critic, the actor is initialized with weights from  $\pi_0$ . The parameter  $\beta$  controls the probability of generating counterfactual actions from the hindsight proposer. For most experiments, we set  $\beta = 0.0$ . When we use hindsight-guided action experation, we set  $\beta = 0.5$ 

Due to budget constraints, we distill the GPT-40 hindsight proposer into a Llama-3.2-3B-Instruct. To create the training data, we first consolidate all the counterfactual actions generated during Section C.2 and use SFT to fully fine-tune Llama-3.2-3B-Instruct for 2 epoches using learning rate 3e-5. Then, we further fine-tune the model on only counterfactuals that lead to success for another epoch using learning 3e-6. The other training hyperparameters are the same as the one for SFT in Table 3.

#### 731 C.4.1 Evaluation

For all domains, we evaluate the policy once per game on the training set. For TwentyQuestions and CarDealer, we evaluate the policy three times per game on the validation and test set. For GuessMyCity, we evaluate the policy ten times per game on the validation and test set. We always select the intermediate checkpoint that has the average highest validation success rate.

<sup>2</sup>https://github.com/sgl-project/sglang

						LEAP		Reject-S		PRM+RL		HOPE	
			gpt4o	3B	$\pi_0$	$\pi_1$	$\pi_2$	$\pi_1$	$\pi_2$	$\pi_1$	$\pi_2$	$\pi_1$	$\pi_2$
TwentyQuestions	Train	Return Success	1.00 (0.09) 0.60 (0.04)	0.06 (0.02) 0.16 (0.03)	1.04 (0.09) 0.66 (0.04)	1.13 (0.09) 0.75 (0.04)	1.15 (0.09) 0.77 (0.04)	1.25 (0.08) 0.79 (0.04)	1.24 (0.09) 0.79 (0.04)	1.53 (0.08) 0.86 (0.03)	1.34 (0.09) 0.81 (0.04)	1.28 (0.09) 0.77 (0.04)	1.29 (0.09) 0.80 (0.04)
	Val	Return Success	1.00 (0.09) 0.65 (0.05)	0.04 (0.02) 0.16 (0.04)	1.01 (0.05) 0.62 (0.02)	1.21 (0.09) 0.83 (0.04)	1.00 (0.11) 0.64 (0.05)	1.22 (0.10) 0.79 (0.04)	1.07 (0.09) 0.79 (0.04)	1.38 (0.05) 0.79 (0.02)	1.10 (0.05) 0.79 (0.02)	1.18 (0.04) 0.77 (0.02)	1.34 (0.05) 0.78 (0.02)
	Test	Return Success	1.00 (0.11) 0.61 (0.06)	0.09 (0.04) 0.28 (0.05)	1.08 (0.11) 0.67 (0.05)	1.16 (0.12) 0.68 (0.06)	0.52 (0.09) 0.45 (0.06)	0.63 (0.11) 0.55 (0.06)	0.90 (0.13) 0.63 (0.06)	1.19 (0.16) 0.55 (0.06)	1.05 (0.09) 0.82 (0.05)	1.15 (0.13) 0.77 (0.05)	1.91 (0.09) 0.97 (0.02)
	Total	Return Success	1.00 (0.06) 0.62 (0.03)	0.06 (0.01) 0.19 (0.02)	1.03 (0.04) 0.63 (0.02)	1.16 (0.06) 0.76 (0.03)	0.95 (0.06) 0.65 (0.03)	1.09 (0.06) 0.73 (0.03)	1.10 (0.06) 0.75 (0.03)	<b>1.39 (0.04)</b> 0.78 (0.02)	1.14 (0.04) 0.80 (0.02)	1.19 (0.04) 0.77 (0.02)	1.39 (0.04) 0.80 (0.02)
ty	Train	Return Success	1.00 (0.11) 0.58 (0.04)	0.03 (0.01) 0.04 (0.02)	0.93 (0.09) 0.72 (0.04)	0.60 (0.06) 0.70 (0.04)	0.39 (0.04) 0.66 (0.04)	0.50 (0.04) 0.76 (0.04)	0.52 (0.04) 0.81 (0.04)	0.36 (0.03) 0.66 (0.04)	0.63 (0.07) 0.65 (0.04)	0.37 (0.03) 0.64 (0.04)	0.87 (0.09) 0.71 (0.04)
GuessMyC	Val	Return Success	1.00 (0.11) 0.66 (0.05)	0.07 (0.03) 0.10 (0.03)	0.74 (0.05) 0.74 (0.03)	0.40 (0.05) 0.68 (0.05)	0.26 (0.03) 0.59 (0.05)	0.41 (0.02) 0.79 (0.02)	0.35 (0.02) 0.78 (0.02)	0.27 (0.01) 0.72 (0.02)	0.42 (0.02) 0.73 (0.02)	0.32 (0.03) 0.74 (0.05)	0.69 (0.07) 0.78 (0.04)
	Test	Return Success	1.00 (0.12) 0.56 (0.05)	0.01 (0.01) 0.01 (0.01)	0.80 (0.09) 0.70 (0.05)	0.35 (0.04) 0.59 (0.05)	0.29 (0.03) 0.62 (0.05)	0.46 (0.02) 0.69 (0.02)	0.31 (0.03) 0.62 (0.05)	0.35 (0.05) 0.66 (0.05)	0.54 (0.03) 0.66 (0.02)	0.45 (0.05) 0.74 (0.05)	0.91 (0.10) 0.77 (0.04)
	Total	Return Success	<b>1.00 (0.07)</b> 0.59 (0.03)	0.03 (0.01) 0.05 (0.01)	0.80 (0.04) 0.73 (0.02)	0.47 (0.03) 0.66 (0.03)	0.32 (0.02) 0.63 (0.03)	0.44 (0.01) 0.74 (0.01)	0.37 (0.01) 0.77 (0.02)	0.29 (0.01) 0.70 (0.02)	0.50 (0.02) 0.69 (0.02)	0.38 (0.02) 0.70 (0.03)	0.83 (0.05) 0.75 (0.03)
CarDealer (Tool)	Train	Return Success	1.00 (0.10) 0.46 (0.04)	0.33 (0.07) 0.30 (0.04)	1.27 (0.11) 0.54 (0.04)	1.17 (0.11) 0.50 (0.05)	1.27 (0.10) 0.57 (0.04)	1.23 (0.11) 0.54 (0.04)	1.27 (0.10) 0.55 (0.04)	1.31 (0.11) 0.54 (0.04)	1.31 (0.11) 0.54 (0.04)	1.38 (0.11) 0.58 (0.04)	1.34 (0.10) 0.55 (0.04)
	Val	Return Success	1.00 (0.12) 0.40 (0.05)	0.68 (0.11) 0.34 (0.05)	1.23 (0.13) 0.47 (0.05)	1.35 (0.13) 0.54 (0.05)	1.44 (0.12) 0.62 (0.05)	1.40 (0.13) 0.56 (0.05)	1.52 (0.12) 0.62 (0.05)	1.51 (0.13) 0.58 (0.05)	1.57 (0.13) 0.60 (0.05)	1.61 (0.13) 0.64 (0.05)	1.62 (0.13) 0.62 (0.05)
	Test	Return Success	1.00 (0.09) 0.59 (0.05)	0.40 (0.07) 0.51 (0.05)	1.03 (0.08) 0.66 (0.05)	0.77 (0.08) 0.52 (0.05)	0.86 (0.08) 0.57 (0.05)	1.07 (0.08) 0.68 (0.05)	1.00 (0.08) 0.62 (0.05)	1.00 (0.08) 0.67 (0.05)	1.14 (0.08) 0.72 (0.04)	1.07 (0.08) 0.69 (0.05)	1.20 (0.08) 0.77 (0.04)
	Total	Return Success	1.00 (0.06) 0.48 (0.03)	0.46 (0.05) 0.38 (0.03)	1.18 (0.06) 0.55 (0.03)	1.10 (0.07) 0.52 (0.03)	1.19 (0.06) 0.59 (0.03)	1.24 (0.06) 0.59 (0.03)	1.27 (0.06) 0.60 (0.03)	1.28 (0.06) 0.60 (0.03)	1.34 (0.06) 0.62 (0.03)	1.36 (0.06) 0.63 (0.03)	1.38 (0.06) 0.64 (0.03)

Table 4: **Policy performance on all data splits of three task-oriented conversational environments.** We report the average normalized return and success rate with standard error, and we highlight the top-performing approaches. The return is normalized with respect to gpt4o's performance.

#### 736 D Additional Experimental Results

#### 737 D.1 What is the optimal ratio between on-policy and hindsight data in the replay buffer?

We study how the proportion of hindsight data influences critic training by varying the data mixture 738 used in Algorithm 2. For each setting, we construct datasets with a fixed size of 10k transitions but 739 740 vary the ratio of on-policy to hindsight-generated data. As shown in Figure 3, incorporating even a small amount of hindsight data substantially improves performance: adding just 20–40% hindsight 741 data increases  $\pi_1$ 's success rate from 0.55 (the PRM+RL baseline with no hindsight) to over 0.70. 742 Performance quickly plateaus, and pure hindsight data does not yield additional gains. In practice, we 743 use a 60% hindsight and 40% on-policy mixture for the Twenty Questions domain, which achieves 744 the highest observed success rate of 0.77. Note that this setting is not purely off-policy-the critic 745 still sees some on-policy data from  $\pi_0$ , which stabilizes learning and anchors value estimates around 746 feasible behaviors. 747

#### 748 D.2 Full Results

Table 4 shows the complete result of all approaches 749 performance on different data split. For TwentyQues-750 tions and CarDealer, HOPE outperform all baselines 751 including gpt4o, achieving the highest normalized 752 return and success rate in total. For GuessMyCity, 753 although HOPE has slightly lower normalized returns 754 (0.83) compared to gpt4o (1.00), it still outperforms 755 756 the remaining baselines. Similarly, it has the second 757 highest success rate  $(0.75 \pm 0.03)$  that is within stan-758 dard error of the highest success rate  $(0.77 \pm 0.02)$ .

Overall, other approaches tend to overfit on the train 759 set, while HOPE is able to maintain high performance 760 on the test set. For example, for TwentyQuestions, 761 although PRM+RL  $\pi_1$  achieves the highest normalized 762 return (1.53) and success rate (1.38) on the training 763 set, it significantly underperforms on the test set with 764 normalized return of 1.19 and a lower success rate 765 0.55 than  $\pi_0$  (0.67). In contrast, although HOPE has 766



Figure 3: Effect of critic data mixture on downstream policy performance. We vary the percentage of hindsight data via the hyperparameter  $\alpha$  and train the corresponding critic and  $\pi 1$ . We report the average success with standard error in TwentyQuestions.

- <sup>767</sup> the third highest normalized return (1.29) and suc-
- <sup>768</sup> cess rate (0.80), it maintains the highest normalized
- return (1.91) and success rate (0.97), significantly

outperforming all other baselines.



Figure 4: **Qualitative Example of how** HOPE **improves exploration during critic data collection.** The three columns shows traces of different exploration strategies collecting data for the same game (where the secret word is "Painting". In the 1st column, H-Temp (PRM+RL) simply sample actions from the policy with high temperature. In the 2nd column, HOPE first summarize the completed rollout before using the hindsight proposer to generate counterfactual actions. In the 3rd column, Explore-Prompt explicitly prompt the policy to generate actions different from the common one.

# 771 D.3 Qualitative: How does hindsight help772 exploration?

Figure 4 demonstrations how HOPE effectively explore the state and action space, leading to successful trajectories, compared to other exploration strategies.

In the 1st column, H-Temp (PRM+RL) simply samples the current policy (in this case  $\pi_0$ ) at a high temperature of 1.0. The policy is ineffective at exploring the high-value state space and action space. Although it attempts to eliminate high level categories, it fails to consider categories such as "art supply", leading the policy to exhaust the maximum number of questions that it can ask.

In the 3rd column, Explore-Prompt first samples 5 actions from the current policy. Then, given this set of high-probability actions, prompt asks the policy to generate some alternative actions that differ from the list. The policy is able to identify the correct category, but the policy still takes too long and uses all the questions before being able to guess "Painting." Although Explore-Prompt can increase the diversity of states and actions visited, it fails to intentionally explore high-value regions of the state and action space.

In contrast, in the 2nd column, HOPE first accurately identifies the original policy's mistake in its summary  $\phi$  of the completed trajectory: "it did not focus on the art category... this oversight ... contributes to the agent's failure." Then, the hindsight proposer generates an effective counterfactual action "is it a decorative item?", which leads to the policy quickly guess the correct word at timestep t = 13.

### 790 E Domain: TwentyQuestions Details

#### 791 E.1 Environment Setup

Training set contains 15 object categories with 110 objects in total. Validation set has unseen objects in the train categories with 28 objects in total. Test set has 2 unseen object categories with 20 objects.

When the agent successfully guess the word, the reward is 0. Otherwise, the reward is -1 to encourage the agent to complete the game as soon as possible. The agent has maximum of 20 steps to complete the game.

#### 797 E.2 Agent Prompt

800

When there is only one question left, we programmatically change the mode to 'input\_final', which require the agent to guess a specific object.

```
{% if mode == 'input' %}
801
    You are an intelligent player playing a game of twenty questions. Your
802
    \hookrightarrow objective is to ask the minimal number of yes-no questions in order to
803
    \hookrightarrow guess the identity of the entity/object chosen by an oracle. You can
804
    \hookrightarrow only ask 20 questions in total.
805
806
807
    The entity/object that the orcacle can choose from are:
    {{ all_obj_list }}
808
809
    Your goal is to generate a yes-no question that either (1) help you narrow
810
    \hookrightarrow down the possible things as much as possible or (2) guess the entity/
811
    \hookrightarrow object directly.
812
813
    Please follow these general instructions:
814
    * You MUST ask a yes-no question.
815
    * If you are guessing the entity/object directly, you MUST ask a question
816
    \hookrightarrow in the format "Is it {your_guess}?". your_guess must be one of the
817
    \hookrightarrow entity/object that the oracle can choose from.
818
    * Before you ask the question, you MUST intelligently reason about the
819
    \hookrightarrow history of previous questions and answers to ask the most informative
820
    \hookrightarrow question. Your reasoning should also be based on the list of entity/
821
    \hookrightarrow object that the oracle can choose from.
822
    * Consult the history of previous questions and answers to see what
823
    \hookrightarrow questions you have asked already so as to not repeat your questions.
824
    * Do NOT repeat the same question. It's going to yield the same result.
825
    * Do NOT get stuck on one idea and try to branch out if you get stuck.
826
827
    Below is the history of previous questions and answers:
828
829
    {{ observation_action_history }}
830
831
    You MUST generate a response in the following format. Please issue only a
    \hookrightarrow single question at a time.
832
    REASON:
833
    Rationale for what question to ask next based on the previous history and
834
    \hookrightarrow the list of things that the oracle can choose from.
835
    QUESTION:
836
    The question to be asked. The question must be a yes-no question.
837
    {% elif mode == 'input_final' %}
838
    You are an intelligent player playing a game of twenty questions. Your
839
    \hookrightarrow objective is to ask the minimal number of yes-no questions in order to
840
    \hookrightarrow guess the identity of the entity/object chosen by an oracle. You can
841
    \hookrightarrow only ask 20 questions in total.
842
843
   The entity/object that the orcacle can choose from are:
844
```

```
{{ all_obj_list }}
845
846
    Your goal is to generate a yes-no question that either (1) help you narrow
847
    \hookrightarrow down the possible things as much as possible or (2) guess the entity/
848
    \hookrightarrow object directly.
849
850
    Please follow these general instructions:
851
    * You MUST ask a yes-no question.
852
    * If you are guessing the entity/object directly, you MUST ask a question
853
    \hookrightarrow in the format "Is it {your_guess}?". your_guess must be one of the
854
    \hookrightarrow entity/object that the oracle can choose from.
855
    * Before you ask the question, you MUST intelligently reason about the
856
    \hookrightarrow history of previous questions and answers to ask the most informative
857
    \hookrightarrow question. Your reasoning should also be based on the list of entity/
858
    \hookrightarrow object that the oracle can choose from.
859
    * Consult the history of previous questions and answers to see what
860
    \hookrightarrow questions you have asked already so as to not repeat your questions.
861
    * Do NOT repeat the same question. It's going to yield the same result.
862
    * Do NOT get stuck on one idea and try to branch out if you get stuck.
863
864
    Below is the history of previous questions and answers:
865
    {{ observation_action_history }}
866
867
    You have already asked 19 questions, so this is your final guess. Your goal
868
    \hookrightarrow is to consult the list of entity/object that the orcale can choose from
869
    \hookrightarrow and quess the entity/object directly.
870
871
    You MUST generate a response in the following format. Please issue only a
872
    \hookrightarrow single question at a time.
873
    REASON:
874
    Rationale for what entity/object to quess based on the previous history. In
875
    \rightarrow your reason, consult the list of entity/object that the oracle can
876
    \hookrightarrow choose from to precisely state VERBATIM what you will guess.
877
    QUESTION:
878
    The question to be asked. The question must be a yes-no question in the
879
    \hookrightarrow format "Is it {your_guess}?". your_guess must be one of the entity/
880
    \hookrightarrow object that the oracle can choose from.
881
882
    {% elif mode == 'output' %}
    REASON:
883
    {{ reason }}
884
    QUESTION:
885
    {{ action }}
886
    {% elif mode == 'output_no_reason' %}
887
    QUESTION:
888
    {{ action }}
889
    {% endif %}
899
```

#### 892 E.3 Summary Prompt

```
{% if system %}
894
    You are an intelligent assistant summarizing a game of twenty questions,
895
    \hookrightarrow where an agent is trying to guess a word by asking at most 20 yes-no
896
    \hookrightarrow questions.
897
898
    ## Overall information about the game
899
   Here is the list of possible secret words:
900
    {{ all_obj_list }}
901
902
```

```
## What you receive as input
903
    You are given (1) a chat history of the questions and answers and (2) the
904
    \hookrightarrow secret word that the agent is trying to guess.
905
906
    ## Your goal and rules to follow
907
    You must summarize the chat history in 2-4 sentences (what the agent asked,
908
    \hookrightarrow whether the agent succeeded or not). You must also mention what the
909
    \hookrightarrow actual secret word is and what general category the secret word is in.
910
911
    You should also discuss in your summary the strategy of the agent. You can
912
    \hookrightarrow refer to the list of possible secret words to help with your discussion:
913
    * Did the agent repeat similar questions that seem less helpful?
914
    * Did the agent move on to ask about specific things/entities too quickly?
915
    \hookrightarrow Or did the agent keep asking broader, higher-level questions even though
916
    \hookrightarrow they can directly guess the word with the information it had?
917
    * If the agent succeeded, what type of questions did it ask to help it
918
    \hookrightarrow succeed?
919
    * If the agent failed, what are some possible reasons on why it failed? Did
920
    \hookrightarrow it ask questions that violate previous questions and answers? Did it
921
    \hookrightarrow ask about things not in the list of possible secret words?
922
923
    You can directly generate the summary as a paragraph. You should not add
924
    \hookrightarrow any formatting (e.g., markdown) when generating the summary.
925
    {% endif %}
926
    {% if not system %}
927
    {% if mode == 'input' %}
928
    ## Chat History
929
    {{ observation_action_history }}
930
931
    ## Secret Word
932
    {{ goal }}
933
    {% endif %}
934
    {% endif %}
935
```

#### 937 E.4 Hindsight Proposer Prompt

```
{% if system %}
939
    You are an intelligent teacher who gives guidance on what question to ask
940
    \hookrightarrow in a game of twenty questions. In a game of twenty questions, a player
941
    \hookrightarrow is trying to guess a secret word by asking at most 20 yes-no questions.
942
943
    ## What you receive as input
944
945
    You are given (1) a list of the possible secret words and (2) the chat
946
    \hookrightarrow history of the questions that the player has asked so far and the
947
    \hookrightarrow corresponding answers. To further help you make wise judgements and
948
    \hookrightarrow provide helpful guidance to the player, you are also given (3) a summary
    \hookrightarrow of the complete game (which includes whether the player has succeeded
949
    \hookrightarrow in the end, the actual secret word, and the general category that the
950
    \hookrightarrow secred word is in).
951
952
    ## Your goal and rules to follow
953
    Your goal is to use your hindsight reasoning ability to generate {{
954
    \hookrightarrow num_responses}} alternative questions that the player should have asked
955
    \hookrightarrow given the current chat history. Your process is to: 1. reason about all
956
    \hookrightarrow the information (including the summary); 2. generate some questions that
957
    \hookrightarrow the player could have asked; 3. generate some plausible reasoning that
958
    \hookrightarrow the player could have come up with based on the current chat history.
959
960
```

```
Please follow these general instructions:
961
     - **Ask feasible questions:** The question that you ask MUST be a question
 962
     \hookrightarrow that is possible for the player to ask given the current chat history.
 963
     - **In teacher_reason, reason with all the information:** You are the
 964
     \hookrightarrow teacher. In your teacher reasoning, you MUST make reference to specific
 965
     \hookrightarrow things mentioned under ### List of the possible secret words and the
 966
     \hookrightarrow chat history. You MUST ask a question that is possible and reasonable to
 967
     \hookrightarrow ask given the current chat history. You can make use of the ### Summary
 968
     \hookrightarrow of the entire game to help you identify better but STILL FEASIBLE
 969
     \hookrightarrow question to ask given the current chat history. For example, the general
 970
     \hookrightarrow category that the secret word is in could help inform you what kind of
 971
     \hookrightarrow question to ask.
 972
     - **In player_reason, reasoning should not include secret information from
 973
     \hookrightarrow the summary:"** You are generating what is the possible reasoning that a
 974
     \hookrightarrow player can have based on the chat history in order to generate the
 975
 976
     \hookrightarrow question. The reasoning must not reveal that you know the secret object
     \hookrightarrow or the general category. It must be a feasible reasoning based on the
 977
     \hookrightarrow chat history alone.
 978
     - **Question should only be asking about the possible secret words:** Your
979
     \hookrightarrow question MUST only ask about objects in ### List of the possible secret
 980
     \hookrightarrow words.
 981
     - **Ask a new question.** You MUST NOT simply repeat previously asked
 982
     \hookrightarrow questions. You MUST NOT simply combine multiple previously asked
 983
     \hookrightarrow questions.
 984
     - If you think it is reasonable and feasible to directly guess the object
 985
     \hookrightarrow given the current chat history, you can propose the question to directly
 986
     \hookrightarrow ask: "Is it {your_guess}?". your_guess must be one of the words in ###
987
     \hookrightarrow List of the possible secret words.
 988
 989
     ## Output format
 990
     The output is a list of JSON containing {{num_responses}} different pairs
 991
     \hookrightarrow of reasoning and feasible question that you would have asked.
 992
     '''json
 993
     Γ
 994
         {
 995
              "teacher_reason": "string: your rationale for the question that you
 996
             \hookrightarrow think the player should have asked instead.",
 997
              "question": "string: your question",
998
              "player_reason": "string: If you are the player who does not know
999
             \hookrightarrow the secret object or the general category that the object is in,
1000
             \hookrightarrow what will be your reasoning in order to generate the action? You
1001
             \hookrightarrow MUST only refer to the chat history. You MUST NOT talk about the
1002
             \hookrightarrow summary, which the player does not have access to. You MUST NOT
1003
             \hookrightarrow reveal what the secret object is. You MUST NOT reveal what the
1004
             \hookrightarrow general category that the secret object is in.",
1005
         }
1006
1007
          . . .
1008
     ٦
     ...
1009
1010
     ## Overall information about the current game
1011
     ### List of the possible secret words
1012
     {{ all_obj_list }}
1013
1014
     ### Summary of the entire game
1015
     {{ summary }}
1016
     {% endif %}
1017
     {% if not system %}
1018
     {% if mode == 'input' %}
1019
```

```
1020 Here is the chat history:

1021 {{ observation_action_history }}

1022

1023 At this point in the chat history, what question would you have asked? You

1024 \hookrightarrow MUST generate {{num_responses}} different, diverse questions.

1025 {% endif %}

1026 {% endif %}
```

#### **1028 F Domain: GuessMyCity Details**

#### 1029 F.1 Environment Setup

Training set contains 15 countries, each with 8 cities, so it has 120 cities in total. Validation set has 2 more unseen cities in the 15 training countries, so it has 30 cities in total. Test set has 3 unseen countries, each with 10 cities, so it has 30 cities.

Because the agent can now ask open-ended questions, it is able to extract much more information per turn compared to Twenty Questions. Thus, we limit the maximum step to only 10 steps. In addition, we only allow the agent to only guess the city once. If the agent guesses incorrectly, the episode immediately terminates, and the agent receives a penalty such that its sum of reward is -10. If the agent succeed, it receives a reward of 0. All non-terminating steps get a reward of -1.

#### 1038 F.2 Agent Prompt

```
1039
     {% if mode == 'input' %}
1040
    You are an intelligent player playing a game where you must guess where
1041
     \hookrightarrow someone is from. Your objective is to ask the minimal number of open-
1042
     \hookrightarrow ended questions that are diverse in order to accurately guess their city
1043
        of origin. You can ask 10 questions to determine where they are from.
1044
1045
     \hookrightarrow You are smart, so you will ask the question that will narrow down the
     \hookrightarrow possible cities the person is from as much as possible.
1046
1047
     The cities that the oracle can choose from are:
1048
     {{ all_city_list }}
1049
1050
     Your goal is to generate an open-ended question that (1) helps you narrow
1051
     \hookrightarrow down the possible cities as much as possible or (2) allows you to make a
1052
     \hookrightarrow final guess.
1053
1054
     Please follow these general instructions:
1055
     * You MUST ask an open-ended question.
1056
     * You CANNOT ask the oracle for the name of the city or country.
1057
1058
     * If you are guessing the city directly, you MUST ask a question in the
1059
     \hookrightarrow format "Is the city {your_guess}?". your_guess must be one of the cities
1060
     \hookrightarrow that the oracle can choose from.
     * Before you ask the question, you MUST intelligently reason about the
1061
     \hookrightarrow history of previous questions and answers to ask the most informative
1062
     \hookrightarrow question.
1063
     * Consult the history of previous questions and answers to avoid repeating
1064
1065
     \hookrightarrow the same question.
     * Do NOT repeat the same question. It's going to yield the same result.
1066
     * Do NOT get stuck on one idea and try to branch out if you get stuck.
1067
     * Do NOT ask questions formatted like "Can you ..."
1068
1069
     * If you are certain you know the answer, MAKE A GUESS
1070
    Below is the history of previous questions and answers:
1071
     {{ observation_action_history }}
1072
1073
```

You MUST generate a response in the following format. Please issue only a 1074  $\hookrightarrow$  single question at a time. 1075 1076 REASON: 1077 Rationale for what question to ask next based on the previous history and 1078  $\hookrightarrow$  the list of cities the oracle can choose from. 1079 1080 QUESTION: 1081 The question to be asked. The question must be open ended. 1082 1083 {% elif mode == 'input\_final' %} 1084 You are an intelligent player playing a game where you must guess where 1085  $\hookrightarrow$  someone is from. Your objective is to ask the minimal number of open-1086  $\hookrightarrow$  ended questions that are diverse in order to accurately guess their city 1087  $\hookrightarrow$  of origin. You can ask 10 questions to determine where they are from. 1088 1089  $\hookrightarrow$  You are smart, so you will ask the question that will narrow down the  $\hookrightarrow$  possible cities the person is from as much as possible. 1090 1091 The cities that the oracle can choose from are: 1092 {{ all\_city\_list }} 1093 1094 Your goal is to generate an open-ended question that (1) helps you narrow 1095  $\hookrightarrow$  down the possible cities as much as possible or (2) allows you to make a 1096  $\hookrightarrow$  final guess. 1097 1098 Please follow these general instructions: 1099 \* You MUST ask an open-ended question. 1100 \* You CANNOT ask the oracle for the name of the city or country. 1101 \* If you are guessing the city directly, you MUST ask a question in the 1102  $\hookrightarrow$  format "Is it {your\_guess}?". your\_guess must be one of the cities that 1103  $\hookrightarrow$  the oracle can choose from. 1104 \* Before you ask the question, you MUST intelligently reason about the 1105  $\hookrightarrow$  history of previous questions and answers to ask the most informative 1106  $\hookrightarrow$  question. 1107 \* Consult the history of previous questions and answers to avoid repeating 1108  $\hookrightarrow$  the same question. 1109 \* Do NOT repeat the same question. It's going to yield the same result. 1110 1111 \* Do NOT get stuck on one idea and try to branch out if you get stuck. 1112 Below is the history of previous questions and answers: 1113 {{ observation\_action\_history }} 1114 1115 You have already asked 9 questions, so this is your final guess. Your goal 1116  $\hookrightarrow$  is to consult the list of cities that the orcale can choose from and 1117  $\hookrightarrow$  guess the city directly. 1118 1119 You MUST generate a response in the following format. Please issue only a 1120 1121  $\hookrightarrow$  single question at a time. 1122 1123 REASON: 1124 Rationale for what city to guess based on the previous history. In your  $\hookrightarrow$  reason, consult the list of cities that the oracle can choose from to 1125  $\hookrightarrow$  precisely state VERBATIM what you will guess. 1126 1127 QUESTION: 1128 The question to be asked. The question must be a yes-no question in the 1129  $\hookrightarrow$  format "Is the city {your\_guess}?". your\_guess must be one of the cities 1130  $\hookrightarrow$  that the oracle can choose from and be based on what you know about the 1131  $\hookrightarrow$  city so far. 1132

1133 {% elif mode == 'output' %} 1134 **REASON:** 1135 {{ reason }} 1136 QUESTION: 1137 {{ action }} 1138 1139 {% elif mode == 'output\_no\_reason' %} 1140 QUESTION: 1141 {{ action }} 1142 {% endif %} 1143

#### 1145 F.3 Summary Prompt

```
{% if mode == 'input' %}
1147
     You are an intelligent player playing a game where you must guess where
1148
     \hookrightarrow someone is from. Your objective is to ask the minimal number of open-
1149
     \hookrightarrow ended questions that are diverse in order to accurately guess their city
1150
     \hookrightarrow of origin. You can ask 10 questions to determine where they are from.
1151
     \hookrightarrow You are smart, so you will ask the question that will narrow down the
1152
1153
     \hookrightarrow possible cities the person is from as much as possible.
1154
     The cities that the oracle can choose from are:
1155
    {{ all_city_list }}
1156
1157
     Your goal is to generate an open-ended question that (1) helps you narrow
1158
     \hookrightarrow down the possible cities as much as possible or (2) allows you to make a
1159
     \hookrightarrow final guess.
1160
1161
1162
     Please follow these general instructions:
     * You MUST ask an open-ended question.
1163
     * You CANNOT ask the oracle for the name of the city or country.
1164
     * If you are guessing the city directly, you MUST ask a question in the
1165
     \hookrightarrow format "Is the city {your_guess}?". your_guess must be one of the cities
1166
     \hookrightarrow that the oracle can choose from.
1167
     * Before you ask the question, you MUST intelligently reason about the
1168
     \hookrightarrow history of previous questions and answers to ask the most informative
1169
     \hookrightarrow question.
1170
     * Consult the history of previous questions and answers to avoid repeating
1171
     \hookrightarrow the same question.
1172
     * Do NOT repeat the same question. It's going to yield the same result.
1173
1174
     * Do NOT get stuck on one idea and try to branch out if you get stuck.
     * Do NOT ask questions formatted like "Can you ..."
1175
1176
     * If you are certain you know the answer, MAKE A GUESS
1177
     Below is the history of previous questions and answers:
1178
     {{ observation_action_history }}
1179
1180
     You MUST generate a response in the following format. Please issue only a
1181
     \hookrightarrow single question at a time.
1182
1183
     REASON:
1184
     Rationale for what question to ask next based on the previous history and
1185
     \hookrightarrow the list of cities the oracle can choose from.
1186
1187
    QUESTION:
1188
    The question to be asked. The question must be open ended.
1189
1190
```

```
{% elif mode == 'input_final' %}
1191
     You are an intelligent player playing a game where you must guess where
1192
     \hookrightarrow someone is from. Your objective is to ask the minimal number of open-
1193
     \hookrightarrow ended questions that are diverse in order to accurately guess their city
1194
     \hookrightarrow of origin. You can ask 10 questions to determine where they are from.
1195
     \hookrightarrow You are smart, so you will ask the question that will narrow down the
1196
     \hookrightarrow possible cities the person is from as much as possible.
1197
1198
     The cities that the oracle can choose from are:
1199
     {{ all_city_list }}
1200
1201
     Your goal is to generate an open-ended question that (1) helps you narrow
1202
     \hookrightarrow down the possible cities as much as possible or (2) allows you to make a
1203
     \hookrightarrow final guess.
1204
1205
1206
     Please follow these general instructions:
     * You MUST ask an open-ended question.
1207
     * You CANNOT ask the oracle for the name of the city or country.
1208
     * If you are guessing the city directly, you MUST ask a question in the
1209
     \hookrightarrow format "Is it {your_guess}?". your_guess must be one of the cities that
1210
     \hookrightarrow the oracle can choose from.
1211
     * Before you ask the question, you MUST intelligently reason about the
1212
     \hookrightarrow history of previous questions and answers to ask the most informative
1213
     \hookrightarrow question.
1214
     * Consult the history of previous questions and answers to avoid repeating
1215
     \hookrightarrow the same question.
1216
     * Do NOT repeat the same question. It's going to yield the same result.
1217
     * Do NOT get stuck on one idea and try to branch out if you get stuck.
1218
1219
     Below is the history of previous questions and answers:
1220
     {{ observation_action_history }}
1221
1222
     You have already asked 9 questions, so this is your final guess. Your goal
1223
     \hookrightarrow is to consult the list of cities that the orcale can choose from and
1224
     \hookrightarrow guess the city directly.
1225
1226
     You MUST generate a response in the following format. Please issue only a
1227
1228
     \hookrightarrow single question at a time.
1229
     REASON:
1230
     Rationale for what city to guess based on the previous history. In your
1231
     \hookrightarrow reason, consult the list of cities that the oracle can choose from to
1232
     \hookrightarrow precisely state VERBATIM what you will guess.
1233
1234
     QUESTION:
1235
1236
     The question to be asked. The question must be a yes-no question in the
     \hookrightarrow format "Is the city {your_guess}?". your_guess must be one of the cities
1237
1238
     \hookrightarrow that the oracle can choose from and be based on what you know about the
     \hookrightarrow \  \  \text{city so far}.
1239
1240
     {% elif mode == 'output' %}
1241
     REASON:
1242
     {{ reason }}
1243
     QUESTION:
1244
     {{ action }}
1245
1246
     {% elif mode == 'output_no_reason' %}
1247
     QUESTION:
1248
    {{ action }}
1249
```

1253

#### 1252 F.4 Hindsight Proposer Prompt

```
{% if mode == 'input' %}
1254
     You are an intelligent player playing a game where you must guess where
1255
     \hookrightarrow someone is from. Your objective is to ask the minimal number of open-
1256
     \hookrightarrow ended questions that are diverse in order to accurately guess their city
1257
     \hookrightarrow of origin. You can ask 10 questions to determine where they are from.
1258
     \hookrightarrow You are smart, so you will ask the question that will narrow down the
1259
1260
     \hookrightarrow possible cities the person is from as much as possible.
1261
     The cities that the oracle can choose from are:
1262
     {{ all_city_list }}
1263
1264
     Your goal is to generate an open-ended question that (1) helps you narrow
1265
     \hookrightarrow down the possible cities as much as possible or (2) allows you to make a
1266
     \hookrightarrow final guess.
1267
1268
     Please follow these general instructions:
1269
1270
     * You MUST ask an open-ended question.
     * You CANNOT ask the oracle for the name of the city or country.
1271
     * If you are guessing the city directly, you MUST ask a question in the
1272
     \hookrightarrow format "Is the city {your_guess}?". your_guess must be one of the cities
1273
     \,\hookrightarrow\, that the oracle can choose from.
1274
     * Before you ask the question, you MUST intelligently reason about the
1275
     \hookrightarrow history of previous questions and answers to ask the most informative
1276
     \hookrightarrow question.
1277
     * Consult the history of previous questions and answers to avoid repeating
1278
1279
     \hookrightarrow the same question.
     * Do NOT repeat the same question. It's going to yield the same result.
1280
     * Do NOT get stuck on one idea and try to branch out if you get stuck.
1281
     * Do NOT ask questions formatted like "Can you ..."
1282
     * If you are certain you know the answer, MAKE A GUESS
1283
1284
     Below is the history of previous questions and answers:
1285
     {{ observation_action_history }}
1286
1287
     You MUST generate a response in the following format. Please issue only a
1288
     \hookrightarrow single question at a time.
1289
1290
1291
     REASON:
1292
     Rationale for what question to ask next based on the previous history and
1293
     \hookrightarrow the list of cities the oracle can choose from.
1294
1295
     QUESTION:
     The question to be asked. The question must be open ended.
1296
1297
     {% elif mode == 'input_final' %}
1298
     You are an intelligent player playing a game where you must guess where
1299
1300
     \hookrightarrow someone is from. Your objective is to ask the minimal number of open-
     \hookrightarrow ended questions that are diverse in order to accurately guess their city
1301
     \hookrightarrow of origin. You can ask 10 questions to determine where they are from.
1302
1303
     \hookrightarrow You are smart, so you will ask the question that will narrow down the
1304
     \hookrightarrow possible cities the person is from as much as possible.
1305
    The cities that the oracle can choose from are:
1306
    {{ all_city_list }}
1307
```

```
1308
     Your goal is to generate an open-ended question that (1) helps you narrow
1309
     \hookrightarrow down the possible cities as much as possible or (2) allows you to make a
1310
     \hookrightarrow final guess.
1311
1312
    Please follow these general instructions:
1313
1314
     * You MUST ask an open-ended question.
     * You CANNOT ask the oracle for the name of the city or country.
1315
     * If you are guessing the city directly, you MUST ask a question in the
1316
     \hookrightarrow format "Is it {your_guess}?". your_guess must be one of the cities that
1317
     \hookrightarrow the oracle can choose from.
1318
     * Before you ask the question, you MUST intelligently reason about the
1319
     \hookrightarrow history of previous questions and answers to ask the most informative
1320
     \hookrightarrow question.
1321
     * Consult the history of previous questions and answers to avoid repeating
1322
     \hookrightarrow the same question.
1323
     * Do NOT repeat the same question. It's going to yield the same result.
1324
     * Do NOT get stuck on one idea and try to branch out if you get stuck.
1325
1326
     Below is the history of previous questions and answers:
1327
     {{ observation_action_history }}
1328
1329
     You have already asked 9 questions, so this is your final guess. Your goal
1330
     \hookrightarrow is to consult the list of cities that the orcale can choose from and
1331
     \hookrightarrow guess the city directly.
1332
1333
     You MUST generate a response in the following format. Please issue only a
1334
1335
     \hookrightarrow single question at a time.
1336
     REASON:
1337
     Rationale for what city to guess based on the previous history. In your
1338
     \hookrightarrow reason, consult the list of cities that the oracle can choose from to
1339
     \hookrightarrow precisely state VERBATIM what you will guess.
1340
1341
     QUESTION:
1342
     The question to be asked. The question must be a yes-no question in the
1343
     \hookrightarrow format "Is the city {your_guess}?". your_guess must be one of the cities
1344
1345
     \hookrightarrow that the oracle can choose from and be based on what you know about the
1346
     \hookrightarrow city so far.
1347
     {% elif mode == 'output' %}
1348
     REASON:
1349
     {{ reason }}
1350
     QUESTION:
1351
     {{ action }}
1352
1353
     {% elif mode == 'output_no_reason' %}
1354
     QUESTION:
1355
     {{ action }}
1356
     {% endif %}
1358
```

#### **1359** G Domain: CarDealer Details

#### 1360 G.1 Environment Setup

We improve the original CarDealer environment in LMRL [1] to include tool use and simulated users whose preferences/constraints evolve in the interaction. There are 10 steps to complete the task. At each step, the agent must perform a two-part action: making API calls and generating response to theuser.

- Action (Part 1): API calls. At each step, the agent must choose one of the following API calls:
- 1366 search\_car\_by\_brand\_type(car\_brand:str, car\_type:str)
- 1367 search\_car\_by\_brand(car\_brand:str)
- 1368 search\_car\_by\_type(car\_type:str)
- search\_car\_that\_have\_features(features\_list:List[str]). Note that this finds the cars that have at least all the features specified in the features\_list. These cars could have additional features.
- 1372 no\_op()

To limit the input length to the agent, we only show a maximum of 8 cars. Each shown car has information about its brand, type, features, market price (MSRP), and suggested discounts that the agent can use.

Action (Part 2): Reply to Buyer. After the API calls, the agent has access to a list of cars from the database. The agent must generate a reply to the buyer and select a car from the list of cars if it wants to propose a car to the user. Because we empirically observe that open models, such as Llama-3.2-3B-Instruct, has a tendency to hallucinate cars not in the database, we also require the agent to copy down all the information about the car that it is proposing. When negotiating with the user, the agent can offer discount, but they can at most give 10% discounts.

Simulated Users and Data Split. We use Qwen2.5-14B-Instruct, a more powerful model, to simulate the user because it must roleplay as user with significantly different constraints and preferences. Each user has begins with an ideal car (that often does not exist in the database), and they gradually reveal more information about their hard constraints as the agent talks with the user.

The training set and the validation set has the same types of users with randomly generated ideal car and budget. The validation set uses a different set of car brand and car type compared to the training set. The user types are:

- 1389 1. They will only buy a car that matches their ideal car's brand and type. They require the 1390 seller to give them a least one discount, but they are ok with going slightly above budget.
- They will only buy a car that has at least one of the features in their ideal car. They must be under budget. They are impatient during negotiation, and they will terminate the conversation immediately if the seller takes too long finding a car/price they like.
- They will only buy a car that matches their ideal car's brand and type. They are flexible with
   their budget. They are distrustful, so they will never accept the first car suggested by the
   seller.
- 1397 The test set has a different set of users with the following types:
- They will only buy a car that matches their ideal car's brand. They are flexible with their
   budget. They are impatient during negotiation, and they will terminate the conversation
   immediately if the seller takes too long finding a car/price they like.
- 1401
  1401
  2. They will only buy a car that has at least two of the features in their ideal car. They must be under budget. They are distrustful, so they will never accept the first car suggested by the seller.
- They will only buy a car that matches their ideal car's type. They want an expensive car.
   They are impatient with how many cars the seller show to them, and they will terminate the conversation immediately if the seller takes too long.
- 1407 **Reward Function.** All non-terminating steps receive a reward of 0.
- <sup>1408</sup> If the buyer agrees to buy a car, the reward function first verify that the transaction is valid:
- 1409 1. The car sold is in the database.
- 1410 2. The seller has not offered over 10% discount.

1411 3. The car satisfies all the buyer's preferences/constraints.

If the transaction is invalid, the agent receives a reward of 0. Otherwise, if the transition is valid, the reward is  $(purchase\_cost)^2/(budget \times market\_price)$ .

1414 If the buyer has not agreed to buy anything after the final step, the agent receives a negative reward of 1415 -(budget - market\_price)/market\_price.

```
1416 G.2 Agent Prompt
```

```
1417 G.2.1 API Call
```

1418

{% if mode == 'input' %} 1419 You are roleplaying as a seller in a car dealership. You are talking to a 1420  $\hookrightarrow$  buyer, and your objective is to call the database APIs so that you can 1421  $\hookrightarrow$  get information from the database to answer the user's question. 1422 1423 ### Car information 1424 Here is all the possible car brands in the database: 1425 {{ all\_car\_brands }} 1426 1427 Here is all the possible car types in the database: 1428 {{ all\_car\_types }} 1429 1430 Here is all the possible car features in the database: 1431 {{ all\_car\_features }} 1432 1433 ### API calls that you can use 1434 All the API calls (except no-op) will find cars that satisfy your specified 1435  $\hookrightarrow$  criteria. Each car will have information about its brand, type, 1436 1437  $\hookrightarrow$  features, and estimated car price (msrp). You MUST specify the criteria  $\hookrightarrow$  in the following format: 1438 1. search\_car\_by\_brand\_type 1439 This will search for all the cars that satisfy both the "API BRAND" and the 1440  $\hookrightarrow$  "API TYPE" in the database. 1441 1442 API NAME: 1443 search\_car\_by\_brand\_type 1444 API BRAND: 1445 { TODO: name of the brand that you want to search for } 1446 API TYPE: 1447 { TODO: name of the type that you want to search for } 1448 1449 **API FEATURES:** 1450 [] 1451 1452 2. search\_car\_by\_brand This will search for all the cars that satisfy the "API BRAND" in the 1453  $\hookrightarrow$  database. 1454 1455 API NAME: 1456 search\_car\_by\_brand 1457 API BRAND: 1458 { TODO: name of the brand that you want to search for } 1459 API TYPE: 1460 1461 None API FEATURES: 1462 [] 1463 1464 3. search\_car\_by\_type 1465

```
This will search for all the cars that satisfy the "API TYPE" in the
1466
     \hookrightarrow database.
1467
1468
     API NAME:
1469
     search_car_by_type
1470
     API BRAND:
1471
1472
     None
     API TYPE:
1473
     { TODO: name of the type that you want to search for }
1474
    API FEATURES:
1475
     Г٦
1476
1477
     4. search_car_that_have_features
1478
     This will search for all the cars that have the features in the database.
1479
     \hookrightarrow Features must be a list of strings with double quotes around each
1480
1481
     \hookrightarrow feature.
1482
     API NAME:
1483
     search_car_that_have_features
1484
     API BRAND:
1485
    None
1486
     API TYPE:
1487
     None
1488
     API FEATURES:
1489
     ["feature1", "feature2", ...]
1490
1491
     4. no_op
1492
     If you don't need any information from the database, you can do nothing.
1493
1494
     API NAME:
1495
     no_op
1496
     API BRAND:
1497
     None
1498
     API TYPE:
1499
     None
1500
     API FEATURES:
1501
     []
1502
1503
     ### Your goal and instructions
1504
     Your goal is to use database API to get information about cars in your
1505
     \hookrightarrow dealership. You will decide what API to call based on the chat history
1506
     \hookrightarrow and the previous API call that you have done.
1507
1508
1509
     Please follow these general instructions:
     * You MUST follow the API call format above.
1510
1511
     * You MUST ask about specific brand and/or type of car based on the buyer's
1512
     \hookrightarrow request in the chat history.
     * If the previous api request and response already answer the buyer's
1513
     \hookrightarrow request (e.g., the buyer is just negotiating prices), you can choose
1514
1515
     \hookrightarrow no_op since you already have all the necessarry information.
1516
     * If the buyer is asking for a brand of car, you can choose
1517
     \hookrightarrow search_car_by_brand.
     * If the buyer is asking for a type of car, you can choose
1518
     \hookrightarrow search_car_by_type.
1519
     * If the buyer is asking for a car that has certain features, you can
1520
     \hookrightarrow choose search_car_that_have_features. Note that this will show all the
1521
     \hookrightarrow cars that have these features, and they might contain other features as
1522
    \hookrightarrow well.
1523
```

\* If you are not able to find a car that have all the features that the 1524  $\hookrightarrow$  buyer wants, you MUST choose to search on a subset of the features 1525  $\hookrightarrow$  mentioned by the buyer, still using the search\_car\_that\_have\_features 1526  $\hookrightarrow$  API. 1527 \* If the buyer has not made any request it, you can choose no\_op. 1528 \* DO NOT repeat the exact same API call as the previous one. It will NOT 1529  $\hookrightarrow$  lead to a better outcome. 1530 1531 Below is the history of the conversation so far: 1532 {{ observation\_action\_history }} 1533 1534 All previous {{ past\_N }} API calls (include no-ops): 1535 {{ prev\_api\_call\_history }} 1536 1537 Previous API call (that is not no-op): 1538 {{ previous\_api\_call }} 1539 Database's response to the previous API call: 1540 {{ previous\_api\_response }} 1541 1542 ### Output format (You MUST ALWAYS have 4 fields: REASON, API NAME, API 1543  $\hookrightarrow$  BRAND, API TYPE, API FEATURES. You MUST begin directly with the REASON 1544  $\hookrightarrow$  field) 1545 **REASON:** 1546 Rationale for what API call to make based on the previous history. In your 1547  $\hookrightarrow$  reason, consult the list of API calls that you can make to precisely 1548  $\hookrightarrow$  state VERBATIM what you will do. 1549 API NAME: 1550 name of the API call that you will make 1551 API BRAND: 1552 brand of the car that you will search for (else None) 1553 API TYPE: 1554 type of the car that you will search for (else None) 1555 API FEATURES: 1556 A list of features of the car that you will search for. Each feature must 1557  $\hookrightarrow$  be a string with double quotes around it. (else []) 1558 {% elif mode == 'output' %} 1559 **REASON:** 1560 1561 {{ reason }} API NAME: 1562 {{ api\_name }} 1563 1564 API BRAND: {{ api\_brand }} 1565 API TYPE: 1566 {{ api\_type }} 1567 API FEATURES: 1568 1569 {{ api\_features }} {% endif %} 1579

#### 1572 G.2.2 Reply to Buyer

```
{% if mode == 'input' %}
1574
     You are roleplaying as a seller in a car dealership. You are talking to a
1575
     \hookrightarrow buyer, and your objective is to get the buyer to buy the car from you
1576
1577
     \hookrightarrow with as high price as possible.
1578
    ### Car information
1579
    Here is all the possible car brands in the database:
1580
     {{ all_car_brands }}
1581
1582
```

Here is all the possible car types in the database: 1583 {{ all\_car\_types }} 1584 1585 ### Your goal and instructions 1586 You goal is to (1) You need to output the response to talk to the buyer, 1587  $\hookrightarrow$  which can be getting information about what type of car they want, 1588 1589  $\hookrightarrow$  discussing what car your dealership has, or negotiating the price of the  $\hookrightarrow$  car. (2) You need to select one car from the list of cars that you look 1590  $\hookrightarrow$  up that you are proposing to the buyer. 1591 1592 Please follow these general instructions: 1593 \* You MUST pay close attention to the chat history and what you looked up 1594  $\hookrightarrow$  in the database (if any) to decide what to say to the buyer. 1595 \* If you have not looked up any car yet (maybe you just started the 1596  $\hookrightarrow$  conversation), you should ask the buyer what type of car they want. 1597 1598 \* If you have already looked up some cars, you should discuss ONLY ONE car 1599  $\hookrightarrow$  with the buyer. Ideally, this should be the most expensive car that  $\hookrightarrow$  satisfy the buyer's request. In addition to writing down the car index, 1600  $\hookrightarrow$  you MUST also copy down the car information and write it in the ' 1601  $\hookrightarrow$  PROPOSED CAR BRAND', 'PROPOSED CAR TYPE', 'PROPOSED CAR FEATURES', and ' 1602  $\hookrightarrow$  PROPOSED CAR MSRP' fields. 1603 \* You MUST NOT make up a car that is not in the database. You MUST only 1604  $\hookrightarrow$  mention cars that are in the database. 1605 \* Some car brand do not have the type of car that the buyer wants. If that 1606  $\hookrightarrow$  happens, you should suggest the type of car that the brand has or 1607  $\hookrightarrow$  suggest another brand. 1608 \* Unless the buyer explicitly asks for a discount, you should NOT offer a 1609  $\hookrightarrow$  discount. You should just state the market price. 1610 \* If the buyer asks for a discount, you MUST start by a discount that is 1611  $\hookrightarrow$  less than 10% discount (for example, 2% discount). If the buyer is not 1612  $\hookrightarrow$  satisfied, you can try offering a better discount. You should aim to 1613  $\hookrightarrow$  offer as little discount as possible because you want to maximize your 1614  $\hookrightarrow$  profit. The maximum discount you can offer is 10% discount. 1615 \* If the discounted car is still not satisfying the buyer's request, you 1616  $\hookrightarrow$  can try suggesting another similar car. 1617 \* Your dealership only has the cars that are in the database. You cannot 1618  $\hookrightarrow$  add additional features to the car. 1619 1620 ### Input 1621 Below is the history of the conversation so far: 1622 {{ observation\_action\_history }} 1623 1624 API call (that is not no-op): 1625 {{ api\_call }} 1626 Database's response to the API call: 1627 {{ api\_response }} 1628 1629 Here is the most recent message from the buyer: 1630 {{ buyer\_response }} 1631 1632 1633 ### Output format (You MUST always have 3 fields: REASON, RESPONSE, CAR  $\hookrightarrow$  INDEX, PROPOSED CAR BRAND, PROPOSED CAR TYPE, PROPOSED CAR FEATURES, 1634  $\hookrightarrow$  PROPOSED CAR MSRP. You MUST begin directly with the REASON field) 1635 1636 REASON: Rationale for what reply you will give to the buyer and which car you will 1637  $\hookrightarrow$  propose. 1638 **RESPONSE:** 1639 Your 1-5 sentence response to the buyer. 1640 CAR INDEX: 1641

The index of the car you are proposing to the buyer (0 if you haven't 1642  $\hookrightarrow$  looked up any car yet). 1643 PROPOSED CAR BRAND: 1644 Copy down the brand of the car you are proposing to the buyer via the CAR 1645  $\hookrightarrow$  INDEX. (None if you haven't looked up any car yet) 1646 PROPOSED CAR TYPE: 1647 Copy down the type of the car you are proposing to the buyer via the CAR 1648  $\hookrightarrow$  INDEX. (None if you haven't looked up any car yet) 1649 PROPOSED CAR FEATURES: 1650 Copy down the features of the car you are proposing to the buyer via the 1651  $\hookrightarrow$  CAR INDEX. MUST be a list of strings. ([] if you haven't looked up any 1652  $\hookrightarrow$  car yet) 1653 PROPOSED CAR MSRP: 1654 Copy down the msrp (and integer) of the car you are proposing to the buyer 1655  $\hookrightarrow$  via the CAR INDEX. (O if you haven't looked up any car yet) 1656 {% elif mode == 'output' %} 1657 **REASON:** 1658 {{ reason }} 1659 **RESPONSE:** 1660 {{ response }} 1661 CAR INDEX: 1662 {{ car\_idx }} 1663 PROPOSED CAR BRAND: 1664 {{ proposed\_car\_brand }} 1665 PROPOSED CAR TYPE: 1666 {{ proposed\_car\_type }} 1667 PROPOSED CAR FEATURES: 1668 {{ proposed\_car\_features }} 1669 PROPOSED CAR MSRP: 1670 {{ proposed\_car\_msrp }} 1671 {% endif %} 1673

#### 1674 G.3 Summary Prompt

```
{% if system %}
1676
     You are an intelligent assistant summarizing a conversation between a
1677
     \hookrightarrow seller and a buyer, where the seller is trying to sell a car to the
1678
     \hookrightarrow buyer.
1679
1680
     ## Overall information about the game
1681
     Here is all the possible car brands in the database:
1682
     {{ all_car_brands }}
1683
1684
1685
     Here is all the possible car types in the database:
1686
     {{ all_car_types }}
1687
     ## What you receive as input
1688
     You are given (1) a history of what the seller has said and done and (2)
1689
     \hookrightarrow failure reason if the seller failed to sell the car. At each step of the
1690
     \hookrightarrow conversation, the seller first look up a car in the dataset. Then, the
1691
     \hookrightarrow seller will examine what is available in the dataset and potentially
1692
     \hookrightarrow suggest a car to the buyer based on the buyer's request.
1693
1694
1695
    ## Your goal and rules to follow
    You must summarize the chat history by answering the questions below.
1696
    1. What kind of requirements are not negotiable for the buyer? What does
1697
    \hookrightarrow the buyer care about the most? You MUST not ignore what the buyer said
1698
    \mapsto in the 1st step because they are just describing their ideal car. You
1699
```

 $\hookrightarrow$  MUST pay attention to what the buyer keeps repeating in the remaining 1700  $\hookrightarrow$  conversation as those requirements are non-negotiable. If features are 1701  $\hookrightarrow$  not negotiable, did the seller want at least some features or all the 1702  $\hookrightarrow$  features? 1703 2. What is the mood of buyer? Are they neutral, excited, impatient, etc.? 1704 1705  $\hookrightarrow$  Were they rushing the negotiation? Were they dubious about the car? 1706 3. What type of API calls are made? Why did the seller make these API calls 1707  $\hookrightarrow$  based on the conversation history? What are the results of the API 1708 1709  $\hookrightarrow$  calls? 4. If the API is searching for car based on features ' 1710  $\hookrightarrow$  search\_car\_that\_have\_features', and 'api\_features' contains multiple 1711  $\hookrightarrow$  features, the API will look for cars that have ALL the listed features. 1712  $\hookrightarrow$  Did the seller just look up one feature, or multiple features? You MUST 1713  $\hookrightarrow$  be very specific what features the seller was looking for. 1714 1715 5. What type of car have the seller suggested to the buyer? If the API does 1716  $\hookrightarrow$  not find any car, there will be no car under "Car in the database" 1717  $\hookrightarrow$  section. You MUST note down if the API call fails to find any car. 1718  $\hookrightarrow$  However, the seller might make up a car that is not in the database 1719  $\hookrightarrow$  under the "Copied car information" section. If the seller suggest a car 1720  $\hookrightarrow$  that is not in the database, you must include that in your summary. 1721 6. Did the seller suggest the most expensive car that satisfies the buyer's 1722  $\hookrightarrow$  requirements? Did the seller offers a discount? 1723 7. When the seller reject a car suggested by the seller, what is the reason 1724  $\hookrightarrow$  ? 1725 1726 1727 8. What happened at the last conversation step? Did the buyer end the  $\hookrightarrow$  negotiation? If so, what is the reason? Was there any car found in the 1728  $\hookrightarrow$  database? It is problematic if the "Car suggested in the database" 1729  $\hookrightarrow$  section is empty. 1730 1731 9. If the final transaction is successful, what is the price of the car? 1732 10. If the final transaction is not successful, you must carefully explain 1733  $\hookrightarrow$  what the seller said at the end, and you might also examine the failure 1734  $\hookrightarrow$  reason and include that in your summary. 1735 1736 1737 ## Output format You must generate the summary as the following markdown format. You must 1738  $\hookrightarrow$  answer the questions above in the right section, and you can add other 1739  $\hookrightarrow$  information that you think is important. 1740 ### Summary of the buyer 1741 { You MUST include your answer to 1. and 2. here. } 1742 1743 ### Summary of the API calls 1744 1745 { You MUST include your answer to 3. and 4. here. } 1746 1747 ### Summary of the car suggested { You MUST include your answer to 5. and 6. and 7. here} 1748 1749 1750 ### Summary of what happened at the last step { You MUST include your answer to 8. here. } 1751 1752 ### Summary of the overall transaction (success or failure) 1753 { You MUST include your answer to 9. and 10. here. } 1754 {% endif %} 1755 {% if not system %} 1756 {% if mode == 'input' %} 1757 ## Chat History 1758

```
1759 {{ chat_history }}
1760
1761 ## Failure Reason (if any)
1762 {{ failure_reason }}
1763 {% endif %}
1765 {% endif %}
```

#### 1766 G.4 Hindsight Proposer Prompt

```
1767 G.4.1 API Call
```

```
1768
     {% if system %}
1769
     You are an intelligent teacher who gives guidance on a seller who is trying
1770
     \hookrightarrow to look up a car in the database in response to conversation with a
1771
     \hookrightarrow buyer. The seller only has 10 steps to sell the car to the buyer.
1772
1773
     ## What you receive as input
1774
     You are given (1) a list of car brands, types, and features in the
1775
     \hookrightarrow dealership; (2) the chat history of between the buyer and the seller;
1776
     \hookrightarrow (3) the last {{ past_N }} API calls made by the seller and whether the
1777
     \hookrightarrow API calls have found any car in the database; (4) the most recent no-op
1778
1779
     \hookrightarrow API call, which affects what car the seller can propose to the buyer.
1780
     To further help you make wise judgements and provide helpful guidance to
1781
     \hookrightarrow the player, you are also given (5) a summary of the complete interaction
1782
     \hookrightarrow between the buyer and the seller (which include what the buyer cares
1783
     \hookrightarrow about, what API has the seller made, what car the seller has proposed,
1784
     \hookrightarrow and whether the seller has successfully sold the car to the buyer in the
1785
     \hookrightarrow end.
1786
1787
1788
     ## Your goal and rules to follow
     Your goal is to use your hindsight reasoning ability to generate {{
1789
     \hookrightarrow num_responses}} alternative API calls that the seller should have made
1790
     \hookrightarrow given the current chat history. Your process is to: 1. reason about all
1791
     \hookrightarrow the information (including the summary); 2. generate some API calls that
1792
     \hookrightarrow the seller could have made; 3. generate some plausible reasoning that
1793
     \hookrightarrow the seller could have come up with based on the current chat history.
1794
1795
     ### API calls that the seller can make
1796
     All API calls are in the json format. All the API calls (except no-op) will
1797
     \hookrightarrow find cars that satisfy your specified criteria. Each car will have
1798
     \hookrightarrow information about its brand, type, features, and estimated car price (
1799
     \hookrightarrow msrp).
1800
1801
     1. search_car_by_brand_type
1802
     This will search for all the cars that satisfy both the 'api_brand' and the
1803
     \hookrightarrow 'api_type' in the database. You must format you API call as the
1804
     \hookrightarrow following:
     '''json
1805
     {
1806
         "api_name": "search_car_by_brand_type",
1807
         "api_brand": "{ TODO: name of the brand that you want to search for }",
1808
         "api_type": "{ TODO: name of the type that you want to search for }",
1809
         "api_features": [] # THIS MUST BE EMPTY
1810
     }
1811
     ...
1812
1813
     2. search_car_by_brand
    This will search for all the cars that satisfy the 'api_brand' in the
1814
     \hookrightarrow database. You must format you API call as the following:
1815
     ''json
1816
```

```
1817
     {
         "api_name": "search_car_by_brand",
1818
         "api_brand": "{ TODO: name of the brand that you want to search for }",
1819
         "api_type": "",
1820
         "api_features": [] # THIS MUST BE EMPTY
1821
1822
     }
     ...
1823
     3. search_car_by_type
1824
     This will search for all the cars that satisfy the 'api_type' in the
1825
     \hookrightarrow database. You must format you API call as the following:
1826
     '''json
1827
     {
1828
         "api_name": "search_car_by_type",
1829
         "api_brand": "",
1830
         "api_type": "{ TODO: name of the type that you want to search for }",
1831
1832
         "api_features": [] # THIS MUST BE EMPTY
1833
     }
     4. search_car_that_have_features
1834
     This will search for all the cars that have the features in the database.
1835
     \hookrightarrow You must format you API call as the following:
1836
     '''json
1837
     {
1838
         "api_name": "search_car_that_have_features",
1839
         "api_brand": "",
1840
         "api_type": "",
1841
         "api_features": ["feature1", "feature2", ...]
1842
     7
1843
     ...
1844
1845
     4. no_op
     If you don't need any information from the database, you can do nothing.
1846
     ""
        json
1847
     {
1848
         "api_name": "no_op",
1849
         "api_brand": ""
1850
         "api_type": "",
1851
         "api_features": []
1852
     }
1853
     ...
1854
1855
     ### Rules that you must follow
1856
     - **Make feasible API calls:** The API call that you make MUST be a
1857
     \hookrightarrow feasible API call that the seller can make given the current chat
1858
     \hookrightarrow history.
1859
     - **In teacher_reason, reason with all the information:** You are the
1860
     \hookrightarrow teacher. In your teacher reasoning, you MUST make reference what is
1861
1862
     \hookrightarrow discussed in the chat history. You MUST generate API calls that are
     \hookrightarrow possible and reasonable to make given the current chat history. You can
1863
     \hookrightarrow make use of the "### Summary of the Entire Negotation" to help you
1864
     \hookrightarrow identify better but STILL FEASIBLE API calls to make given the current
1865
1866
     \hookrightarrow chat history.
1867
     - **In seller_reason, reasoning should not include secret information from
     \hookrightarrow the summary:"** You are generating what is the possible reasoning that a
1868
     \hookrightarrow seller can have based on the chat history in order to generate the API
1869
     \hookrightarrow call. It must be a feasible reasoning based on the chat history, all
1870
     \hookrightarrow previous {{past_N}} API calls, Previous API call, and the "Database's
1871
     \hookrightarrow response to the previous API call". You MUST NOT talk about the summary,
1872
        which the seller does not have access to.
     \hookrightarrow
1873
1874
```

```
Here are some information about the dealership to help you make better API
1875
     \hookrightarrow calls:
1876
     - You must always make sure that the "Database's reseponse to the previous
1877
     \hookrightarrow API call" is not empty. Even if you have made the same API call in "All
1878
     \hookrightarrow previous {{ past_N }} API calls", if the "Database's response to
1879
     \hookrightarrow previous API call" is currently empty, you must make the same API call
1880
1881
     \hookrightarrow that would give you the right set of cars requested by the buyer in the
     \hookrightarrow chat history.
1882
     - In the chat history, if the buyer has revealed what are their non-
1883
     \hookrightarrow negotiable requirements (you can get help from the summary), you should
1884
     \hookrightarrow make the API call to search for calls that satisfy the buyer's non-
1885
     \hookrightarrow negotiable requirements. However, in the seller_reason, you MUST make
1886
     \hookrightarrow sure that you only talk about what the seller discussed in the chat
1887
     \hookrightarrow history, not the summary.
1888
     - If "Database's response to the previous API call" has cars matching the
1889
     \hookrightarrow buyer's non-negotiable requirements, you can just make a no-op API call.
1890
     - Sometimes, "Database's response to the previous API call" might have too
1891
     \hookrightarrow many cars. You can try to make a 'search_car_by_brand_type' API call to
1892
     \hookrightarrow narrow down the set of cars.
1893
     - 'search_car_that_have_features' will return cars that have ALL the
1894
     \hookrightarrow features in the 'api_features' list. If there are no cars that have all
1895
     \hookrightarrow the features, you can try searching FOR A SUBSET of the features.
1896
1897
     ## Output format
1898
     The output is a list of JSON containing {{num_responses}} different pairs
1899
     \hookrightarrow of reasoning and feasible API calls that you would have made.
1900
     '''json
1901
     [
1902
         {
1903
              "teacher_reason": "string: your rationale for the API call that you
1904
              \hookrightarrow think the seller should have made instead.",
1905
              "api_call": {
1906
                  "api_name": "string: the name of the API call that you will make
1907
                  \hookrightarrow ",
1908
                  "api_brand": "string: the brand of the car that you will search
1909
                  \hookrightarrow for (if applicable)",
1910
                  "api_type": "string: the type of the car that you will search
1911
1912
                  \hookrightarrow for (if applicable)",
                  "api_features": ["string: the features that you will search for
1913
                  \hookrightarrow (if applicable)"]
1914
             },
1915
              "seller_reason": "string: If you are the seller who does not know
1916
             \hookrightarrow the entire negotiation history, what will be your reasoning in
1917
              \hookrightarrow order to generate the API call based on the chat history alone?
1918
              \hookrightarrow You MUST only refer to the chat history, 'All previous {{ past_N
1919
1920
              \hookrightarrow }} API calls (include no-ops)', 'Previous API call (that is not
             \hookrightarrow no-op)', and 'Database's response to the previous API call'. You
1921
             \hookrightarrow MUST NOT talk about content in the summary, which the seller does
1922
             \hookrightarrow not have access to. ",
1923
         }
1924
1925
          . . .
1926
     1
     ""
1927
1928
     ## Overall information about the current negotiation
1929
     Here is all the possible car brands in the database:
1930
     {{ all_car_brands }}
1931
1932
    Here is all the possible car types in the database:
1933
```

```
{{ all_car_types }}
1934
1935
     Here is all the possible car features in the database:
1936
     {{ all_car_features }}
1937
1938
     ## Summary of the Entire Negotiation
1939
1940
     {{ summary }}
     {% endif %}
1941
     {% if not system %}
1942
     {% if mode == 'input' %}
1943
     Here is the chat history until step {{step_idx}}:
1944
     {{ observation_action_history }}
1945
1946
     All previous {{ past_N }} API calls (include no-ops):
1947
     {{ prev_api_call_history }}
1948
1949
1950
     Previous API call (that does not return empty response):
1951
     {{ previous_api_call }}
     Database's response to the previous API call (This affects what car the
1952
     \hookrightarrow seller can propose to the buyer, so it is important to try to keep this
1953
1954
     \hookrightarrow non-empty if possible):
     {{ previous_api_response }}
1955
1956
     At this point in the chat history, what API call would you have made? You
1957
     \hookrightarrow MUST generate {{num_responses}} API calls. The API calls should try to
1958
     \hookrightarrow be diverse, but it is ok if sometimes they are the same due to the chat
1959
     \hookrightarrow history.
1960
     {% endif %}
1961
     {% endif %}
1963
```

#### 1964 G.4.2 Rely to Buyer

1965

{% if system %} 1966 You are an intelligent teacher who gives guidance on a seller who negotiate 1967  $\hookrightarrow$  with a buyer and propose a car that they can buy. The seller only has 1968  $\hookrightarrow$  10 steps to sell the car to the buyer. 1969 1970 ## What you receive as input 1971 1972 You are given (1) a list of car brands, types in the dealership; (2) the  $\hookrightarrow$  chat history of between the buyer and the seller; (3) the most recent 1973  $\hookrightarrow$  API call with non-empty response, which affects what car the seller can 1974  $\hookrightarrow$  propose to the buyer. 1975 1976 1977 To further help you make wise judgements and provide helpful guidance to 1978  $\hookrightarrow$  the player, you are also given (4) a summary of the complete interaction 1979  $\hookrightarrow$  between the buyer and the seller (which include what the buyer cares 1980  $\hookrightarrow$  about, what API has the seller made, what car the seller has proposed,  $\hookrightarrow$  and whether the seller has successfully sold the car to the buyer in the 1981  $\hookrightarrow$  end). 1982 1983 ## Your goal and rules to follow 1984 Your goal is to use your hindsight reasoning ability to generate {{ 1985  $\rightarrow$  num\_responses} alternative responses and proposed car (if applicable) 1986  $\hookrightarrow$  that the seller should have made given the current chat history. Your 1987 1988  $\hookrightarrow$  process is to: 1. reason about all the information (including the 1989  $\hookrightarrow$  summary); 2. generate some responses and proposed car that the seller  $\hookrightarrow$  could have made; 3. generate some plausible reasoning that the seller 1990  $\hookrightarrow$  could have come up with based on the current chat history. 1991 1992

```
### Rules that you must follow
1993
     - **Generate feasible responses:** The response that you generate MUST be a
1994
     \hookrightarrow feasible response that the seller can make given the current chat
1995
     \hookrightarrow history. You MUST NOT make up a car that is not in the database in the
1996
1997
     \hookrightarrow response.
     - **Select feasible proposed car:** If "Database's response to the previous
1998
1999
     \hookrightarrow API call" is not empty, you CAN select the index of a car from the list
     \hookrightarrow to propose to the buyer. However, if that field is empty, you MUST
2000
2001
     \hookrightarrow leave the car_index as 0.
     - **Copy down the proposed car information:** If you selected a car (which
2002
     \hookrightarrow CAN ONLY happen if "Database's response to the previous API call" is not
2003
     \hookrightarrow empty", you MUST copy down the information of the car that you are
2004
     \hookrightarrow proposing to the user. You MUST NOT make up a car that is not in the
2005
     \hookrightarrow database.
2006
     - **In teacher_reason, reason with all the information:** You are the
2007
2008
     \hookrightarrow teacher. In your teacher reasoning, you MUST make reference to what is
     \hookrightarrow discussed in the chat history. You MUST generate responses that are
2009
     \hookrightarrow possible and reasonable to make given the current chat history. You can
2010
     \hookrightarrow make use of the "### Summary of the Entire Negotation" to help you
2011
     \hookrightarrow identify better but STILL FEASIBLE responses to make given the current
2012
     \hookrightarrow chat history.
2013
     - **In seller_reason, reasoning should not include secret information from
2014
     \hookrightarrow the summary:"** You are generating what is the possible reasoning that a
2015
     \hookrightarrow seller can have based on the chat history in order to generate the
2016
     \hookrightarrow response, car_idx, and proposed_car. It must be a feasible reasoning
2017
     \hookrightarrow based on the chat history alone.
2018
2019
     Here are some information about the dealership to help you make better
2020
2021
     \hookrightarrow responses and select better cars:
     - You MUST NOT make a car that is not in the database.
2022
     - You can ONLY propose a car to the user if "Database's response to the
2023
     \hookrightarrow previous API call" is not empty. You MUST NOT make up general claims
2024
     \hookrightarrow about the car that is not in the database.
2025
     - When you select a car, ideally, you should select the most expensive car
2026
     \hookrightarrow that satisfy the buyer's request.
2027
     - If the buyer asks for a discount, you MUST start by a discount that is
2028
     \hookrightarrow less than 10% discount (for example, 2% discount). If the buyer is not
2029
2030
     \hookrightarrow satisfied, you can try offering a better discount. You should aim to
     \hookrightarrow offer as little discount as possible because you want to maximize your
2031
     \hookrightarrow profit. The maximum discount you can offer is 10% discount.
2032
2033
     - Once you have exhausted the maximum discount, you MUST try suggesting
     \hookrightarrow another car based on what is the buyer's non-negotiable requirements
2034
     \hookrightarrow based on the chat history (and you can also get help from the summary).
2035
2036
     - You MUST NOT get stuck in a loop of suggesting the same car over and over
          again especially if you have already given the maximum discount and the
2037
2038
     \rightarrow
          buyer is asking for a cheaper price again.
2039
     ## Output format
2040
     The output is a list of JSON containing {{num_responses}} different pairs
2041
2042
     \hookrightarrow of reasoning and feasible responses that you would have made.
     ···json
2043
2044
     Γ
         {
2045
             "teacher_reason": "string: your rationale for the response that you
2046
             \hookrightarrow think the seller should have made instead.",
2047
             "response": "string: your 1-5 sentence response to the buyer",
2048
              "car_idx": "integer: the index of the car you are proposing to the
2049
             \hookrightarrow buyer (0 if you haven't looked up any car yet)",
2050
              "proposed_car": {
2051
```

```
"brand": "string: the brand of the car. Empty string if you
2052
                  \hookrightarrow haven't looked up any car yet",
2053
                  "type": "string: the type of the car. Empty string if you haven'
2054
                  \hookrightarrow t looked up any car yet",
2055
                  "features": ["string: the features of the car. Empty list if you
2056
2057
                 \hookrightarrow haven't looked up any car yet",
                  "msrp": "integer: the manufacturer's suggested retail price of
2058
                 \hookrightarrow the car. O if you haven't looked up any car yet"
2059
             }
2060
             "seller_reason": "string: If you are the seller who does not know
2061
             \hookrightarrow the entire negotiation history, what will be your reasoning in
2062
             \hookrightarrow order to generate the response, car_idx, and proposed_car based
2063
             \hookrightarrow on the chat history alone? You MUST only refer to the chat
2064
             \hookrightarrow history, 'Previous API call (that is not no-op)', and 'Database's
2065
             \hookrightarrow response to the previous API call'. You MUST NOT talk about
2066
2067
             \hookrightarrow content in the summary, which the seller does not have access to.
2068
             \rightarrow ",
         }
2069
2070
          . .
2071
     ""
2072
2073
     ## Overall information about the current negotiation
2074
     Here is all the possible car brands in the database:
2075
     {{ all_car_brands }}
2076
2077
     Here is all the possible car types in the database:
2078
2079
     {{ all_car_types }}
2080
     ## Summary of the Entire Negotiation
2081
     {{ summary }}
2082
     {% endif %}
2083
     {% if not system %}
2084
     {% if mode == 'input' %}
2085
     Here is the chat history until step {{step_idx}}:
2086
     {{ observation_action_history }}
2087
2088
     Previous API call (that does not return empty response):
2089
     {{ api_call_used }}
2090
     Database's response to the previous API call (This affects what car the
2091
     \hookrightarrow seller can propose to the buyer, so it is important to try to keep this
2092
     \hookrightarrow non-empty if possible):
2093
     {{ api_response_used }}
2094
2095
     At this point in the chat history, what response would you have made? You
2096
2097
     \hookrightarrow MUST generate {{num_responses}} responses. The responses should try to
2098
     \hookrightarrow be diverse.
     {% endif %}
2099
     {% endif %}
3189
```