

One-shot Policy Elicitation via Semantic Reward Manipulation

Aaquib Tabrez · Ryan Leonard · Bradley Hayes

Received: date / Accepted: date

Abstract Synchronizing expectations and knowledge about the state of the world is an essential capability for effective collaboration. For robots to effectively collaborate with humans and other autonomous agents, it is critical that they be able to generate intelligible explanations to reconcile differences between their understanding of the world and that of their collaborators. In this work we present Single-shot Policy Explanation for Augmenting Rewards (SPEAR), a novel sequential optimization algorithm that uses semantic explanations derived from combinations of planning predicates to augment agents' reward functions, driving their policies to exhibit more optimal behavior. We provide an experimental validation of our algorithm's policy manipulation capabilities in two practically grounded applications and conclude with a performance analysis of SPEAR on domains of increasingly complex state space and predicate counts. We demonstrate that our method makes substantial improvements over the state-of-the-art in terms of runtime and addressable problem size, enabling an agent to leverage its own expertise to communicate actionable information to improve another's performance.

Keywords Human-agent collaboration · Policy Explanation · Cooperating Robots · Explainable AI

1 Introduction

Autonomous systems have been shown to be capable of motivating behavior changes and conveying new knowledge to improve human or agent performance on a multitude of tasks [5,24,25]. Unfortunately, the process of generating concise and informative explanations

Department of Computer Science, University of Colorado Boulder, Boulder, CO 80309, USA

Aaquib Tabrez (Corresponding Author)
E-mail: mohd.tabrez@colorado.edu

Ryan Leonard
E-mail: ryan.leonard@colorado.edu

Bradley Hayes
E-mail: bradley.hayes@colorado.edu

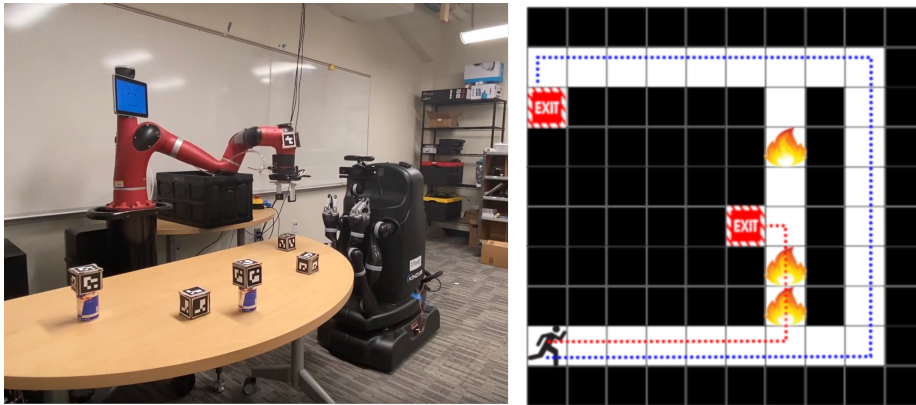


Fig. 1: (Left) The red robot is engaged in a tabletop cleaning task but has a malformed policy that will throw away all of the objects, despite some not being garbage. (Right) A human agent attempts to navigate to a building exit in an emergency evacuation scenario lacking knowledge about the location of the fires. Using SPEAR, these policies can be repaired, leveraging semantic updates to produce more optimal behavior.

capable of eliciting desired changes is a difficult task, as it requires both insights into a collaborator’s decision making process and the ability to determine and convey important information. Further, repairing a human policy requires the generation of explanations which are intelligible and concise, especially in situations with high costs of failure or time pressure [20,8,39]. Similarly, for autonomous agents operating with different state representations, policy repair requires a common ground (language) to communicate updates for efficient behavior modification during the task.

In this work we provide a characterization of the problem of semantically manipulating an agent’s policy through external advice, and propose a solution that provides autonomous agents with the ability to: 1) Infer the reward function that most likely drives another agent’s behavior; 2) Identify divergences between another agent’s reward function and their own; and 3) Communicate advice to repair relevant differences concisely and efficiently at an appropriate level of detail, removing harmful behaviors from the other agent’s policy.

An illustrative scenario that we use to showcase the importance and practicality of this work is routing people during an emergency evacuation, particularly in cases where inhabitants of the building are aware of neither the precise nature of the emergency nor the entire building’s floorplan. Even if people are capable of navigating towards the nearest exit, uncertainty about the nature of the environment and its hazards could be disastrous. Adding to the complexity of time-critical, high-consequence scenarios like this one, it is essential that any advice or instructions account for the recipient’s knowledge of the world, and therefore must also consider tradeoffs between accuracy, specificity, and interpretability [33,23]. As an example, someone visiting a building for a meeting may not know how to change their evacuation plan when told “There’s a fire near Conference Room 3”, but may be able to adapt their plan if told “the north half of the building is on fire” for the same scenario. Even though the latter phrase may not communicate the true representation of the hazard, it is more easily comprehensible to someone who is less familiar with the building than the former. Thus, autonomous systems that aim to generate useful feedback to help others will

need to explicitly consider the complexity of their own explanations and the knowledge held by those they attempt to help.

To achieve this, we present Single-shot Policy Explanation for Augmenting Rewards (SPEAR), a novel Integer Programming-based algorithm for generating reward updates in the form of semantic advice to improve the policies of human and robot collaborators. SPEAR enables an autonomous robot to utilize knowledge about the beliefs and goals of its collaborators to identify inaccuracies in their models, generating targeted, human-interpretable guidance for updating their reward functions (and thus, policies) during a task. Key to its effectiveness, SPEAR generates feedback with levels of specificity appropriate to the agent’s understanding of the world. The four primary contributions of our work are:

- A characterization of the policy elicitation problem.
- SPEAR, a novel algorithm for improving task performance through semantic elicitation of others’ policies.
- A linear programming formulation enabling semantic communication of world state regions, achieving a several order of magnitude improvement over the nearest comparable state-of-the-art [20] in terms of computation time and applicable domain size.
- An experimental validation and performance analysis of SPEAR, grounded in practical application scenarios.

2 Background and Related Work

As autonomous systems become increasingly capable decision makers, Explainable AI (xAI) has emerged as a necessary component for fielding safe autonomous systems. Many of today’s systems with learned control policies continue to only satisfy Donald Michie’s weak criterion for learning systems*, despite explainability being recognized as crucial for improving transparency, trust, and team performance [41,28,43], further underscoring the need for the development of xAI techniques that work alongside popular, generally opaque methods.

Furthermore, these black-box models are difficult to validate and debug when developers do not know what is happening inside the system, leading to further distrust and resistance towards the widespread usage of these algorithms [3,36]. Most of the research in explainable AI has been explicitly focused on making algorithms transparent to developers, allowing them the ability to debug or predict model behavior in restricted domains (e.g., how model parameters affect the final classification decision) [32,35]. These approaches are fundamentally limiting to non-expert stakeholders and/or end-users who interact with the models or products regularly, and thus directly experience the consequences of failures [11, 29]. Therefore, to generate explanations that are comprehensible and useful for both experts as well as non-experts, they should draw from characteristics of everyday human-centric explanations focused on the ‘why’ of particular events, properties, or decisions, rather than general scientific relationships [28,21].

*In 1980, Donald Michie proposed three criteria to evaluate machine learning systems [27]:

1. Weak Criterion: A system increases performance on unseen data by learning from sample data.
2. Strong Criterion: The system contains the weak criterion plus the ability to communicate its learned hypotheses function in symbolic form.
3. Ultra-strong Criterion: The system contains the strong criterion plus the ability to teach a user the learned hypothesis function.

One popular approach in human-robot collaboration has been to use explanations to speed up the learning process or to explore alternative policies. These approaches primarily focus on active learning frameworks for modelling user preferences by asking questions [31] or taking instruction [14] in natural language. While these systems focus on improving a robotic agent’s transparency and behavior using explanation, they don’t account for collaborators’ level of knowledge or need for the information. Therefore, these traditional goal-based agents lack the ability to consider a shared mental model in their planning, a crucial element for building effective teammates and interpretable behavior [7, 40].

Another trend in the direction of “Value of Information” (VOI) in human-robot interaction has been to decide what information should be communicated and when it should be delivered in the collaborative decision making process [22]. Kaupp et al. used VOI theory to propose a framework in which robots query human operators in uncertain environments if the expected benefit of the humans’ feedback exceeds the cost of the query. They show that accounting for human cost in querying has an advantage in performance, operator workload, usability, and the users’ perception of the robot. More recently, active learning has been utilized to understand user preference and provide valuable information while keeping trust and adaptability of humans in account [9, 31, 34].

Identifying and resolving the model differences of a collaborator (i.e., establishing a shared mental model) is a key aspect of establishing interpretability [19, 7, 12]. Having interpretable and predictable behavior is essential for achieving team fluency and avoiding catastrophic failures in planning and decision making [18, 30].

An effective approach for establishing shared mental models in human-robot collaboration has been to use natural language to explain robots’ behavior or underlying logic [20, 37]. In [20], the authors formulated the problem of efficiently describing state regions as a set cover problem, trying to find the smallest logical expression of predicates that succinctly describe a target state region. Their approach is innovative in performing explanation generation using communicable predicates, but their solution is exponential in memory and runtime with respect to the size of the domain and predicate set, preventing its use in most real-world problems.

Others have explored the generation of different types of explanations based on user preference [44]. Work by Briggs and Scheutz explored adverbial cues informed by Grice’s maxims [16] of effective conversational communication (quality, quantity, and relation) to transparently track and update mental models of the collaborators [6]. Van et al. developed transparent Reinforcement Learning (RL) models for the generation of contrastive explanations of agent behavior based on the user’s query-derived policy, and of the learned policy of the agent [42]. They show that people are more interested in explanations about the behavior of an agent rather than a single action taken by it. People have also been shown to prefer contrastive, selective, and social properties in explanation, as previously hinted in the psychology literature [8, 10, 28].

In this work we focus on *policy elicitation*, a process through which feedback is crafted and given to another agent, in the form of reward function updates during task execution, such that they change their behavior to match the desired policy. By providing reward information for targeted regions of state space through explanations (symbolic updates), we can modify a collaborator’s reward function using semantic descriptions. This allows our method to operate at a level of abstraction that does not depend on the recipient’s underlying state space representation, instead only requiring a similar vocabulary of planning predicates.

3 Policy Elicitation via Social Manipulation

The goal of policy elicitation is to cause a behavior change (policy update) in another agent through some form of communicative act. To effectively collaborate with others and coach them towards more optimal policies, it is essential that these communications are intelligible [2, 1, 38], but directly communicating states (i.e. the feature vector itself) relies on both the feature space between agents being the same and there being an efficient mechanism to communicate that information quickly. These criteria are unlikely to hold with either humans or heterogeneous autonomous agents (not all agents will have the same state representations) as the intended recipients. Our work generates state space-agnostic natural language descriptions of state regions and corresponding reward information, allowing agents to update their reward functions (and policies).

We define a **base predicate** to be a pre-defined boolean state classifier (as found in traditional STRIPS planning [13]) with associated string explanation (e.g., `in_central_hallway(x)` \rightarrow “X is in the central hallway.”). To represent intersections of predicates (e.g., “in the central hallway” AND “has fire extinguisher”), we introduce **composite predicates**, which consist of multiple base predicates and evaluate to true if and only if all base predicate members evaluate to true. Predicates may also have a cost associated with them (e.g., how long they take to communicate) to assist the optimizer with generating more desirable solutions. In SPEAR, predicates are used to communicate reward updates to improve collaborator task performance (Figure 3).

We characterize the problem of repairing or otherwise manipulating an agent’s policy as one of identifying and reconciling divergence between a source and target reward function. Our proposed algorithm relies upon the following assumptions: 1) agents can understand the natural language description of the predicates; 2) agents are operating using a rational planner informed by a reward function (or something analogous); and 3) agents’ suboptimal action selections are attributable to a malformed reward function rather than a malformed policy search algorithm or dynamics model.

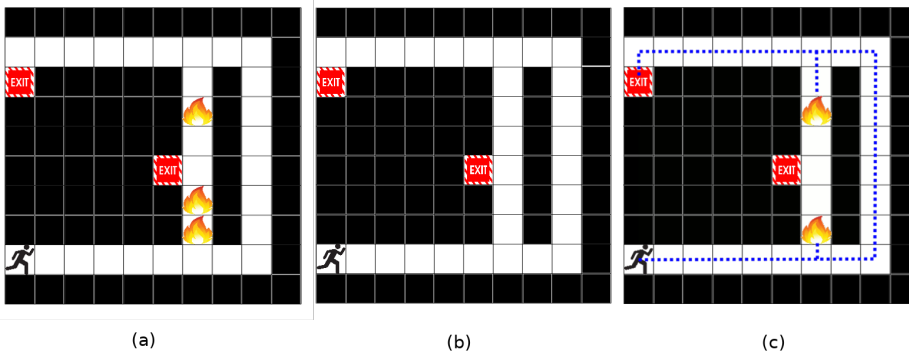


Fig. 2: Rollouts of the agent’s estimated policy are used to identify problematic behavior. (a) The true environment that the agent is acting in. (b) The environment that the agent believes it is acting in. (c) A rollout of the agent’s estimated policy reveals the minimal amount of hazardous states that need to be communicated for policy repair.

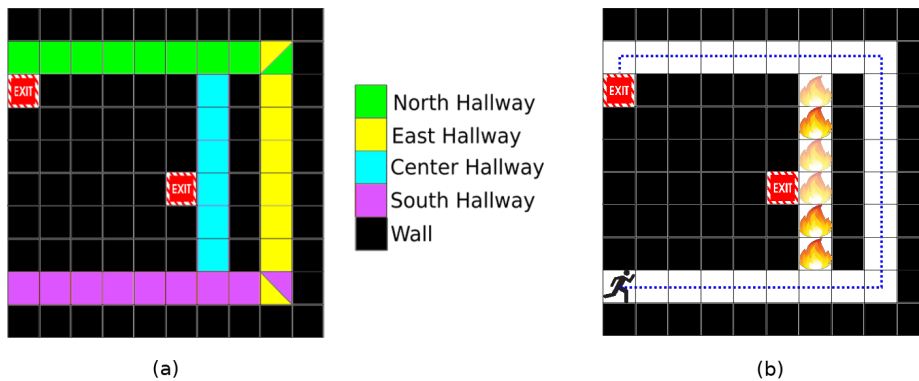


Fig. 3: A Boolean algebra over predicates with string descriptions is used as the basis for grounding reward function updates in language. (a) A map of the environment with an overlay indicating states where various predicates are true. (b) A potential belief update for the human’s reward function, realized after communicating “The center hallway is on fire”. Some of the fires, indicated by being faded, don’t actually exist in the environment and are an artifact of imprecise language. Despite this imprecision, the optimal policy can be elicited from this update.

4 Approach

In this section, we specify the problem of reparative policy elicitation: repairing or manipulating an agent’s policy by communicating corrections for model divergences. Our approach to policy elicitation uses three components that: 1) estimate the agent’s reward function; 2) determine important reward function disparities that prevent desirable behavior; and 3) determine which states to provide reward updates for and communicate a corrective explanation.

4.1 Belief Estimation with Active Observation

We formulate our domain as a Markov Decision Process (MDP)[4], wherein an agent acts to maximize an expected reward. M is a Markov Decision Process defined by the 4-tuple (S, A, T, R) where S is the set of states in the MDP, A is the set available actions, T is a stochastic transition function describing the model’s action-based state transition dynamics, and R is the reward function $R : S \times A \times S \rightarrow \mathbb{R}$. Intuitively, M serves as a simulator for an agent in the task domain.

Belief estimation (inferring an agent’s plan) is performed by leveraging the assumption that the agent has a rational planner and any sub-optimal behavior can be attributed to a misinformed understanding of the task’s reward function R_h , where R_h is the malformed reward function the agent is using [37]. We assume the agent’s dynamics model is correct ($T_h = T$).

Inferring an agent’s reward function R_h (e.g., the agent’s beliefs about the task w.r.t. the world) is non-trivial and an active field of research. Tabrez et al. [37] provides an approach to this problem by solving for whether the agent knows about particular components of the

reward function; therefore if one can maintain a belief over which rewards another agent is aware of, their reward function may be able to be inferred.

Intuitively, estimating R_h is achieved by observing agent behavior and the goal states they approach — in the emergency evacuation domain, this corresponds to observing their coverage of the building and which exits they are heading toward. Once we obtain R_h , we can simulate the MDP M using R_h to get the agent’s expected policy π_h .

4.2 Finding Important Model Divergences

Once we have our belief of the agent’s possible reward function R_h , we identify divergences between the optimal policy π^* and the policy of the agent π_h that would cause a reduction in the agent’s expected cumulative reward. We do this by comparing agent policies trained on R and R_h , where R is the true reward function of the domain and R_h is the reward function used by the agent. We find a set of states \bar{S} about which we communicate updated reward information to augment R_h , such that a more optimal policy (closer to the expected reward of π^*) is elicited (Figure 3b).

4.3 Communicating State Regions

We approach the problem of efficiently describing state regions as a set cover problem, trying to find the smallest logical expression of communicable predicates to succinctly describe target states as in Hayes and Shah [20]. Unlike prior work, we solve for the minimum set cover of the targeted state region using an Integer Program. The inputs to our IP formulation include:

- A set of state indices \bar{S} that correspond to states with expected reward function divergence that needs to be communicated about for policy repair.
 $\bar{S} = \{s_1, s_2, \dots, s_{|\bar{S}|}\},$
- A set of communicable predicates $\bar{P} = \{p_1, p_2, \dots, p_{|\bar{P}|}\}$ requiring the following for a solution to exist:
 - $\exists \bar{Q} \subseteq \bar{P}$ such that $s \in \bar{S}$ is covered by a predicate in $\bar{Q}, \forall s \in \bar{S}$ — for every state that needs to be covered (\bar{S}), there exists a non-empty subset of predicates that can cover it.
- A set of costs $\bar{C} = \{c_1, c_2, \dots, c_{|\bar{P}|}\}$, such that $c \neq 0 \forall c \in \bar{C}$ — every predicate has non-zero cost associated with using it to update the agent’s reward function R_h .
- The desired trajectory for the agent $\bar{O} = \{o_1, o_2, \dots, o_{|\bar{O}|}\}$, where o_i describes the state achieved after taking the i^{th} action following the optimal policy from the start state.

The cost for each predicate can be customized per task. This is an important characteristic, as many factors may influence the cost of a predicate. One such criteria for defining cost can be the length of the string describing the predicate. Such a criteria could generate more easily understood explanations by imposing penalties for being too verbose.

A solution to the policy elicitation problem consists of selecting predicates to communicate reward information about specific state regions such that a more optimal policy is produced within some ϵ bound of the optimal policy’s expected reward, $|E_R(\pi|R) - E_R(\pi|R_h)| \leq \epsilon$. To minimize this objective while satisfying all the constraints, we define the mathematical formulation of our IP, which we refer to as **SPEAR-IP**, below:

$$\min \sum_{j=1}^{|\bar{P}|} c_j x_j + L \sum_{k=1}^{|\bar{O}|} \sum_{j=1}^{|\bar{P}|} v_{kj} x_j \quad (1)$$

subject to

$$\sum_{j=1}^{|\bar{P}|} u_{ij} x_j \geq 1 \quad \forall i \in [1, |\bar{S}|] \quad (2)$$

where we define the $|\bar{S}| \times |\bar{P}|$ matrix U by

$$u_{ij} = \begin{cases} 1, & \text{if } s_i \text{ is covered by } p_j \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

and, the $|\bar{O}| \times |\bar{P}|$ matrix V by

$$v_{kj} = \begin{cases} 1, & \text{if } o_k \text{ is covered by } p_j \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

where i and k are the indices into sets \bar{P} and \bar{O} respectively. $x_j \in \{0, 1\}$ is a binary variable with $x_j = 1$ meaning predicate p_j is included in the cover and $x_j = 0$ indicating exclusion from the set cover. Equation 1 is the objective function which needs to be minimized for the desired set cover based on cost. The first term of Equation 1 ensures that the net cost of the predicates selected is as low as possible.

Equation 1's second term heavily penalizes the objective function if the chosen set cover overlaps with the desired path, in the case of communicating a negative reward (for positive reinforcement this can be reversed). L is a large positive number, and is used to encapsulate near-optimal behavior in the objective function (providing a soft constraint). The second term can be used separately to provide a hard constraint for finding a solution which will definitely elicit the desired path, but this approach inherently restricts the ability of an integer program to find solutions which would provide a *near* optimal policy update. This second term in the formulation leads to **three possible cases**: 1) Set cover solution with low objective value; 2) No solution for the set cover; and 3) Set cover solution with high objective value.

Cases 1 and 2 are simple and describe the ability of SPEAR-IP to solve the set cover. Case 3 is interesting and provides the option of further exploration to find an alternate solution with the given set of predicates. We can use the set cover from Case 3 to determine which states in the cover overlapped with the desired path. Using these states as a reference, we can simulate an alternate desired policy by penalizing these states in the true reward function as well, effectively coming up with a contingency for not having precise enough language available. This process of penalizing states through R is repeated until either Case 1 or 2 is achieved (i.e. an immediate solution with low objective value or no solution for given a set of predicates).

Equation 2 provides a hard constraint for the inclusion of states from \bar{S} in the set cover. Equation 3 defines the elements of matrix U , which encapsulates cover constraints from \bar{S} (inclusion of all states). Equation 4 defines the elements of matrix V using the desired trajectory \bar{O} , encapsulating the requirements for eliciting the desired policy.

Algorithm 1: Single-shot Policy Explanation for Augmenting Rewards (SPEAR)

```

Input: MDP  $(S, A, T, R)$ , Min. Reward Threshold  $R_L$ , Agent Reward Function  $R_h$ , Current state  $s_c$ , Num. Rollouts  $k$ 
Output: Semantic Reward Correction
1  $\bar{S} \leftarrow \emptyset$ ;  $L \leftarrow$  large scalar value;
2  $R_h^* \leftarrow R_h$ ; // Best possible agent reward function
3 for rollout in range(1 to  $k$ ) do
4    $r_c \leftarrow 0$ ; // Cumulative reward
5    $\pi_h \leftarrow$  policy trained on  $R_h$ ;
6   // Find states to update for best possible  $R_h^*$ 
7    $s \leftarrow s_c$ ;
8   while  $s$  is not a terminal state do
9     // Perform forward rollout of  $\pi_h$ 
10     $s' \leftarrow T(s, \pi_h(s))$ ;
11     $r_c \leftarrow r_c + R(s, \pi_h(s), s')$ ;
12     $s \leftarrow s'$ ;
13    if  $r_c \leq R_L$  then // Reward too low
14       $\bar{S} = \bar{S} \cup s'$ ; // Track state for later
15       $R_h^*(s, \pi_h(s), s') \leftarrow R(s, \pi_h(s), s')$ ;
16      break;
17  $\pi_h^* \leftarrow$  policy trained on  $R_h^*$ ;  $\pi^* \leftarrow$  policy trained on  $R$ ;
18  $Set\_Cover, Objective \leftarrow$  predicate_selection( $\bar{S}, \pi_h^*, \dots$ );
19 if objective is no_solution then exit;
20 if objective  $\geq L$  (from Eq. 1) then
21   for rollout in range(1 to  $k$ ) do
22      $s \leftarrow s_c$ ;
23     while  $s$  is not terminal do
24       // Perform forward rollout of  $\pi^*$ 
25        $s' \leftarrow T(s, \pi^*(s))$ ;
26       if  $s' \in Set\_Cover$  then  $R(s, \pi^*(s), s') \leftarrow -L$ ;
27        $s \leftarrow s'$ ;
28   go to 1;
29 else
30   feedback  $\leftarrow$  "There's a bad reward in  $Set\_Cover$ ."
31   return feedback

```

4.4 Algorithm

Here we outline the details of SPEAR for finding the predicates to communicate which will elicit optimal behavior, as presented in Algorithm 1. We continue using our running example of an emergency evacuation scenario (Figure 1-right) to build intuition and illustrate the inner workings of our algorithm for a human-robot collaborative scenario.

Given an estimate of the human agent’s reward function, SPEAR communicates a reward update in an attempt to elicit the best possible policy. To achieve this, (Line 3-16) we perform multiple forward rollouts of the agent policy π_h derived from R_h and (Line 13-16) compare the accumulated expected reward to the reward threshold R_L . The moment this threshold is crossed, the reward from that transition is determined to be relevant for updating R_h and the state is added to the set cover.

This can be easily illustrated for our example, where Figure 2b shows the belief of an agent trying to evacuate the building. Figure 3b gives insight into how updating the agent’s reward function can result in an optimal policy even if the agent’s belief about the fires doesn’t truly match the environment.

After finding the states responsible for meaningful reward divergence, we find a minimal set cover of communicable predicates that map onto these states. This is achieved through **SPEAR-IP**, discussed in Section 4.3. We use a third-party optimizer [17] to solve SPEAR-IP (Line 18), where the *predicate_selection* method takes in the states to cover (\bar{S}) and the best agent policy π_h^* . This gives a set of communicable predicates and a final objective value using Algorithm 2.

Line 19 checks whether or not a solution exists for a given set of predicates. In lines 20-28, our algorithm evaluates case 3 to determine if alternate solutions exist. In lines 23-27, SPEAR performs multiple forward rollouts of the optimal policy to find states responsible for a high objective value. In line 26, these states are penalized in the true reward function (R) to incentivize the algorithm to find an alternate solution which avoids these states. This enables SPEAR to explore alternate solutions, making it more robust in applications where the available predicates are insufficient, overcoming barriers due to imprecise or unavailable language.

In line 28, now that all the appropriate states are penalized, the algorithm repeats the whole procedure from the beginning with a modified R , continuing this process until it finds a low objective solution or no solution (case 1 or 2). Finally, in line 30, the update is serialized as semantic feedback using the *Set_Cover*. This feedback generation strategy uses negative reward to drive an agent’s policy away from undesirable states, as improved policies can then be elicited through the exclusion of states along the agent’s (hypothesized) originally intended path. While similar outcomes can be achieved via positive reward, a state exclusion-based strategy generally allows for the use of less precise predicates.

In Algorithm 2, we produce the set cover for communicating the reward update. Lines 4-7 evaluate states we want the agent to traverse (the desired trajectory \bar{O}) by performing a forward rollout of the best attainable agent policy π_h^* . In lines 8-13, the matrices U and V from Equation 3-4 are defined, which form the basis of the constraints governing the inclusion of states to cover and exclusion of optimal states (when giving information about negative reward) in the set cover respectively. Finally, in line 15, *SPEAR-IP* is solved using the matrices U and V to give *Set_Cover* and *Objective*.

5 Evaluation and Discussion

To demonstrate the utility of our algorithm we present two applications in which an agent is required to augment the policies of other agents (autonomous or human) via policy elicitation from SPEAR. We also provide a characterization of SPEAR’s performance as a function of domain size and predicate count. Results are generated on an Intel(R) Core(TM) i7-7700K CPU @ 4.20GHz.

5.1 Robotic Cleaning Task

In the first scenario, one agent needs to correct the policy of a second robotic agent to prevent it from removing specific items during a pick and place task (Figure 4). Importantly, these robots do not operate in the same state space, and therefore cannot directly communicate reward function updates to each other. Here, the Rethink Robotics Sawyer has been tasked with clearing *trash* from a table, and is operating with the malformed belief that all objects on the table’s surface are *trash*. The Kinova Movo is observing this scene and has accurate knowledge about the environment (i.e., it knows that certain objects in the scene aren’t

Algorithm 2: Predicate selection (Minimal set cover)

Input: Set of States to cover \bar{S} , Agent policy π , MDP (S, A, T) , Set of predicates \bar{P} , Current state s_c

Output: Set_Cover and Objective (high, low, or no solution)

- 1 Set_Cover $\leftarrow \emptyset$; // Predicates in min set cover
- 2 $O \leftarrow \emptyset$; // States we don't want to cover
- 3 $s \leftarrow s_c$; $O \leftarrow O \cup s$;
- 4 **while** s is not terminal **do**
- 5 // Perform 'optimistic' forward rollout of π
- 6 $s \leftarrow$ most likely transition from $T(s, \pi(s))$;
- 7 $O \leftarrow O \cup s$; //append occupied states
- 8 //Define matrix U to be $|\bar{S}| \times |\bar{P}|$ matrix s.t.
- 9 **for** $i \in [1, |\bar{S}|]$, $j \in [1, |\bar{P}|]$
- 10 $u_{ij} = \begin{cases} 1, & \text{if } s_i \in \bar{S} \text{ is covered by } p_j \in \bar{P} \\ 0, & \text{otherwise} \end{cases}$
- 11 //Define matrix V to be $|\bar{O}| \times |\bar{P}|$ matrix s.t.
- 12 **for** $k \in [1, |\bar{O}|]$, $j \in [1, |\bar{P}|]$
- 13 $v_{kj} = \begin{cases} 1, & \text{if } o_k \in \bar{O} \text{ is covered by } p_j \in \bar{P} \\ 0, & \text{otherwise} \end{cases}$
- 14 //For SPEAR-IP: refer to Equations 1 - 2
- 15 Set_Cover, Objective \leftarrow SPEAR-IP(U, V);
- 16 **return** Set_Cover; Objective

trash). As Sawyer's reach expresses intent to remove one of the items that isn't trash, Movo detects that Sawyer's policy is incorrect. Movo corrects Sawyer's behavior by providing an update for its reward function, allowing Sawyer to compute a better policy. We accomplish this reward update step by generating semantic expressions about the reward function in parts of the state space (i.e., predicate-grounded natural language communicating about a region of negative reward) for Sawyer using SPEAR.

Within this task, Sawyer used a series of ArUco markers [15] to track the 6-dof poses of various objects (states) on the tabletop. Using these poses, Sawyer systematically transferred all items classified as *trash* to the waste bin using an interruptable pick and place action. As Movo observes Sawyer, it detects Sawyer's end effector reaching for one of the energy drinks on the table, thus indicating the intent to remove that object, and signalling to Movo that Sawyer's policy was malformed. Using SPEAR, Movo is able infer Sawyer's erroneous belief from this observed action (that a non-negative reward is associated with putting the energy drink in the trash). Movo performs a forward rollout with the inferred policy of Sawyer to predict which state transitions with *negative reward* Sawyer will reach. Next, Movo computes an updated policy for Sawyer by determining which states should be assigned negative reward and which of the available predicates are needed for communicating it with Algorithm 2. A DNF formula of predicates that covers the desired states is then computed and communicated by Movo through natural language: "*There is a Bad Reward when energy drink is in the trash*". Finally, Sawyer uses its own set of predicates to map the communication into its own state space, update its reward function accordingly, and reconverge a repaired policy.

5.2 Emergency Evacuation

In our second scenario, our autonomous agent must help a human agent escape from a smart building in which a series of fires has broken out, as shown in Figure 2. While people in the

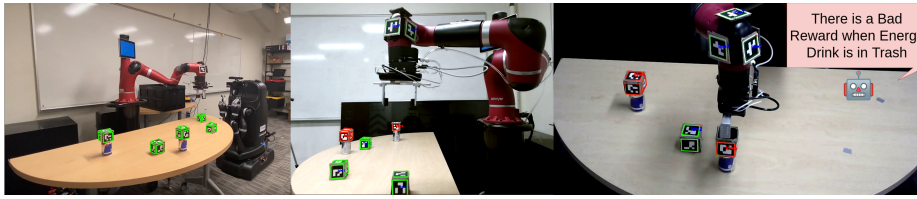


Fig. 4: (Left) Sawyer (red robot) views a misinformed policy for a tabletop cleaning task as Movo (black robot) moves into position to assist. (Center) Movo, from its perspective views the scene with the correct policy and watches Sawyer remove the correct (green) objects from the table. Movo then observes Sawyer reach for one of the incorrect (red) objects and delivers a verbal reward update. (Right) Sawyer computes an update to its planner thus aligning with the desired policy.

building are not aware of the fire locations, they are assumed to understand the generated predicate-based language.

The building layout for each trial is generated with a stochastic placement of rooms, hallways, and exits over a gridworld of fixed size. For example, a gridworld of size 40×40 (1600 states) may have parameters: rooms = 10, hallways = 40, number of exits = 5. These can be changed depending on the desired architecture. We then use these areas to create predicates which are grounded in the layout of the building. In our validation, we also use randomly generated predicates to prove the generalizability of our approach.

To accommodate the full range of possible state regions to cover using predicates, we create composite predicates by creating the power set of base predicates. By combining base predicates into composite predicates, each scenario has an upper bound of $2^n - 1$ total predicates.

Once the building layout and predicates are generated, we observe a randomly placed human agent explore a hazard-free building over a fixed number of episodes. The human agent’s policy is then trained to always seek the shortest path to the closest exit that was discovered during these exploration episodes. Initially, SPEAR has no knowledge about which exits the human is aware of, but gradually, its belief about the human’s reward function is updated from these observations. Next, we begin the evaluation by adding fire in the building via random placement and expansion. Once this process completes, we start the SPEAR evaluation. Predicates from SPEAR-IP are used to update the human agent’s reward function during the episode. We update the policy of the human using the repaired reward function after each update.

5.3 Analysis

We evaluated the performance of our algorithm as a function of state space size and the number of predicates using different building layouts. We also demonstrated that our algorithm can generalize to arbitrary state mappings by using randomly generated predicates. To fairly evaluate this, we randomly generate ball-shaped predicates (evaluating to true within a given radius of a random state vector). When in a structured environment, the predicates grounded in the layout of the building tend to naturally make it easier to find potentially hazardous predicates. For example, in Figure 3a, predicates are grounded in the layout of the building and fire can be avoided through communicating the “*Center Hallway*” predicate as it is directly associated with pathway of the human. However, if predicates are placed stochastically throughout the map (corresponding to cases where the language isn’t tailored

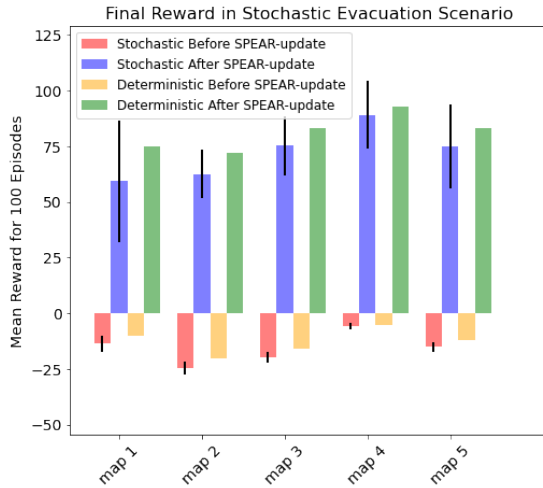


Fig. 5: The evaluation of SPEAR algorithm in stochastic and deterministic evacuation domains (25x25). Substantial increases in episode reward occur due to the generated symbolic reward updates.

to the domain), finding a solution becomes more difficult. Hence, by using randomly placed predicates we make it difficult for the algorithm to find a low cost solution while also directly increasing the computational time, providing performance analysis for cases of predicate-domain mismatch.

We qualitatively assess our algorithm in an emergency evacuation scenario for both deterministic and stochastic environments (Figure 5). In the stochastic domain, a stochastic transition function was applied describing the model’s action-based state transition dynamics. For each map (25 X 25) in both environments, we evaluated our algorithm over the course of 100 episodes. In each episode the reward was computed both before and after the SPEAR update (exit: +100, fire: -100, and each step: -1). At each time step, our agent would take the prescribed action with a probability of 85% and would alternatively take a step in a different random direction. For our stochastic evaluation, we set our forward rollout count parameter, k , to 10. The results from our simulations show a substantial improvement in the episodic reward after the SPEAR update, (Figure 5) demonstrating its utility in both deterministic and stochastic domains.

Furthermore, we do performance analysis as a function of predicate count by dividing into two success cases described in Section 4.3: 1) Maps with a low objective value solution (Case 1); and 2) Maps with high objective value solution (Case 3).

For each map of a given size, we start with 100 base predicates and increase in increments of 10. Figure 6-right shows how algorithm performance changes with increasing predicate count on low objective value maps. For Case 1 solutions our algorithm exhibits linear computational costs as a function of the number of predicates, and the set cover can be computed within 50 seconds with nearly 10,000 communicable predicates to choose from. To provide context for the significance of this result, we can consider a comparison against the state-of-the-art method for succinctly describing state regions (using the Quine-McCluskey algorithm) by Hayes and Shah [20]. Our Integer Programming approach to this set cover problem achieves a several order of magnitude improvement over the QM based method (where set covers with 10 possible predicates takes 60-120 seconds to compute on

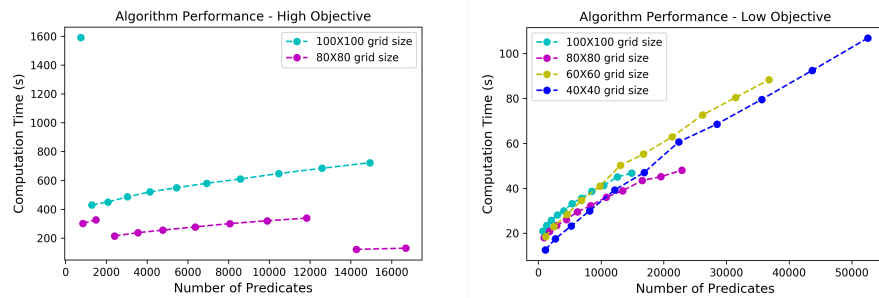


Fig. 6: (Left) SPEAR’s performance for high-objective (case 3) maps scales linearly as predicate count increases. Sharp decreases in computation time (purple) occur when new predicates cause SPEAR to find solutions in fewer algorithm loops. (Right) SPEAR’s performance for low objective (case-1) maps. Plots reveal a linear relationship between runtime and predicate count.

similar hardware, deteriorating exponentially as a function of predicate count). SPEAR is made possible by this performance improvement, effectively operationalizing the insight of Hayes and Shah’s work as a means of communicating information about those state regions — enabling SPEAR’s policy update method.

Similarly, we plot the performance for Case 3 solutions as predicate count increases (Figure 6-left). We observe that the plots again scale linearly in practice with respect to the number of predicates, but with higher computation time due to multiple SPEAR runs. Computation time is higher here because the algorithm has to explore alternative solutions for the desired policy, solving for set cover solutions multiple times (Section 4.3).

Our method achieves a dramatic improvement in computation time over the QM based set cover approach for two reasons: offline computation of predicates and allowing approximate solutions. QM uses a two-step procedure: 1) finding prime implicants, and 2) using those prime implicants to find the exact minimal set cover (see [26]). In SPEAR, we alleviate a significant time sink by pre-computing predicate values over the state space (which the QM approach had to solve online in step 1). Secondly, since SPEAR is often able to use approximate set covers to communicate its reward updates (i.e., sometimes it’s okay to over-state or under-state the reward region, as in Figure 3b), we are able to relax the requirement of having an exact set cover. In contrast, QM solves for exact set cover making it slower in execution and limited in application. For example, referring to Figures 2c and 3a, QM will try to find the precise set cover for covering only the two fire states in the central hallway (and fail as there is no exact solution). Whereas, SPEAR will try to find the approximate set cover that encompasses those fire states, accepting inaccuracy if necessary so long as it still achieves the same outcome as if there were a perfect solution (Case 1) and preferring a suboptimal solution to no solution (Case 3). This added layer of approximation heuristic makes SPEAR faster and more flexible than the QM based approach.

Interestingly, we found some irregularities in performance as illustrated within the 80 X 80 world (Figure 6-left-purple). We found that a sudden dip in computation time can occur when new language enables the algorithm to find a Case 1 (single-loop of algorithm) or cheaper Case 3 (still multiple loops of the algorithm, but fewer) solution.

The final part of our analysis looks at the performance of our algorithm as domain size increases with fixed predicate count (100 base predicates). We generate a set of maps with similar parameters for a fixed state space size, sampling from these maps to get the mean and confidence bound for 10 simulation runs as shown in Figure 7. A significant take-away from

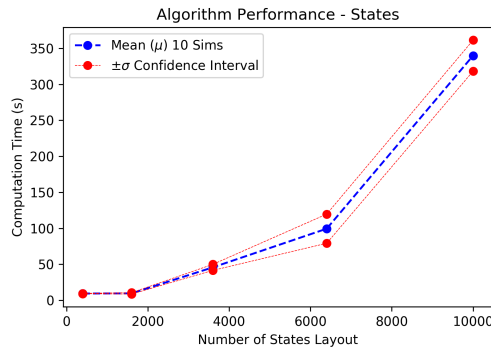


Fig. 7: The evaluation of the SPEAR algorithm in stochastic and deterministic evacuation domains (25x25). SPEAR’s performance as state space size increases using 100 randomly generated base predicates. Computation time increases polynomially as a function of state space size

this analysis is the insight that an attention mechanism is more important for abstracting and reducing the domain’s state space than for limiting the number of predicates to consider, as prior work anticipated [1, 20, 37].

We have shown that *SPEAR* enables dramatic improvements in agent performance through policy elicitation, while also contributing a novel method for communicating about state regions that substantially outperforms prior work. These contributions enable semantically guided policy manipulation for a much broader class of problems than was previously possible, providing a method that scales linearly with predicate count as opposed to exponentially, advancing the state-of-the-art in autonomous coaching through new algorithms and improved foundational capability.

6 Conclusion

In this work, we describe an approach for *policy elicitation*, the manipulation of an agent’s behavior through the use of semantically grounded reward updates. We present a novel Integer Programming-based algorithm for rendering policy explanation [20] and policy manipulation [37] techniques feasible for use in applications substantially larger than previously possible. We demonstrate a several order of magnitude improvement over the state-of-the-art in terms of runtime and allowable domain size when communicating about state regions, and utilize this to achieve significant improvements in agent performance through policy elicitation from an autonomous coaching agent.

References

1. Abel, D., Arumugam, D., Lehnert, L., Littman, M.: State abstractions for lifelong reinforcement learning. In: International Conference on Machine Learning, pp. 10–19 (2018)
2. Andre, D., Russell, S.J.: State abstraction for programmable reinforcement learning agents. In: AAAI/IAAI, pp. 119–125 (2002)
3. Arnold, M., Bellamy, R.K., Hind, M., Houde, S., Mehta, S., Mojsilović, A., Nair, R., Ramamurthy, K.N., Olteanu, A., Piorkowski, D., et al.: Factsheets: Increasing trust in ai services through supplier’s declarations of conformity. *IBM Journal of Research and Development* **63**(4/5), 6–1 (2019)
4. Bellman, R.: A markovian decision process. *Journal of Mathematics and Mechanics* pp. 679–684 (1957)

5. Belpaeme, T., Kennedy, J., Ramachandran, A., Scassellati, B., Tanaka, F.: Social robots for education: A review. *Science robotics* **3**(21), eaat5954 (2018)
6. Briggs, G., Scheutz, M.: Facilitating mental modeling in collaborative human-robot interaction through adverbial cues. In: *Proceedings of the SIGDIAL 2011 Conference*, pp. 239–247 (2011)
7. Chakraborti, T., Kambhampati, S., Scheutz, M., Zhang, Y.: Ai challenges in human-robot cognitive teaming. *arXiv preprint arXiv:1707.04775* (2017)
8. Chakraborti, T., Sreedharan, S., Grover, S., Kambhampati, S.: Plan explanations as model reconciliation. In: *2019 14th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, pp. 258–266. IEEE (2019)
9. Chitnis, R., Kaelbling, L.P., Lozano-Pérez, T.: Learning what information to give in partially observed domains. In: *Conference on Robot Learning*, pp. 724–733 (2018)
10. Ciocirlan, S.D., Agrigoroaie, R., Tapus, A.: Human-robot team: Effects of communication in analyzing trust. In: *2019 28th IEEE International Conference on Robot and Human Interactive Communication (RO-MAN)*, pp. 1–7. IEEE (2019)
11. Doshi-Velez, F., Kim, B.: Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608* (2017)
12. Dragan, A.D.: Robot planning with mathematical models of human state and action. *arXiv preprint arXiv:1705.04226* (2017)
13. Fikes, R.E., Nilsson, N.J.: Strips: A new approach to the application of theorem proving to problem solving. *Artificial intelligence* **2**(3-4), 189–208 (1971)
14. Fukuchi, Y., Osawa, M., Yamakawa, H., Imai, M.: Autonomous self-explanation of behavior for interactive reinforcement learning agents. In: *Proceedings of the 5th International Conference on Human Agent Interaction*, pp. 97–101. ACM (2017)
15. Garrido-Jurado, S., Muñoz-Salinas, R., Madrid-Cuevas, F.J., Marín-Jiménez, M.J.: Automatic generation and detection of highly reliable fiducial markers under occlusion. *Pattern Recognition* **47**(6), 2280–2292 (2014)
16. Grice, H.P.: *Logic and conversation*. In: *Speech acts*, pp. 41–58. Brill (1975)
17. "Gurobi Optimization LLC": Gurobi optimizer (2019). URL <http://www.gurobi.com>
18. Hayes, B., Scassellati, B.: Challenges in shared-environment human-robot collaboration. In: "Collaborative Manipulation" Workshop at the 8th ACM/IEEE International Conference on Human-Robot Interaction., p. 8 (2013)
19. Hayes, B., Scassellati, B.: Effective robot teammate behaviors for supporting sequential manipulation tasks. In: *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* (2015)
20. Hayes, B., Shah, J.: Improving robot controller transparency through autonomous policy explanation. In: *Proceedings of the 12th ACM/IEEE International Conference on Human-Robot Interaction* (2017)
21. Henne, P., Niemi, L., Pinillos, Á., De Brigard, F., Knobe, J.: A counterfactual explanation for the action effect in causal judgment. *Cognition* **190**, 157–164 (2019)
22. Kaupp, T., Makarenko, A., Durrant-Whyte, H.: Human-robot communication for collaborative decision making—a probabilistic approach. *Robotics and Autonomous Systems* **58**(5), 444–456 (2010)
23. Lage, I., Chen, E., He, J., Narayanan, M., Kim, B., Gershman, S., Doshi-Velez, F.: An evaluation of the human-interpretability of explanation. *arXiv preprint arXiv:1902.00006* (2019)
24. Leite, I., Martinho, C., Paiva, A.: Social robots for long-term interaction: a survey. *International Journal of Social Robotics* **5**(2), 291–308 (2013)
25. Leyzberg, D., Spaulding, S., Scassellati, B.: Personalizing robot tutors to individuals' learning differences. In: *Proceedings of the 2014 ACM/IEEE international conference on Human-robot interaction*, pp. 423–430. ACM (2014)
26. McCluskey, E.J.: Minimization of boolean functions. *The Bell System Technical Journal* **35**(6), 1417–1444 (1956)
27. Michie, D.: Machine learning in the next five years. In: *Proceedings of the 3rd European Conference on European Working Session on Learning*, pp. 107–122. Pitman Publishing, Inc. (1988)
28. Miller, T.: Explanation in artificial intelligence: Insights from the social sciences. *Artificial Intelligence* **267**, 1–38 (2019)
29. Mittelstadt, B., Russell, C., Wachter, S.: Explaining explanations in ai. In: *Proceedings of the conference on fairness, accountability, and transparency*, pp. 279–288 (2019)
30. Nikolaidis, S., Shah, J.: Human-robot cross-training: computational formulation, modeling and evaluation of a human team training strategy. In: *Proceedings of the 8th ACM/IEEE international conference on Human-robot interaction*, pp. 33–40. IEEE Press (2013)
31. Racca, M., Kyrki, V.: Active robot learning for temporal task models. In: *Proceedings of the 2018 ACM/IEEE International Conference on Human-Robot Interaction*, pp. 123–131. ACM (2018)
32. Ribeiro, M.T., Singh, S., Guestrin, C.: "why should i trust you?" explaining the predictions of any classifier. In: *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 1135–1144 (2016)

33. Robinette, P., Li, W., Allen, R., Howard, A.M., Wagner, A.R.: Overtrust of robots in emergency evacuation scenarios. In: The Eleventh ACM/IEEE International Conference on Human Robot Interaction, pp. 101–108. IEEE Press (2016)
34. Sadigh, D., Dragan, A., Sastry, S., Seshia, S.A.: Active preference-based learning of reward functions. In: Robotics: Science and Systems (RSS) (2017)
35. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Grad-cam: Visual explanations from deep networks via gradient-based localization. In: Proceedings of the IEEE international conference on computer vision, pp. 618–626 (2017)
36. Sokol, K., Flach, P.: Explainability fact sheets: a framework for systematic assessment of explainable approaches. In: Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency, pp. 56–67 (2020)
37. Tabrez, A., Agrawal, S., Hayes, B.: Explanation-based reward coaching to improve human performance via reinforcement learning. In: 2019 14th ACM/IEEE International Conference on Human-Robot Interaction (HRI), pp. 249–257. IEEE (2019)
38. Tabrez, A., Hayes, B.: Improving human-robot interaction through explainable reinforcement learning. In: HRI Pioneers workshop at the ACM/IEEE International Conference on Human-Robot Interaction (HRI) (2019)
39. Tabrez, A., Luebbers, M.B., Hayes, B.: Automated failure-mode clustering and labeling for informed car-to-driver handover in autonomous vehicles. arXiv preprint arXiv:2005.04439 (2020)
40. Tabrez, A., Luebbers, M.B., Hayes, B.: A survey of mental modeling techniques in human-robot teaming. *Current Robotics Reports* pp. 1–9 (2020)
41. Viganò, L., Magazzeni, D.: Explainable security. arXiv preprint arXiv:1807.04178 (2018)
42. van der Waa, J., van Diggelen, J., Bosch, K.v.d., Neerincx, M.: Contrastive explanations for reinforcement learning in terms of expected consequences. arXiv preprint arXiv:1807.08706 (2018)
43. Wang, N., Pynadath, D.V., Hill, S.G.: Trust calibration within a human-robot team: Comparing automatically generated explanations. In: The Eleventh ACM/IEEE International Conference on Human Robot Interaction, pp. 109–116. IEEE Press (2016)
44. Zahedi, Z., Olmo, A., Chakraborti, T., Sreedharan, S., Kambhampati, S.: Towards understanding user preferences for explanation types in model reconciliation. In: 2019 14th ACM/IEEE International Conference on Human-Robot Interaction (HRI), pp. 648–649. IEEE (2019)