INQUIRE: Interactive Querying for User-aware Informative REasoning

Anonymous Author(s)

Affiliation Address email

Abstract: Research on Interactive Robot Learning has yielded several modalities for querying a human for training data, including demonstrations, preferences, and corrections. While prior work in this space has focused on optimizing the robot's queries within each interaction type, there has been little work on optimizing over the selection of the interaction type itself. We present INQUIRE, the first algorithm to implement and optimize over a generalized representation of information gain across multiple interaction types. Our evaluations show that INQUIRE can dynamically optimize its interaction type (and respective optimal query) based on its current learning status and the robot's state in the world, resulting in more robust performance across tasks in comparison to state-of-the art baseline methods. Additionally, INQUIRE allows for customizable cost metrics to bias its selection of interaction types, enabling this algorithm to be tailored to a robot's particular deployment domain and formulate cost-aware, informative queries.

Keywords: Active Learning, Learning from Demonstration, Human-Robot Interaction

1 Introduction

2

3

5

10

11

12

13

14

15

16

23

24

25

26

27

30

31

As we envision robots that adapt to novel tasks and environments after deployment, it is important to consider *how* they can efficiently obtain training data to address this novelty. For robots that operate in human spaces, the people around it can provide training data. Research in Interactive Robot Learning has yielded many effective methods for obtaining training data via interaction between a robot and a human teacher. While *demonstrations* are a popular type of interaction, other research has examined robot learning from *preferences* [1, 2], *corrections* [3, 4], and *binary feedback* [5].

All of these interaction types differ according to how the agent queries the teacher, the constraints placed on the teacher's feedback, and how the agent should interpret the teacher's feedback as training data [6]. In a demonstration, the teacher provides the trajectory that the robot should take starting from a particular state. By contrast, a preference query involves the robot proposing a number of possible trajectories, from which the teacher selects one. When seeking a correction, the robot demonstrates a single trajectory that the teacher can modify in whole or in part. Binary feedback again involves the robot demonstrating a single trajectory, but it only receives a positive or negative reward. As a result, the *type* of interaction used to query the teacher influences the training data available to the learning agent and subsequently influences the robot's learning performance [6].

Prior work in Active Learning has investigated how to formulate informative queries by maximizing the expected information gain resulting from the teacher's feedback. However, this work typically assumes that the robot uses a single interaction type for all queries, and thus does not incorporate feedback obtained via multiple interaction types. Furthermore, the optimal interaction type depends on the robot's task knowledge (which changes over time), the robot's query state (i.e., its state within its environment at the time it queries the teacher), and domain-specific considerations (e.g., the time or effort required for a teacher to respond to each interaction type [7]).

Our work is motivated by this question: How can a robot optimize both the *type* and *content* of its queries to a human teacher based on what information it needs at any given moment? We intro-

duce INOUIRE: a robot learning system that performs this optimization by representing multiple interaction types in a single unified framework, enabling the robot to directly estimate and com-42 pare the expected information gain of its queries across multiple interaction types. We evaluated 43 INQUIRE against two state-of-the-art interactive learning methods that use a single or fixed pattern 44 of interaction types. We analyzed the effect of domain on INQUIRE's performance and selection 45 of interaction types over time by simulating four reward-learning problems in robotics domains. 46 We found that INQUIRE learned reward functions that were more accurate and resulted in better task performance than either baseline, with particular strength in adapting to low-information query 48 states (i.e., repeated states in which the robot has already received feedback). 49

2 Related Works

50

51

53

54

55

56

57

58

One popular approach to learning from human feedback is learning from demonstrations, either via imitation learning [8] or inverse reinforcement learning [9]. However, there is a wealth of other approaches to learning from human feedback, including learning from preferences [1, 2], labels [10], and corrections [4, 3]. Notably, these approaches optimize queries *within* a particular interaction type, typically by choosing queries that maximize volume removal [11] or maximize the information gain from the teacher's response to that query [1]. Prior work has also investigated the use of fixed strategies for selecting interaction types; for example, requesting a fixed number of demonstrations before requesting preferences for the remaining queries [12, 13]. [14] incorporates more interaction types (demonstrations, labels, and feature queries) and contributes both rule-based and decision-theoretic strategies for query selection.

The set of possible techniques to learn from human feedback is vast and varied, and several attempts 61 have been made to impose a unifying and consistent framework across them. [7] identifies four pri-62 mary interaction archetypes (Showing, Categorizing, Sorting, and Evaluating) based on how people 63 perceive them; [15] describes interactions in terms of the explicit and implicit information they convey; [6] situates interactions in the realm of human-in-the-loop learning and surveys how interaction 65 types result in different effects on a teacher's ability to provide informative feedback. Our work similarly contributes a perspective on the unifying and differentiating features of interaction types: we 67 propose a generalized framework for computing information gain across multiple interaction types. 68 We focus on four different interaction types corresponding to the archetypes in [7] and empirically 69 show the effects of dynamically selecting interaction types in robot learning.

71 3 Approach

We define a query as a set of possible choices presented to the teacher, and feedback as the teacher's selected choice in response to a query. Our goal is to enable a robot to (1) efficiently query a teacher 73 using multiple interaction types, and (2) learn from feedback obtained via these interactions. We ground this goal in the problem of learning a distribution W over feature weight vectors $\omega \in W$, 75 each resulting in a linear reward function $r(t) = \phi(t) \cdot \omega$, where $\phi(t)$ is the feature vector of a 76 77 trajectory t. Thus, our goal translates into (1) selecting queries and interaction types that minimize uncertainty over W, and (2) updating W over feedback from multiple interaction types. 79 We present INQUIRE (Alg. 1), an algorithm comprised of three key steps for each query: (1) selecting the optimal interaction type i and corresponding query q_i^* that maximizes the information 80 gain over the weight distribution W (approximated as the sample set Ω), (2) recording the teacher's 81 82 feedback to that query in a feedback set \mathbf{F} , and (3) updating the weight distribution \mathcal{W} such that it maximizes the likelihood of all feedback in F. To generalize across multiple interaction types, we 83 84 must contend with the differing formulations of query and feedback corresponding to each type. We follow the framing presented in [6], where each interaction type consists of a query space Q(s) (the 85 set of possible queries that may be posed by the agent in state s) and a choice space C(q) (the set of 86 possible teacher feedback, i.e., the choices available to the teacher in response to a query $q \in Q(s)$). For a **demonstration**, let $\mathcal{T}(s)$ represent the set of all possible trajectories originating from the initial 88 state s. The robot (implicitly) enables the teacher to demonstrate any trajectory in this set, and thus its query space is $Q(s) = \{T(s)\}$ (i.e., a single query consisting of the entire trajectory space). The 90 teacher's choice space is $C = \mathcal{T}(s)$ (any trajectory within that space). For a **preference**, the robot

queries the teacher with two trajectories $q = \{t_0, t_1 \mid t_0, t_1 \in \mathcal{T}(s)\}$ who then chooses either t_0

Table 1: Each interaction involves separate query spaces, choice spaces, and choice implications.

	Query Space $Q_i(s)$	$\begin{array}{l} \textbf{Query} \\ q \in Q_i(s) \end{array}$	Choice Space $C_i(q)$	Choice Implication $c \in C_i(q) \implies (c^+, c^-)$
Demo.	$\{T\}$	T	T	$c^+:t\in T c^-:T\setminus t$
Pref.	$T \times T$	$\{t_0, t_1\}, t_0, t_1 \in T$	$\{t_0,t_1\}$	$c^+: t \in \{t_0, t_1\}$ $c^-: \{t_0, t_1\} \setminus c^+$
Corr.	T	$t \in T$	T	$c^+:t'\in T c^-:q$
Binary	T	$t \in T$	{0,1}	$c = 0 \implies c^{+}: T \setminus q c^{-}: q$ $c = 1 \implies c^{+}: q c^{-}: T \setminus q$

Algorithm 1 INQUIRE - Overview

Input: Set of query states S

Parameters: K (# of queries), \mathcal{I} (interaction

Output: Weight vector ω^*

1: $\mathbf{F} \leftarrow \{\}$

2: $\Omega \leftarrow M$ random initial weight vectors

3: **for** K iterations **do**

 $s \leftarrow \text{next query state in } S$

 $q_i^* \leftarrow \text{generate_query}(s, \mathcal{I}, \mathbf{\Omega}) \text{ (Alg. 2)}$

 $\mathbf{F} \leftarrow \mathbf{F} \cup \{\text{query_teacher}(q_i^*)\}$

 $\Omega \leftarrow update_weights(\mathbf{F})$

8: $\omega^* \leftarrow \text{mean}(\mathbf{\Omega})$

9: return ω^*

99

100

101

103

104

Algorithm 2 INQUIRE - Generate Query

Input: s (state), \mathcal{I} (interaction types), Ω (weight samples) Output: Query q*

1: $\mathbf{T} \leftarrow \text{uniformly_sample_trajectories}(s)$

2: Compute $\mathbf{E}: \{\mathbf{E}_{t,t',\omega}, \forall t, t' \in \mathbf{T}, \omega \in \Omega\}$ (Eq. 4)

3: **for** each interaction type $i \in \mathcal{I}$ **do**

 $\mathbf{Q} \leftarrow Q_i(s)$ (See Table 1)

 $\mathbf{C} \leftarrow \{C_i(q), \forall q \in \mathbf{Q}\}$ (See Table 1)

Compute info gain matrix $\mathbf{G^{(i)}}$ from \mathbf{E} (Eq. 9)

 $q \leftarrow \arg\max_{q'} \sum_{c \in \mathbf{C}_q, \omega \in \mathbf{\Omega}} \mathbf{G}_{q', c, \omega}^{(i)}$

 $\begin{array}{l} g \leftarrow \frac{1}{\log(\lambda_i)} \sum_{c \in \mathbf{C}_q, \omega \in \mathbf{\Omega}} \mathbf{G}_{q,c,\omega}^{(i)} \\ \text{if information gain } g > g^* \text{ then} \end{array}$

9:

 $g^* \leftarrow g \\ q^* \leftarrow q$ 10:

11: {Store query with highest info. gain}

12: return q^*

or t_1 . The query space is $Q(s) = \mathcal{T}(s) \times \mathcal{T}(s)$ and the teacher's choice space is $C(q) = \{t_0, t_1\}$. For a **correction**, the robot executes one trajectory $q \in \mathcal{T}(s)$ which the teacher then modifies to a preferable behavior. The agent's query space is $Q(s) = \mathcal{T}(s)$ and the teacher's choice space is $C(q) = \mathcal{T}(s)$. For **binary feedback**, the robot executes a single trajectory $q \in \mathcal{T}(s)$, and the teacher indicates a positive or negative reward. The agent's query space is $Q(s) = \mathcal{T}(s)$ and the teacher's 97 choice space is $C(q) = \{0, 1\}.$ 98

The implication of the teacher's choice $c \in C(q)$ is a set of accepted trajectories c^+ and set of rejected trajectories c^- , which we define in Table 1 and use later to calculate information gain. Since the set of all possible trajectories originating from s (represented by $\mathcal{T}(s)$) is potentially infinite, we approximate it as the set T containing N trajectory samples originating from the state sand consisting of randomly selected actions.

3.1 **Query Optimization**

When optimizing the agent's query, our goal is to greedily select one that maximizes the agent's expected information gain over \mathcal{W} after receiving any feedback from the choice set (summarized in Alg. 2). Selecting a query involves optimizing over information gain (IG) as follows:

$$q_i^* = \underset{q \in Q_i(s)}{\arg \max} \mathbb{E}_{c|C_i(q)} \left[\text{IG}(\mathcal{W} \mid c) \right]$$
 (1)

$$= \underset{q \in Q_i(s)}{\operatorname{arg max}} \sum_{c \in C_i(q)} \sum_{w \in \Omega} \left[P(c|w) \cdot \log \frac{M \cdot P(c|w)}{\sum_{w' \in \Omega} P(c|w')} \right]$$
 (2)

where Ω contains M samples of the distribution W. The expansion from Eq. 1 to 2 follows the derivation presented in [1]; see Appendix A.1 for intermediate steps. We adopt the commonly-used Boltzmann-rational equation to define $P(c|\omega)$:

$$P(c|\omega) = \frac{\sum_{t \in c^{+}} e^{\beta \cdot \phi(t) \cdot \omega}}{\sum_{t \in c^{+} \cup c^{-}} e^{\beta \cdot \phi(t) \cdot \omega}}$$
(3)

where $\phi(t)$ returns the feature trace of the trajectory t; that is, the sum over the feature vectors of all 111 states visited in trajectory t. β is a parameter representing the expected optimality of the teacher's 112 feedback with respect to ω . We assign a value of $\beta = 20$ across all interaction types (selected 113 through empirical evaluation). 114

To minimize the computational complexity of solving for Eq. 2, we reformulate it as a series of 115 operations over a $|Q| \times |C| \times |\Omega|$ probability matrix **P**, where $\mathbf{P}_{q,c,\omega}$ represents the probability 116 (according to weight sample $\omega \in \Omega$) that the teacher will select choice c in response to query q. To 117 construct **P**, let **E** be a $N \times N \times M$ (i.e., $|T| \times |T| \times |\Omega|$) matrix representing exponentiated rewards: 118

$$\mathbf{E}_{t,t',\omega} = e^{\beta \cdot \phi(t') \cdot \omega} \qquad \Longrightarrow \qquad \left[\mathbf{E} + \mathbf{E}^{\mathbf{T}} \right]_{t,t',\omega} = e^{\beta \cdot \phi(t') \cdot \omega} + e^{\beta \cdot \phi(t) \cdot \omega} \tag{4}$$

With E in hand, we next define the probability matrices of each interaction type as follows:

$$\mathbf{P}_{q,c,\omega}^{(\text{demo})} = \left[\mathbf{E}_0 \otimes \sum_{t \in T} \mathbf{E}^{\mathbf{T}}_t \right]_{C(\omega)}$$
 (since $|Q| = 1$ for demonstrations) (5)

$$\mathbf{P}_{q,c,\omega}^{(\text{pref})} = \left[\left(\mathbf{E} \otimes (\mathbf{E} + \mathbf{E}^{\mathbf{T}}) \right)^{\mathbf{T}}, \mathbf{E} \otimes (\mathbf{E} + \mathbf{E}^{\mathbf{T}}) \right]_{c,q_0,q_1,\omega} \qquad (c \in \{0,1\} \text{ for prefs.})$$
 (6)

$$\mathbf{P}_{q,c,\omega}^{(\text{corr})} = \left[\mathbf{E} \oslash (\mathbf{E} + \mathbf{E}^{\mathbf{T}}) \right]_{q,c,\omega}$$
(7)

$$\mathbf{P}_{q,c,\omega}^{(\text{bnry})} = \left[1 - \left(\mathbf{E}_0 \oslash \alpha \sum_{t \in T} \mathbf{E}^{\mathbf{T}}_t\right), \mathbf{E}_0 \oslash \alpha \sum_{t \in T} \mathbf{E}^{\mathbf{T}}_t\right]_{c,q,\omega} \quad (c \in \{0,1\} \text{ for binary rewards}) \quad (8)$$

where \oslash represents an element-wise division of two matrices (i.e., $(\mathbf{A} \oslash \mathbf{B})_{ij} = \mathbf{A}_{ij}/\mathbf{B}_{ij}$) and α is a normalization factor such that $\sum_c \mathbf{P}_{q,c,\omega}^{(\text{bnry})} = 1$. For derivations, see Appendix A.3. The main effect of this formulation is that it enables tractable optimization over multiple interaction types by 122 sharing a common representation E. To solve for the optimal query q_i^* using interaction type i, we use $P^{(i)}$ to construct a $|Q| \times |C| \times |\Omega|$ information gain matrix $G^{(i)}$:

$$\mathbf{G}_{q,c,\omega}^{(i)} = \mathbf{P}_{q,c,\omega}^{(i)} \cdot \log \left(\frac{M \cdot \mathbf{P}_{q,c,\omega}^{(i)}}{\sum_{\omega' \in \Omega} \mathbf{P}_{q,c,\omega'}^{(i)}} \right) \qquad q_i^* = \arg\max_{q} \sum_{c,\omega} \mathbf{G}_{q,c,\omega}^{(i)}$$
(9)

We then solve for the optimal interaction type itself. To perform a *cost-weighted* optimization, with the aim of optimizing over both interaction cost and informativeness, λ_i may be set according to domain-specific cost factors (e.g., the time or mental load involved in answering a query) for each 128 interaction type.² To perform an unweighted optimization and maximize solely over the informa-129 tiveness of each query, let λ_i be a constant value over all interaction types $i \in \mathcal{I}$. 130

$$i^* = \operatorname*{arg\,max}_{i \in \mathcal{I}} \frac{1}{\log(\lambda_i)} \sum_{c,\omega} \mathbf{G}_{q_i^*,c,\omega}^{(i)} \tag{10}$$

Update Weights from Feedback 3.2 131

119

121

123

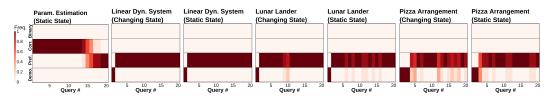
124

After presenting the optimal query to the teacher, the agent receives feedback and appends it to a 132 feedback set F—a cumulative set that contains all feedback received by the agent thus far. Our goal is to then update the weight estimate such that it maximizes the likelihood of all feedback in F:

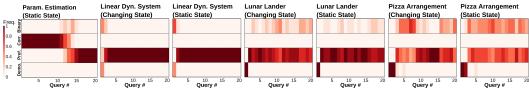
$$\omega^* = \arg\max_{\omega} \prod_{c \in \mathbf{F}} P(c|\omega) = \arg\max_{\omega} \prod_{c \in \mathbf{F}} \frac{\sum_{t \in c^+} e^{\beta \cdot \phi(t) \cdot \omega}}{\sum_{t \in c^+ \cup c^-} e^{\beta \cdot \phi(t) \cdot \omega}}$$
(11)

¹See Appendix B.1 for each domain's definition of ϕ .

²In our evaluations, we assign a cost of 20 to each demonstration, 15 to each correction, 10 to each preference, and 5 to each binary query.



(a) Selected interaction types without cost-weighting



(b) Selected interaction types with cost-weighting

Figure 1: Heatmaps illustrating how INQUIRE selects different interaction types as it learns more over time. These selections differ when deriving unweighted (top) or cost-weighted (bottom) information gain estimations. In the cost-weighted setting (bottom), INQUIRE selects more low-cost binary queries than it does in the unweighted setting (top).

We calculate the gradient over ω by differentiating over its log-likelihood given **F**:

$$\frac{\partial \ell(\omega)}{\partial \omega_{j}} = \sum_{c \in \mathbf{F}} \left[\frac{\sum_{t \in c^{+}} \beta \cdot \phi_{j}(t) \cdot e^{\beta \cdot \phi(t) \cdot \omega}}{\sum_{t \in c^{+}} e^{\beta \cdot \phi(t) \cdot \omega}} - \frac{\sum_{t \in c^{+} \cup c^{-}} \beta \cdot \phi_{j}(t) \cdot e^{\beta \cdot \phi(t) \cdot \omega}}{\sum_{t \in c^{+} \cup c^{-}} e^{\beta \cdot \phi(t) \cdot \omega}} \right]$$

$$= \sum_{c \in \mathbf{F}} \left[\beta \cdot \phi_{j}(c_{0}^{+}) - \frac{\sum_{t \in c^{+} \cup c^{-}} \beta \cdot \phi_{j}(t) \cdot e^{\beta \cdot \phi(t) \cdot \omega}}{\sum_{t \in c^{+} \cup c^{-}} e^{\beta \cdot \phi(t) \cdot \omega}} \right]$$
(iff $|c^{+}| = 1$) (13)

$$= \sum_{c \in \mathbf{F}} \left[\beta \cdot \phi_j(c_0^+) - \frac{\sum_{t \in c^+ \cup c^-} \beta \cdot \phi_j(t) \cdot e^{\beta \cdot \phi(t) \cdot \omega}}{\sum_{t \in c^+ \cup c^-} e^{\beta \cdot \phi(t) \cdot \omega}} \right]$$
 (iff $|c^+| = 1$) (13)

After receiving feedback from each query and updating \mathbf{F} , we approximate Ω by randomly initializing and then performing gradient ascent on each weight sample $\omega \in \Omega$.

Results

138

139

140

141

142

143

144

145

146

147

148

150

151

152

153

154

155

156

157

158

159

We simulate four types of learning problems in robotics using an oracle teacher to obtain controlled evaluations. The oracle teacher, similar to INQUIRE, requires its own set of trajectory samples T'. It then selects a response to a query via one of three mechanisms: returning the highest-reward trajectory from its choice space (demonstrations/preferences), rejection sampling of trajectories followed by selection of the trajectory with the highest reward-to-distance ratio from the queried trajectory (corrections), and returning whether a query meets or exceeds a reward threshold (binary feedback). Implementation details can be found in Appendix B.2.

The Parameter Estimation domain involves directly estimating a randomly-initialized, ground truth weight vector ω^* containing 8 parameters. The **Linear Dynamical System** domain, inspired by [1], simulates a controls problem and involves learning 8 parameters. The **Lunar Lander** domain simulates a controls problem involving 4 parameters. The Pizza Arrangement domain simulates a preference-learning problem involving 4 parameters. Each domain (except for Parameter Estimation) has a static-state and changing-state condition indicating whether the robot must formulate all queries from the same query state or not, respectively. For the full evaluation procedure and oracle implementation details for each domain see Appendix B.

4.1 INQUIRE Query Selection

We first analyze how INOUIRE selects queries. Figure 1 reflects the changes in interaction type selected by INQUIRE over time. Figure 1a first reports these interaction selections in an unweighted query optimization setting, where all interaction types are assumed to be equally costly. In the parameter optimization domain, INQUIRE requests corrections in the first 14-18 queries and then requests preferences as the remaining queries. Demonstrations were not enabled in this domain. In all other domains, INQUIRE requests a demonstration as its first query, then immediately switches

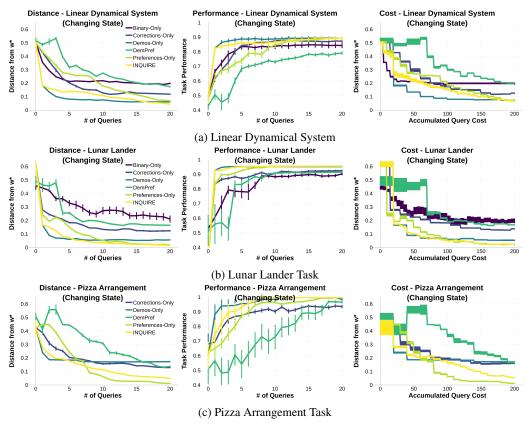


Figure 2: Metrics for the *changing state* condition in which the robot's initial state changes with each query. Error bars/regions represent variance across multiple evaluation runs with randomized query states and initial weights.

to requesting preferences for the remaining queries (occasionally alternating between preferences and demonstrations in the Lunar Lander domain).

After assigning different cost values to each interaction type, INQUIRE chooses more diverse interaction types in order to maximize its information-to-cost ratio. As shown in Figure 1b, this typically results in INQUIRE posing more binary queries due to their relatively low cost. This pivot toward binary queries may occur at the start (as seen in the linear dynamical system), middle (as seen in the parameter estimation domain), or interspersed throughout the learning process (as seen in the lunar lander domain).

4.2 Learning Performance

We now analyze the effect of INQUIRE's interaction type selections on its learning performance and compare to two types of baselines. The first, DemPref [12], learns from 3 demonstrations and then learns from preference queries by using a volume removal objective function. As our second baseline, we compare INQUIRE against agents that use only one form of interaction: demonstrations, preferences, corrections, or binary feedback. Note that the preference-only agent is formulated according to [1] and thus represents this baseline method.

We first consider the changing-state formulation of each domain, where the robot is presented with a new state for each query. Since the Parameter Estimation domain does not contain states, we exclude it from this first set of results. Figure 2 illustrates this learning performance in the Linear Dynamical System and Lunar Lander domains according to three key metrics. **Distance** measures the angular distance between the ground truth feature weights (ω^*) and the algorithm's estimated feature weight $\tilde{\omega}$ after each query. **Performance** measures the task reward achieved using a trajectory optimized according to $\tilde{\omega}$ (the algorithm's estimated feature weight after each query). Performance is scaled between 0-1, with 0 and 1 representing the worst and best possible task rewards according to ω^* ,

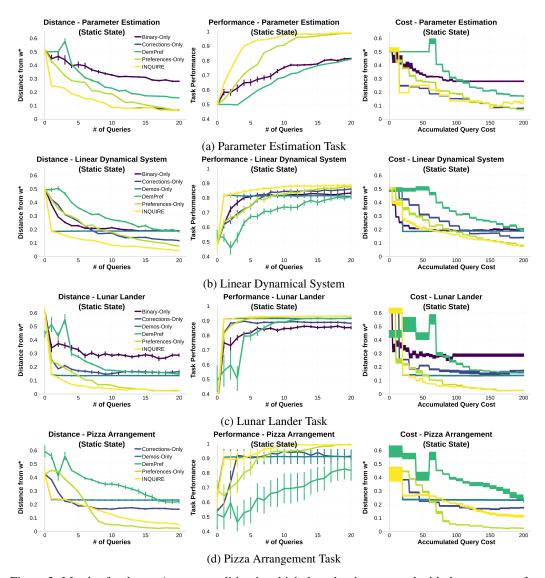


Figure 3: Metrics for the *static state* condition in which the robot is presented with the same state for all 20 queries. Error bars/regions represent variance across multiple evaluation runs with randomized query states and initial weights.

respectively. Note that INQUIRE's distance and performance metrics are achieved in the unweighted condition. **Cost-vs-Distance** measures the relationship between the cumulative cost of each query and the resulting distance between $\tilde{\omega}$ and ω^* after each query. INQUIRE's metrics in this graph are achieved in the cost-weighted condition.

Figure 3 presents the same three metrics for the *static-state* condition in which all 20 queries must be selected from the same initial state. Finally, we quantify these graphs by reporting the area-under-the-curve (AUC) metrics for the distance, performance, and cost curves across all tasks. These metrics are available in Appendix C. The AUC metrics indicate that, compared to the baseline methods, INQUIRE results in the best average learning performance (measured both by the distance and performance plots in Figures 2-3) across all domains and dominates learning performance in the static-state domains. INQUIRE also results in the best average distance-to-cost ratio across all domains.

5 Discussion

The results show the importance of dynamically selecting interaction types according to the robot's current state. For example, demonstrations can be highly informative when provided in novel states, but when the robot may only query a teacher from a single state, multiple demonstrations are likely to be very similar (if not identical). As a result, receiving multiple demonstrations in a static query state is uninformative. We see the benefits of dynamically selecting interaction types in Figure 3, where INQUIRE outperforms all single-interaction methods by optimizing both query type and content to maximize the informativeness of the query feedback.

INQUIRE selects the interaction type that, after receiving feedback, minimizes the entropy over its 204 distribution of weight estimates \mathcal{W} . This distribution thus serves as a representation of the robot's 205 current model of the task reward. Figure 1a illustrates how INOUIRE changes the query type as 206 it learns over time (represented by # of queries). This is particularly evident in the Parameter Es-207 timation task, where the algorithm originally requests corrections until it has refined its model of 208 the task reward to a point where preferences become more informative (after 14-18 queries). Over-209 all, dynamically adapting to the robot's model of the task reward results in better performance than 210 adopting a fixed strategy for selecting interaction types (i.e., DemPref, which always requests one 211 demonstration before selecting preferences). 212

An added benefit of INQUIRE is that it can incorporate a cost metric to identify cost-aware, informative queries. The AUC metrics for the cost graphs indicates that INQUIRE selects queries that, on average, minimize the cost-to-distance ratio across all domains. We expect that this cost metric is domain-specific, and can represent a number of human factors that the algorithm should take into account (e.g., the effort involved for a human to respond to each query type [6]). The cost metric used in our study thus serves as an example of how INQUIRE can factor in interaction costs.

219 6 Limitations

Our evaluation is performed using feedback from an optimal oracle. Real human feedback, however, 220 is likely to be at least somewhat sub-optimal, and its severity likely depends on the interaction type. 221 For example, a non-optimal demonstration may be one that is sufficient but not ideal for completing 222 the task. In contrast, binary rewards offer only two feedback choices to the user, and thus a non-223 optimal binary reward may indicate the opposite information from what the user intended to convey. These examples illustrate how non-optimal feedback may need to be handled differently depending 225 on the interaction type, and thus, should affect INOUIRE's estimation of information gain. Further-226 more, INQUIRE does not yet have the ability to select the state in which it queries the teacher. Prior 227 work in Active Learning has shown that state selection can improve the informativeness of resulting 228 demonstrations [16], and we expect that optimizing over the query state in addition to query type 229 and content would improve the performance of INQUIRE. 230

231 7 Conclusion

232

233

234

235

236

238 239 We introduced INQUIRE, an algorithm enabling a robot to dynamically optimize its queries and interaction types according to its task knowledge and its state within the environment. We showed that using information gain to select not just optimal queries, but optimal interaction *types*, results in consistently high performance across multiple tasks and state configurations. Future work will include formal user studies to investigate our method's efficacy with people of varied skillsets and comfort with robots; incorporation of novel interaction types and other communication modalities; and alternative representations of the reward function and feature spaces. Moreover, we are excited at the possible extensions others might present by using our open-source framework³ for evaluating and comparing active-learning agents across multiple environments and simulated teachers.

³Link removed for anonymous review

References

- [1] E. Biyik, M. Palan, N. C. Landolfi, D. P. Losey, and D. Sadigh. Asking easy questions: A
 user-friendly approach to active reward learning. In *Conference on Robot Learning (CoRL)*,
 pages 1177–1190, 2020.
- ²⁴⁵ [2] C. Wirth, R. Akrour, G. Neumann, and J. Fürnkranz. A survey of preference-based reinforcement learning methods. *The Journal of Machine Learning Research*, 18(1):4945–4990, 2017.
- [3] T. Fitzgerald, E. Short, A. Goel, and A. Thomaz. Human-guided trajectory adaptation for tool transfer. In *Intl. Conf. on Autonomous Agents and MultiAgent Systems*, pages 1350–1358, 2019.
- ²⁵⁰ [4] A. Bajcsy, D. Losey, M. O'Malley, and A. Dragan. Learning robot objectives from physical human interaction. *Proceedings of Machine Learning Research*, 78:217–226, 2017.
- [5] C. Celemin and J. Ruiz-del Solar. An interactive framework for learning continuous actions policies based on corrective feedback. *Journal of Intelligent & Robotic Sys.*, 95(1):77–97, 2019.
- Y. Cui, P. Koppol, H. Admoni, S. Niekum, R. Simmons, A. Steinfeld, and T. Fitzgerald. Understanding the relationship between interactions and outcomes in human-in-the-loop machine learning. 2021.
- P. Koppol, H. Admoni, and R. Simmons. Interaction considerations in learning from humans.
 Under Review, 2021.
- [8] A. Hussein, M. M. Gaber, E. Elyan, and C. Jayne. Imitation learning: A survey of learning methods. *ACM Comput. Surv.*, 50(2), apr 2017. ISSN 0360-0300. doi:10.1145/3054912. URL https://doi.org/10.1145/3054912.
- [9] A. Y. Ng, S. J. Russell, et al. Algorithms for inverse reinforcement learning. In *ICML*, volume 1, page 2, 2000.
- [10] C. Daniel, M. Viering, J. Metz, O. Kroemer, and J. Peters. Active reward learning. In *Robotics:* Science and Systems, 2014.
- [11] D. Sadigh, A. Dragan, S. Sastry, and S. A. Seshia. Active preference-based learning of reward functions. In RSS, 2017.
- ²⁶⁹ [12] M. Palan, N. C. Landolfi, G. Shevchuk, and D. Sadigh. Learning reward functions by integrating human demonstrations and preferences. *RSS*, 2019.
- [13] B. Ibarz, J. Leike, T. Pohlen, G. Irving, S. Legg, and D. Amodei. Reward learning from human preferences and demonstrations in atari. *NeurIPS*, 31:8011–8023, 2018.
- 273 [14] K. Bullard, A. L. Thomaz, and S. Chernova. Towards intelligent arbitration of diverse active learning queries. In *IROS*, pages 6049–6056, 2018.
- [15] H. J. Jeon, S. Milli, and A. D. Dragan. Reward-rational (implicit) choice: A unifying formalism for reward learning. In *NeurIPS*, 2020.
- [16] M. S. Lee, H. Admoni, and R. Simmons. Machine teaching for human inverse reinforcement learning. *Frontiers in Robotics and AI*, 8:188, 2021.

279 A Approach Details

280 A.1 Information Gain Derivation

$$\mathbb{E}_{c|C_i(q)}\left[\mathrm{IG}(\mathcal{W}\mid c)\right] \tag{14}$$

$$= \mathbb{E}_{\mathcal{W}, c \mid C_i(q)} \left[\log \frac{P(c \mid \mathcal{W})}{P(c)} \right]$$
 (15)

$$= \sum_{c \in C_i(q)} \left[P(c) \sum_{w \in \mathcal{W}} \left[P(w|c) \cdot \log \frac{P(c|w)}{P(c)} \right] \right]$$
 (16)

$$= \sum_{c \in C_i(q)} \left[P(c) \sum_{w \in \mathcal{W}} \left[\frac{P(w)P(c|w)}{P(c)} \cdot \log \frac{P(c|w)}{P(c)} \right] \right]$$
(17)

$$= \sum_{c \in C_i(q)} \sum_{w \in \mathcal{W}} \left[P(w)P(c|w) \cdot \log \frac{P(c|w)}{P(c)} \right]$$
(18)

$$= \sum_{c \in C_{c}(a)} \sum_{w \in \mathcal{W}} \left[P(w)P(c|w) \cdot \log \frac{P(c|w)}{\sum_{w' \in \mathcal{W}} P(w')P(c|w')} \right]$$
(19)

$$\approx \frac{1}{M} \sum_{c \in C_i(g)} \sum_{w \in \Omega} \left[P(c|w) \cdot \log \frac{M \cdot P(c|w)}{\sum_{w' \in \Omega} P(c|w')} \right]$$
 (20)

281 Standard information gain equation:

$$IG(X,y) (21)$$

$$= H(X) - H(X|y) \tag{22}$$

$$= -\mathbb{E}_{x|X}\log P(x) + \mathbb{E}_{x|X,y}\log P(x|y) \tag{23}$$

$$= \sum_{x \in X} P(x|y) \cdot \log P(x|y) - \sum_{x \in X} P(x) \cdot \log P(x)$$
 (24)

282 Best query is returned by:

$$\max_{q \in Q} \mathbb{E}_{c|C_i(q)} \left[\text{IG}(\mathcal{W}, c) \right]$$
 (25)

283 Which we evaluate as follows:

$$\mathbb{E}_{c|C_i(q)}\left[\mathrm{IG}(\mathcal{W},c)\right] \tag{26}$$

$$= H(\mathcal{W}) - \mathbb{E}_{c|C_i(q)} \left[H(\mathcal{W}|c) \right] \tag{27}$$

$$= -\mathbb{E}_{\mathcal{W}}\left[\log P(\mathcal{W})\right] + \mathbb{E}_{\mathcal{W},c|C_i(q)}\left[\log P(\mathcal{W}|c)\right]$$
(28)

$$= \mathbb{E}_{\mathcal{W}, c|C_i(q)} \left[\log P(\mathcal{W}|c) - \log P(\mathcal{W}) \right]$$
 (see proof in Sec. A.1.1) (29)

$$= \mathbb{E}_{\mathcal{W}, c|C_i(q)} \left[\log \frac{P(\mathcal{W}|c)}{P(\mathcal{W})} \right]$$
(30)

$$= \mathbb{E}_{\mathcal{W}, c \mid C_i(q)} \left[\log \frac{P(c \mid \mathcal{W})}{P(c)} \right]$$
 (by Bayes' rule) (31)

$$= \sum_{c \in C_i(q)} \left[P(c) \sum_{w \in \mathcal{W}} \left[P(w|c) \cdot \log \frac{P(c|w)}{P(c)} \right] \right]$$
(32)

$$= \sum_{c \in C_{\delta}(q)} \left[P(c) \sum_{w \in \mathcal{W}} \left[\frac{P(w)P(c|w)}{P(c)} \cdot \log \frac{P(c|w)}{P(c)} \right] \right]$$
(33)

$$= \sum_{c \in C_{c}(a)} \sum_{w \in \mathcal{W}} \left[P(w) \cdot P(c|w) \cdot \log \frac{P(c|w)}{P(c)} \right]$$
(34)

$$= \sum_{c \in C_i(q)} \sum_{w \in \mathcal{W}} \left[P(w) \cdot P(c|w) \cdot \log \frac{P(c|w)}{\sum_{w' \in \mathcal{W}} P(w') \cdot P(c|w')} \right]$$
(35)

$$\approx \frac{1}{M} \sum_{c \in C_i(q)} \sum_{w \in \Omega} \left[P(c|w) \cdot \log \frac{M \cdot P(c|w)}{\sum_{w' \in \Omega} P(c|w')} \right]$$
 (36)

Where Ω contains M samples of the distribution \mathcal{W} .

285 A.1.1 Proof of Eq. 29

$$\mathbb{E}_{\mathcal{W},c|C_i(q)}\left[\log P(\mathcal{W}|c)\right] - \mathbb{E}_{\mathcal{W}}\left[\log P(\mathcal{W})\right] \tag{37}$$

$$= \left[\sum_{w \in \mathcal{W}} P(w) \sum_{c \in C_i(q)} P(c|w) \cdot \log P(w|c) \right] - \left[\sum_{w \in \mathcal{W}} P(w) \cdot \log P(w) \right]$$
(38)

$$= \sum_{w \in \mathcal{W}} P(w) \cdot \left[\left(\sum_{c \in C_i(q)} P(c|w) \cdot \log P(w|c) \right) - \log P(w) \right]$$
(39)

$$= \sum_{w \in \mathcal{W}} P(w) \cdot \left[\sum_{c \in C_i(q)} P(c|w) \cdot \log P(w|c) - \sum_{c \in C_i(q)} P(c|w) \cdot \log P(w) \right]$$
(40)

$$= \sum_{w \in \mathcal{W}} P(w) \cdot \sum_{c \in C_i(q)} P(c|w) \cdot [\log P(w|c) - \log P(w)]$$

$$\tag{41}$$

$$= \mathbb{E}_{\mathcal{W}, c \mid C_i(q)} \left[\log P(\mathcal{W} \mid c) - \log P(\mathcal{W}) \right] \tag{42}$$

286 A.2 KL Divergence

287 Standard KL divergence equation:

$$KL(P||Q) = \sum_{x \in X} \left[P(x) \cdot \log \frac{P(x)}{Q(x)} \right]$$
 (43)

Best query is returned by:

$$\max_{q \in Q} \mathbb{E}_{c|C_i(q)} \left[\text{KL}(P(\mathcal{W}|c)||P(\mathcal{W})) \right]$$
(44)

Which we evaluate as follows:

$$\mathbb{E}_{c|C_i(q)} \mathsf{KL} \tag{45}$$

$$= \mathbb{E}_{c|C_i(q)} \left[\text{KL}(P(\mathcal{W}|c)||P(\mathcal{W})) \right]$$
(46)

$$= \mathbb{E}_{c|C_i(q)} \left[\sum_{w \in \mathcal{W}} \left[P(w|c) \cdot \log \frac{P(w|c)}{P(w)} \right] \right]$$
(47)

$$= \mathbb{E}_{c|C_i(q)} \left[\sum_{w \in \mathcal{W}} \left[P(w|c) \cdot \log \frac{P(c|w)}{P(c)} \right] \right]$$
(48)

$$= \sum_{c \in C_i(q)} \left[P(c) \sum_{w \in \mathcal{W}} \left[P(w|c) \cdot \log \frac{P(c|w)}{P(c)} \right] \right]$$
(49)

Which is equivalent to Eq. 32.

A.3 Probability Matrix Derivations 291

In the demonstration case, we can define P as follows:

$$\mathbf{P}_{q,c,\omega}^{(\text{demo})} = \frac{e^{\beta \cdot \phi(c^{+}) \cdot \omega}}{\sum_{t \in q} e^{\beta \cdot \phi(t) \cdot \omega}}$$
 (50)

$$= \frac{e^{\beta \cdot \phi(c^{+}) \cdot \omega}}{\sum_{t \in T} e^{\beta \cdot \phi(t) \cdot \omega}} \qquad \text{(since } q = T \text{ for demonstrations)}$$
 (51)

$$=\frac{\mathbf{E^{T}}_{0,c^{+},\omega}}{\sum_{t\in\mathcal{T}}\mathbf{E^{T}}_{0,t,\omega}}\tag{52}$$

$$\sum_{t \in T} \mathbf{E}^{T}_{0,t,\omega}$$

$$= \left[\mathbf{E}_{0}^{T} \otimes \sum_{t \in T} \mathbf{E}_{t} \right]_{q,c,\omega}$$
 (since there is a 1-1 correlation between c and c^{+} in demos) (53)

where \oslash represents an element-wise division of two matrices (i.e., $(\mathbf{A} \oslash \mathbf{B})_{ij} = \mathbf{A}_{ij}/\mathbf{B}_{ij}$).

In the preference case: 294

$$\mathbf{P}_{q,c,\omega}^{(\text{pref})} = \frac{e^{\beta \cdot \phi(c^{+}) \cdot \omega}}{\sum_{t \in q} e^{\beta \cdot \phi(t) \cdot \omega}}$$
 (54)

$$= \frac{e^{\beta \cdot \phi(c^{+}) \cdot \omega}}{e^{\beta \cdot \phi(q_{0}) \cdot \omega} + e^{\beta \cdot \phi(q_{1}) \cdot \omega}}$$
 (since $q = (q_{0}, q_{1})$ in preferences) (55)

Since
$$c_0 \implies c^+ = q_0$$
 and $c_1 \implies c^+ = q_1$: (56)

$$\sum_{t \in q} e^{\beta \cdot \phi(t)} \omega$$

$$= \frac{e^{\beta \cdot \phi(c^{+}) \cdot \omega}}{e^{\beta \cdot \phi(q_{0}) \cdot \omega} + e^{\beta \cdot \phi(q_{1}) \cdot \omega}} \qquad \text{(since } q = (q_{0}, q_{1}) \text{ in preferences)} \qquad (55)$$
Since $c_{0} \implies c^{+} = q_{0} \text{ and } c_{1} \implies c^{+} = q_{1}$:
$$= \left[\frac{e^{\beta \cdot \phi(q_{0}) \cdot \omega}}{e^{\beta \cdot \phi(q_{0}) \cdot \omega} + e^{\beta \cdot \phi(q_{1}) \cdot \omega}}, \frac{e^{\beta \cdot \phi(q_{1}) \cdot \omega}}{e^{\beta \cdot \phi(q_{0}) \cdot \omega} + e^{\beta \cdot \phi(q_{1}) \cdot \omega}} \right]_{c} \qquad (57)$$

$$= \left[\frac{\mathbf{E}_{q_0, q_1, \omega}}{\left[\mathbf{E} + \mathbf{E}^T \right]_{q_0, q_1, \omega}}, \frac{\mathbf{E}_{q_0, q_1, \omega}^T}{\left[\mathbf{E} + \mathbf{E}^T \right]_{q_0, q_1, \omega}} \right]_{c}$$
(58)

$$= \left[\mathbf{E} \oslash (\mathbf{E} + \mathbf{E}^T), \mathbf{E}^T \oslash (\mathbf{E} + \mathbf{E}^T) \right]_{c,q_0,q_1,\omega}$$
(59)

In the corrections case:

$$\mathbf{P}_{q,c,\omega}^{(\mathrm{corr})} = \frac{e^{\beta \cdot \phi(c^+) \cdot \omega}}{e^{\beta \cdot \phi(c^+) \cdot \omega} + e^{\beta \cdot \phi(q_0) \cdot \omega}} \qquad \text{(since } c \not\in q \text{ and } |q| = 1 \text{ in corrections)}$$

$$= \frac{\mathbf{E}_{q,c,\omega}^T}{[\mathbf{E} + \mathbf{E}^T]_{q,c,\omega}} \quad \text{(since there is a 1-1 correlation between } c \text{ and } c^+ \text{ in corrections)}$$
 (61)

$$= \left[\mathbf{E}^T \oslash (\mathbf{E} + \mathbf{E}^T) \right]_{q,c,\omega} \tag{62}$$

In the binary feedback case, we compare the likelihood of the teacher demonstrating q to the average likelihood of demonstrating any other trajectory in T:

$$\mathbf{P}_{q,c,\omega}^{(\text{bnry})} = \left[\frac{1 - \mathbf{P}_{T,q,\omega}^{(\text{demo})}}{\alpha \left(|T| - 1 \right)}, \frac{\mathbf{P}_{T,q,\omega}^{(\text{demo})}}{\alpha} \right]_{c} \quad \text{(where } \alpha \text{ is a normalization factor s.t. } \sum_{c} \mathbf{P}_{q,c,\omega}^{(\text{bnry})} = 1 \text{)} \quad (63)$$

$$= \left[1 - \frac{\mathbf{P}_{T,q,\omega}^{(\text{demo})}}{\alpha}, \frac{\mathbf{P}_{T,q,\omega}^{(\text{demo})}}{\alpha} \right]_{c} \quad (64)$$

$$= \left[1 - \left(\mathbf{E}_{0}^{T} \otimes \alpha \sum_{t \in T} \mathbf{E}_{t} \right), \mathbf{E}_{0}^{T} \otimes \alpha \sum_{t \in T} \mathbf{E}_{t} \right]_{c,q,\omega} \quad (65)$$

where $lpha=rac{1-\mathbf{P}_{T,q,\omega}^{(\mathrm{demo})}}{|T|-1}+\mathbf{P}_{T,q,\omega}^{(\mathrm{demo})}$

299 A.4 Gradient Derivation

$$\ell(\omega) = \log \prod_{c \in \mathbf{F}} \frac{\sum_{t \in c^+} e^{\beta \cdot \phi(t) \cdot \omega}}{\sum_{t \in c^+ \cup c^-} e^{\beta \cdot \phi(t) \cdot \omega}} = \sum_{c \in \mathbf{F}} \left[\log \frac{\sum_{t \in c^+} e^{\beta \cdot \phi(t) \cdot \omega}}{\sum_{t \in c^+ \cup c^-} e^{\beta \cdot \phi(t) \cdot \omega}} \right]$$
(66)

$$= \sum_{c \in \mathbf{F}} \left[\log \left(\sum_{t \in c^{+}} e^{\beta \cdot \phi(t) \cdot \omega} \right) - \log \left(\sum_{t \in c^{+} \cup c^{-}} e^{\beta \cdot \phi(t) \cdot \omega} \right) \right]$$
 (67)

300 A.5 Training Parameters

We set a high convergence threshold (10^{-3}) when converging each weight sample in order to maintain sparsity within Ω (which becomes less sparse as \mathbf{F} grows with more queries), and then fully converge (threshold 10^{-6}) to report the weight estimate $\tilde{\omega}$ after each query.

Step size of $5x10^{-4}$. (alpha set to 0.75 in our experiments) We enforce $\forall \omega \in \mathcal{W}, ||\omega|| = 1$. Define N, M, K, etc

306 B Evaluation Details

307 B.1 Domain Implementations

Domain #1: Parameter Estimation This task involves directly estimating a randomly-initialized, 308 ground truth weight vector ω^* containing 8 parameters. This represents a generic learning problem 309 relevant to many robotics tasks, such as learning the relative importance between task outcomes 310 according to a user's preference. There is no "state" in this domain, and each "trajectory" consists 311 of a single sample of the weight vector. As a result, we do not enable demonstration queries in 312 this domain, as the resulting feedback would be akin to directly providing ω^* to the algorithm. The 313 feature representation ϕ of a sample returns the sample itself. Since $||\omega|| = ||\omega^*|| = 1$, the reward 314 of any sampled weight vector directly reflects the cosine similarity between it and the ground truth 315 vector $(r(\omega) = \omega \cdot \omega^* = \cos(\theta))$. 316

Domain #2: Linear Dynamical System We consider a simplified Linear Dynamical System representing a robot that optimizes its controls according to a learned task objective. We represent the dynamics of the robot's state \mathbf{x} as $dx/dt = \mathbf{A}\mathbf{x}(t)$ using a randomized dynamics matrix \mathbf{A} . An optimal control vector is one that results in a trajectory of states maximizing $\frac{1}{|T|} \sum_{s \in T} \phi(s) \cdot \omega^*$. We define the feature representation $\phi(s)$ of a state s as the concatenation of the element-wise, absolute difference between $\mathbf{x}(t)$ and the goal state, and the dynamics $\mathbf{u}(t)$. We experiment with an 8-dimensional feature-space (4 state elements and 4 corresponding controls).

In a demonstration query, the oracle provides a trajectory (produced by simulating the execution of a control vector) from the initial state that maximizes the total reward. In a preference query, the algorithm proposes two trajectories and the oracle selects the higher-rewarding option of the two. In a corrections query, the algorithm proposes a trajectory and the oracle returns a trajectory

that maximizes the reward-to-similarity ratio. In a binary feedback query, the algorithm proposes a 328 trajectory and the oracle indicates whether it results in reward that exceeds its internal threshold. 329

Domain #3: Lunar Lander We define a ω^* that results in a Lunar Lander agent efficiently moving 330 from its start state to an upright pose on the landing pad. We use the same feature representation as 331 in [12], consisting of four features: the lander's angle, velocity, distance from the landing pad, and 332 final position with respect to the landing pad. We implement each query type in the same manner as 333 in the Linear Dynamical System. 334

Domain #4: Pizza Arrangement We approximate a preference-learning task in which the robot 335 learns to place toppings on only the left side of a pizza, with uniform spacing between them. We 336 define each "trajectory" as the *next* action the robot should take from the current pizza state; thus, 337 the trajectory is defined as the (x,y)-coordinate of the next topping to be placed. The feature rep-338 resentation consists of four features: the x and y position of the topping, its distance to its nearest-339 neighboring topping, and the difference between that distance and 4cm.

B.2 Oracle Implementation 341

347

353

When responding to a query, the oracle requires its own set of trajectory samples. Similar to IN-342 QUIRE, we derive this set by uniformly sampling N trajectories; however, the two sample sets are 343 kept separate, and so we distinguish the oracle's trajectory set as T' (resampled for each query state). 344 **Demonstration/Preferences** The oracle returns the highest-reward trajectory (according to ω^*) 345 from a uniformly-sampled trajectory set T' (for demonstrations) or from the pair of queried tra-346 jectories C(q) (for preferences):

$$\operatorname{Oracle_{demo}}(s) = \underset{t \in T'}{\operatorname{arg\,max}} \left(\phi(t) \cdot \omega^* \right) \qquad \operatorname{Oracle_{pref}}(s) = \underset{t \in C(q)}{\operatorname{arg\,max}} \left(\phi(t) \cdot \omega^* \right) \tag{68}$$

Corrections The oracle produces T' by performing rejection sampling; it uniformly samples trajec-348 tories and accepts only those with a reward greater than or equal to the queried trajectory q until T'349 contains N trajectories: 350

$$\forall t \in T', \phi(t) \cdot \omega^* > \phi(q) \cdot \omega^* \tag{69}$$

After producing this trajectory set, the oracle selects the trajectory with the highest ratio of reward-351 to-distance from the queried trajectory: 352

$$\operatorname{Oracle_{corr}}(s) = \underset{t \in T'}{\operatorname{arg\,max}} \frac{\Delta_r(t)}{\Delta_d(t)} \tag{70}$$

 $\Delta_r(t) = \min_{t' \in T'} \frac{\phi(t) \cdot \omega^* - \phi(q) \cdot \omega^*}{\phi(t') \cdot \omega^* - \phi(q) \cdot \omega^*} \qquad \qquad \Delta_d(t) = \min_{t' \in T'} \frac{e^{\delta(t,q)}}{e^{\delta(t',q)}}$ (71)

354 The distance metric δ between two trajectories is domain-specific; see the Appendix for details.

Binary Feedback The oracle produces T' by uniformly sampling N trajectories and produces a 355 cumulative distribution R over ground-truth rewards for T'. It then selects a positive or negative 356 reward indicating whether the agent's query q meets or exceeds a threshold percentile α : 357

$$R = \{\omega^* \cdot \phi(t), \forall t \in T'\} \qquad \text{Oracle}_{\text{bnry}}(s) = \begin{cases} + & R(\omega^* \cdot \phi(q)) \ge \alpha \\ - & \text{otherwise} \end{cases}$$
 (72)

Evaluation Procedure 358

AUC Tables

Algorithm 3 Evaluation Procedure

Input: generate_query and update_weights methods according to algorithm being tested **Parameters:** K (# of queries), I (# of tasks), J (# of runs per task), X (# of test states per task),

- 1: **for** each of *I* tasks **do**
- Generate ground truth reward function ω^*
- 3: Generate X test states
- 4: Generate K query states
- 5: Compute optimal trajectory t_{max} for each test case using ω^*
- 6: Compute least-optimal trajectory t_{\min} for each test case using ω^*
- 7: **for** each of J runs **do**
- 8: **for** each of K queries **do**
- 9: $s \leftarrow \text{next query state}$
- $q^* \leftarrow \text{generate_query}(s, \mathcal{I}, \mathbf{\Omega})$ 10:
- $\mathbf{F} \leftarrow \mathbf{F} + \mathsf{query_oracle}(q^*)$ 11:
- 12: $\Omega \leftarrow update_weights(F)$
- $\tilde{\omega} \leftarrow \text{mean}(\Omega)$ 13:
- Record distance: $\frac{\arccos(\tilde{\omega}\cdot\omega^*)}{-}$ 14:
- **for** each of X test states **do** 15:
- Compute optimal trajectory t from the test state according to $\tilde{\omega}$ Record performance: $\frac{\phi(t)\cdot\omega^*-\phi(t_{\min})\cdot\omega^*}{\phi(t_{\max})\cdot\omega^*-\phi(t_{\min})\cdot\omega^*}$ 16:
- 17:

DISTANCE						
Agent	Parameter Estimation	Dynamical System (Static State)	Dynamical System (Changing State)	Lunar Lander (Static State)	Lunar Lander (Changing State)	Mean w* Distance
Binary-only	7.44	4.87	4.91	6.61	6.06	5.98
Corrections-only	3.01	4.43	4.01	4.02	3.80	3.85
Demo-only	n/a	4.26	2.10	3.35	1.91	2.90
Preferences-only	4.30	4.34	4.45	2.31	2.34	3.55
INQUIRE	2.98	2.46	2.59	1.62	1.67	2.26

(a) AUC values for the distance plots, with darker cells indicating lower (better) values.

PERFORMANCE						
Agent	Parameter Estimation	Dynamical System (Static State)	Dynamical System (Changing State)	Lunar Lander (Static State)	Lunar Lander (Changing State)	Mean Task Performance
Binary-only	15.01	16.54	16.90	16.99	17.90	16.67
Corrections-only	19.14	16.53	17.19	18.02	18.33	17.84
Demo-only	n/a	16.76	18.15	18.61	19.32	18.21
Preferences-only	17.74	16.55	16.74	18.63	19.05	17.74
INQUIRE	19.15	17.81	17.83	18.86	19.33	18.60

(b) AUC values for the performance plots, with darker cells indicating higher (better) values.

COST						
Agent	Parameter Estimation	Dynamical System (Static State)	Dynamical System (Changing State)	Lunar Lander (Static State)	Lunar Lander (Changing State)	Mean Query Cost
Binary-only	64.15	43.18	43.60	60.65	53.58	53.03
Corrections-only	36.39	51.68	46.02	41.91	42.57	43.71
Demo-only	n/a	43.99	27.94	37.09	26.07	33.77
Preferences-only	42.41	42.73	43.92	22.90	23.18	35.03
INQUIRE	37.68	36.39	35.92	22.98	23.71	31.34

(c) AUC values for the cost plots, with darker cells indicating lower (better) values.

Figure 4: Area-Under-the-Curve (AUC) values summarizing Figures 2-3