# Action Categorization for Computationally Improved Task Learning and Planning

### Lakshmi Nair and Sonia Chernova

Georgia Institute of Technology lnair3@gatech.edu, chernova@cc.gatech.edu

#### **Abstract**

This paper explores the problem of task learning and planning, contributing the *Action-Category Representation (ACR)* to improve computational performance of both Planning and Reinforcement Learning (RL). ACR is an algorithm-agnostic, abstract data representation that maps objects to action categories (groups of actions), inspired by the psychological concept of *action codes*. We validate our approach in StarCraft and Lightworld domains; our results demonstrate several benefits of ACR relating to improved computational performance of planning and RL, by reducing the action space for the agent.

## 1 Introduction

Research in Psychology has shown that humans handle the complexity of the real world by biasing or constraining their action choice at a given moment based on known object-related actions. In particular, recent fMRI studies show that the human brain uses *action codes* – automatically evoked memories of prototypical actions that are related to a given object – to bias or constrain expectation on upcoming manipulations [Schubotz *et al.*, 2014]. In effect, given an object, our brain simplifies the action selection process by constraining the decision to a predefined set of known actions. For instance, a knife and an apple seen together evoke the action codes of "cutting apple with knife" and "peeling apple with knife".

Our work presents an analogous mechanism for computational agents, showing that automatically generated action groupings can be used to improve the computational efficiency of both task planning and learning by constraining the action space. We present the *Action Category Representation (ACR)*: an algorithm-agnostic, abstract data representation that encodes a mapping from objects to action categories (groups of actions) for a task. Specifically, we incorporate the idea of *action codes* as the action categorization mechanism. We formally define an action code as the tuple:

$$((o_1, o_2...o_j), (a_1, a_2...a_k))$$

Where  $(o_1,...o_j)$  represents a set of objects and  $(a_1,...a_k)$  represents the set of actions associated with them for the task. For instance, the action code corresponding to the knife and

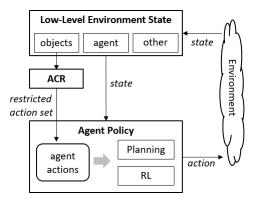


Figure 1: Objects in the low-level environment state are mapped via ACR to action categories to restrict the action set used in the planning or RL techniques

apple example above is ((apple, knife), (peel, cut)). In our work, we use action codes to build the Action Category Representation that can be used to improve computational performance in both task planning and reinforcement learning.

Action codes are closely related to the concept of object affordances [Gibson, 1977; McGrenere and Ho, 2000], which are defined as action possibilities available to the agent for a given object. Affordances function by priming specific actions for the user by virtue of the object's physical properties (shape, size etc.). In contrast, action codes do not derive from the physical properties of objects, rather from the associative memories of what we use the objects for during everyday tasks. Thus, the notion of affordances is often independent of the task [Ellis and Tucker, 2000; Tucker and Ellis, 1998] while action codes take the task into account. ACR builds on the notion of action codes, enabling an agent to learn object-action mappings based on prior experience.

Within a computational framework, the primary benefit of ACR is to reduce the choice of actions the agent must consider. Thus, ACR serves as a layer of abstraction between low-level state information and the learning or planning technique used to control the agent (Figure 1). In this paper we:

- 1. describe the process of constructing ACR from the agent's experience or human demonstration of a task;
- 2. show that ACR has formal computational bounds that

guarantee its use leads to at the worst case the same, and in the common case much improved, computational performance over traditional techniques that consider objects and actions without categorization;

- present the computational benefits of using ACR in conjunction with PDDL planning to reduce planning time; and
- present the computational benefits of using ACR with Q-learning to achieve improved learning performance.

We validate ACR performance in two virtual domains: StarCraft and Lightworld. We conclude the paper with a discussion of our work and potential future uses of ACR.

### 2 Related Work

In this section, we position our paper in relation to existing work.

### 2.1 Affordance Learning

As discussed above, the concept of action codes is closely related to action affordances and affordance learning. Affordances model relationships between individual object *properties* (shape, size, color etc.) to actions and observed effects and are formally defined as

(effect,(object,behavior)) [Şahin et al., 2007]. In contrast, action codes and by extension ACR, relate only semantic labeling and a holistic perception of objects such as "cup" or "box" to appropriate actions for a task.

Traditional approaches to affordance learning often involves "behavioral babbling" [Stoytchev, 2005; Montesano et al., 2008; Lopes et al., 2007] wherein the agent physically interacts with objects in a goal-free manner to discover their affordances. Hence, the resulting affordance representation is dissociated from a task, focusing instead on object properties. Such approaches involve several agent-object interactions affecting the scalability of the learning process, making it unfeasible in situations where there is an implicit cost or time constraint on the robot. ACR helps mitigate this cost by the grouping of actions into categories.

Two works closest to our approach are [Kjellström et al., 2011] and [Sun et al., 2010]. In [Kjellström et al., 2011], Kjellström et al. describe an approach to visual object-action recognition that use demonstrations to categorize semantically labeled objects based on their functionality. This approach bridges the gap between affordance learning and task context since the learning is coupled with a task demonstration. However, it is unclear how the system would incorporate previously unseen objects unless they are observed from additional demonstrations. For instance, given a demonstration of pouring water into a "cup", the agent would require additional demonstrations to identify the similar functionality of a "bowl".

Sun et al. in [Sun et al., 2010] learn visual object categories for affordance prediction (Category-Affordance model), reducing the physical interactions with the objects. They use visual features of objects to categorize them on the basis of their functionality. However, it is unclear how the agent would deal with changing features and categories [Min et al., 2016],

since the model is learned offline as compared to ACR which allows online learning of new objects and categories (Details in Sec 3). Regardless, their approach highlights some of the benefits of categorization on the scalability of learning, which motivates our work.

### 2.2 Precondition Learning

Preconditions can be expressed using predicates which may or may not relate to object affordances. For instance, "At" or "isEmpty" are object states whereas "graspable" is an affordance predicate [Lörken and Hertzberg, 2008]. Object-Action Complexes or OACs [Geib et al., 2006] include instances of affordances as preconditions in the OAC instantiation. Their approach learns an "object" after physical interaction with it, i.e, there is no notion of an object prior to the interaction. For instance, the representation of a cube is learned after the agent grasps a planar surface. Other approaches such as [Ekvall and Kragic, 2008] learn high-level task constraints and preconditions from demonstrations. In contrast to these approaches, ACR categorizes objects on the basis of action codes to improve planning performance as well as the learning performance of RL algorithms.

## 2.3 Learning from Demonstration

Human demonstrations have been used for both high-level task learning and low-level skill learning [Chernova and Thomaz, 2014]; a traditional assumption of LfD is that the human demonstrator is an expert, and the demonstrations are examples of desirable behavior that the agent should emulate. Our work focuses on high level task learning, but considers demonstrations more broadly as examples of what the agent can do, rather than what it should. This interpretation of the data enables our technique to benefit even from non-expert human users. Demonstration errors can be classified to one of 3 categories [Chernova and Thomaz, 2014]: Correct but suboptimal (contains extra steps), conflicting or inconsistent (user demonstrates 2 different actions from the same state) and entirely wrong (user took a wrong action) and we demonstrate the robustness of ACR to suboptimal demonstrations in Sec 6.

### LfD in planning

Abdo et al. in [Abdo et al., 2012] discuss the learning of predicates by analyzing variations in demonstrations. The learned predicates are then applied to plan for tasks and accommodate for environmental changes. Kadir et al. in [Uyanik et al., 2013] demonstrates execution of a task by leveraging human interactions. The agent interacts with all of the objects using all of the precoded behaviors in its repertoire and uses forward chaining planning to accomplish the task goal. However, with increasing number of behaviors and objects, the search space for the planner can become quite large. Our approach using ACR can help reduce the action space making planning easier.

#### LfD in RL

Thomaz and Breazeal in [Thomaz et al., 2006] discuss the effect of human guidance on an RL agent. Similar to our approach with ACR, the teacher guides the action selection process to reduce the action space for the RL agent. While

both expert and non-expert guidance improved performance when compared to unguided learning, the final performance was sensitive to the expertise of the teacher.

Another well-known approach to integrating LfD and RL is Human-Agent Transfer or HAT [Taylor et al., 2011]. Their approach uses a decision list to summarize the demonstration policy with a set of rules. However, it is sensitive to the number and optimality of the demonstrations [Suay et al., 2016; Brys et al., 2015]. We compare ACR to HAT in Sec 6 to demonstrate the benefits of ACR in terms of quantity and quality of the demonstrations.

## **Object Focused Approaches**

In general, recent work in AI and Robotics has increasingly focused on modeling state not simply as a vector of features, but as a set of objects, such as Object-Oriented MDPs (OO-MDP), leading to improved computational performance due to data abstraction and generalization. In the context of reinforcement learning (RL), Object-focused Q learning (Of-Q) represents the state space as a collection of objects organized into object classes, leading to exponential speed-ups in learning over traditional RL techniques [Cobo et al., 2013]. More closer to our work, human input containing object-action associations has been used to effectively guide policy learning in Mario [Krening et al., 2016]. The input advice is analogous to action codes, eg. "Jump over an enemy". Approaches such as [Barth-Maron et al., 2014; Cruz et al., 2014; Wang et al., 2013] have used affordances in RL to prune the action space and improve learning. However, the formalisms in all these approaches differ from ACR and were not extended beyond Reinforcement Learning to planning of tasks.

## **Action-Category Representation (ACR)**

The objective of ACR is to categorize objects based on the action codes of a task. In this section, we first present an example that illustrates the functionality of ACR and contrasts it with object affordance models. We then present the ACR

As an example, consider the task of packing a cardboard box during clean-up. The action codes for the cardboard box in this case are ((box),(close)) and ((box),(move)). Based on these action codes, ACR groups the actions close and move into a single action category associated with the item box. One of the benefits of ACR is that other objects sharing the same action codes, such as cooking pot in a dish-clearing task and suitcase in a travel-packing task, also become associated with the same action category, enabling the agent to reason about groups of similar objects across tasks that share action codes, despite the physical dissimilarities of the objects.

In our human example, a person seeing a knife and an apple may be primed to cut or peel, but may also select to ignore this bias and choose to wash the apple instead. In the agent's case, we similarly have the choice of treating ACR as a hard constraint on the actions available to the agent or as a flexible bias. In the sections below, we show how planning is well suited to use ACR as a hard constraint, and how ACR can be naturally combined with RL as a bias. We discuss possible extensions of this view in the conclusion of the paper.

In this paper, we show how ACR can be utilized in two ways. First, during task planning or learning, ACR improves computational efficiency by pruning the action space. Second, given an object not previously seen by the agent, ACR reduces the number of agent-object interactions required to learn its action associations for the task. Note that, in humans, action codes act as a bias and not a strict restriction on actions. In other words, a person seeing the knife and apple next to each other is primed to perform the actions cut and peel, but may override this bias too and put away the apple instead. In the agent's case, we have the option to treat ACR as a hard constraint on the actions available to the agent or as a flexible bias that also allows re-expanding the action set. In the sections below, we show how planning is well suited to use ACR as a hard constraint, and how ACR can be naturally combined with RL as a bias. We discuss possible extensions of this view in the conclusion of the paper.

To construct ACR, the agent requires observations of objects in its environment and what actions are related to each object. These observations can be gained either through the agent's own exploration of the environment, or, more effectively, from a human teacher performing demonstrations of the task. During the observation phase, the agent maintains a log of action codes based on the actions performed and the objects that the actions were executed upon. We define O as the set of all objects in the task environment, and A as the set of all actions pertaining to the task. An observation log consists of a set of action codes and is represented by  $L = \{\hat{c}_1, \hat{c}_2, ... \hat{c}_n\}$ , where each timestep in the log is represented by an action code  $\hat{c}_i = ((o_1, o_2...o_j), (a_1, a_2...a_k))$ with  $o_i \in O$  and  $a_k \in A$ .

The act of building object-action relations can be formulated as a bipartite graph partitioning problem involving the action set A and the objects set O. Given a graph G(V, E), with vertices V and edges E, the graph is bipartite when the vertices can be separated into two sets, such that  $V = A \cup O$ ,  $A \cap O = \emptyset$ , and each edge in E has one endpoint in A and one endpoint in O. In the context of ACR, A represents actions, O represents objects and an edge  $\{a_i, o_i\}$  exists if action  $a_i \in A$ and object  $o_i \in O$  co-occur within any action code  $\hat{c}_k \in L$ . For instance, the action code ((box), (push)) is represented by an edge from the action *push* to the object *box*. The resulting bipartite graph has a many-to-many association between objects and actions (Figure 2 left). In ACR, the bipartite graph is generated incrementally from the action codes in the observation log.

The main computational units of ACR are action categories, defined by a group or set of actions  $A^c \subseteq A$ . Given the bipartite graph above, for a given action  $a_i \in A$ , let  $\hat{O}_{a_i}$ represent the set of objects for which that action co-occurs in some action code (i.e. the edge  $\{a_i, o_k \in \hat{O}_{a_i}\}$  exists). Then we define an action category  $A^c$  as:

$$A^c = \{a_j : \bigcup \hat{O}_{a_j} = \bigcap \hat{O}_{a_j}\}$$

 $A^c=\{a_j:\bigcup\hat{O}_{a_j}=\bigcap\hat{O}_{a_j}\}$  This is interpreted as, "The set of all actions  $a_j$  such that union over all  $\hat{O}_{a_i}$  is equal to intersection over all  $\hat{O}_{a_i}$ ". In other words, the action category  $A^c$  contains a set of actions that are associated to the same set of objects, allowing us to group all those actions as one set. If we consider action cat-

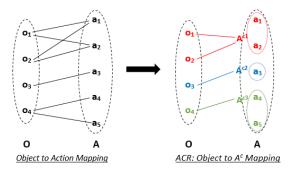


Figure 2: Bipartite graphs representing the relationship from objects to actions (left), and objects to action categories (right)

egories as vertices themselves, then what results is a reduced one-to-many bipartite graph between action categories  $A^c$  and the set of objects O (Figure 2 right), which is the representation we refer to as ACR. Note that the entire set of actions can be grouped into categories such that  $A = A^{c_1} \cup A^{c_2} \cup A^{c_3} \cup ...A^{c_n}$ . We define  $C = \{A^{c_1}, A^{c_2}...A^{c_n}\}$  as the set of all action categories learned from observations. Note that, as in most prior work, we assume single-parametric actions <sup>1</sup> [Montesano  $et\ al.$ , 2008; Ugur  $et\ al.$ , 2011; Şahin  $et\ al.$ , 2007], and that preconditions and effects of those actions are known and can be perceived when planning [Agostini  $et\ al.$ , 2015; Ekvall and Kragic, 2008].

The construction of ACR is an online process, allowing learning of new objects and action categories over time with changes in the environment or task. As new action codes are learned, objects or actions can be incorporated by adding them to the graph, along with corresponding edges. A new action category  $A^{c_i}$  may be added to C when a new combination of associated actions is discovered, such that  $A^{c_i} \neq A^{c_k} \ \forall A^{c_k} \in C$ . The resulting representation provides an automatically-generated online grouping of objects into categories based on action codes.

In this paper, we discuss characteristics of ACR that contribute to its novelty and significance:

- groups actions based on action codes in order to reduce the action space for the agent,
- 2. contains and appropriately represents algorithmagnostic information for planning, as well as RL, to improve their computational performance,
- 3. minimizes agent-object interactions for learning the action associations of a new object; and
- 4. requires one or few human demonstrations and is robust to the optimality of these demonstrations.

## 4 Computational Performance Analysis

In this section, we present performance guarantees of ACR; in the following section we then validate our findings with case studies in StarCraft and Lightworld [Konidaris and Barto, 2007] domains.

## 4.1 Mathematical Analysis

We define the total number of actions in a domain to be |A|=n, allowing us to bound the total possible action categories to be  $|C|\leq 2^n-1$ , representing all possible action combinations from 1 to n actions. Then a given task involves a subset of these action categories  $S\subseteq C$  and a set of objects O. The agent is assigned the task of learning S and categorizing the objects in O from observations of action codes.

One of the benefits of ACR is seen when the agent encounters a new set of objects O' (not previously seen) and must discover which actions in A are related to each object in O' for the task execution. Below we present performance analysis of ACR and the baseline that uses no action categorization, with respect to the number of agent-object interactions prior to learning all the actions related to an object for the task  $(A_{obj})$ . Fewer  $A_{obj}$  is computationally preferred since this reduces the number of agent-object interactions, making the learning or planning faster.

### $A_{obj}$ Without Categorization (Baseline)

Without categorization, each action is considered independently, in which case to determine the set of actions applicable to a new object the agent must test out all |A|=n actions on that object. That is,  $A_{obj}=n$ .

## $A_{obj}$ With Action Categories (ACR)

The use of ACR can improve computational performance through action selection, enabling the agent to more effectively identify (or rule out) object interactions. In the presence of action categories our goal is to, with as few actions as possible, identify the category of a newly discovered object. To do so, we select actions from A, and for every attempted action that is unassociated (or associated) with the object we eliminate any action category in S with (or without) that action from further testing. We use entropy as a measure of the most informative action to test so as to eliminate as many action categories as possible with each action tested. The entropy of an action a is given by:

$$H(a) = -p \log(p) - (1-p) \log(1-p)$$
 Where,  $p = \frac{\sum_{i=1}^{n} |\{a\} \cap A^{c_i}|}{|S|}$ 

The term p denotes the probability that the action categories contain the action a which is used to compute the entropy. Then the action that minimizes entropy is the most informative action. Therefore, the action  $\hat{a}$  chosen for testing is given by:

$$\hat{a} = \arg\min_{a \in A} H(a)$$

In the best case, the category may be learned with a single action and hence the lower bound on  $A_{obj}$  is 1. The worst case upper bound on  $A_{obj}$  remains  $n: 1 \le A_{obj} \le n$ .

That is, with categorization the performance is *never* worse but typically better than without categorization. In practice, S is usually a small subset of C, and therefore  $A_{obj} << n$ . In the case that a new action category must be learned, which occurs rarely in closed world domains such as StarCraft and Lightworld,  $A_{obj} = n$ . In fact, in the experiments described below the agent obtains all possible action codes even from a single demonstration, allowing all the relevant action categories to be known prior to planning and learning.

<sup>&</sup>lt;sup>1</sup>While it is possible to decompose multi-parametric actions to single-parametric actions as described in [Bach *et al.*, 2014], we currently do not model them explicitly within ACR.







Figure 3: Refineries in StarCraft: Terrans, Zergs and Protoss (left to right), showing their distinct physical appearances

## 5 Case Study in StarCraft

In this section, we first briefly describe our domain and highlight the complexity of the problem before discussing the computational benefits of using ACR with planning.

StarCraft is a real-time strategy game which involves managing one of the 3 diverse civilizations (Terrans, Protoss and Zergs), producing buildings and units while destroying all of your opponents. Across the 3 civilizations, there are over 100 diverse units/buildings and reducing the number of agent-object interactions in this case makes the problem of task planning in StarCraft much more tractable. Figure 3 shows a "refinery", one of the buildings in StarCraft that is used to extract a gaseous mineral. Given the distinct appearances of the buildings across the 3 civilizations, it may be challenging to identify their related actions by using only physical features without any actual interaction.

#### 5.1 ACR Extracted from StarCraft

We extracted ACR from one human demonstration in the Terrans civilization where the teacher successfully completed an in-game mission of creating a defense. Replay logs that summarize the action codes within the demonstration are readily available for StarCraft and are used to construct the ACR. In the case of physical systems, it is possible to extract action codes from human demonstrations using verbal communication during the demonstration or approaches such as [Gupta and Davis, 2007].

Table 1 shows the complete ACR built from the human demonstration, highlighting the different action categories and object mappings to the action categories. We use this learned representation in the following section to describe its computational benefits.

## 5.2 Computational Benefits of ACR with Planning

We demonstrate two computational benefits of ACR with planning in terms of:

- 1. reduced number of object interactions or  $A_{obj}$  required to learn all object-action associations prior to planning in StarCraft (Exploration phase)
- improved planning performance due to reduced action space, demonstrated with combat formations in Star-Craft

#### **Benefits of ACR During Exploration Phase**

We demonstrate the benefits of ACR on  $A_{obj}$  using build order planning. Build orders dictate the sequence in which units and structures are produced. Prior to planning for a task, there is usually an exploration phase during which the actions associated with the objects (in this case, for the units/structures

	Actions: Build, Repair, Stop	
A <sup>c</sup> <sub>1</sub> ←	Objects: Barracks, Academy, Refinery, Supply Depot, Engineering Bay, Factory	
	Actions: Train, Move, Stop	
A <sup>C</sup> <sub>2</sub>	Objects: SCV, Marine, Vulture	
	Actions: Gather	
A <sup>C</sup> <sub>3</sub>	Objects: Minerals, Gas	
	Actions: Research, Use_tech	
A <sup>C</sup> <sub>4</sub>	Objects: Stim Packs, Spider Mines	
	Actions: Upgrade	
A <sup>C</sup> 5	Objects: Ion Thrusters	

Table 1: ACR built from human demo

Civilization	Number of objects explored	Total A <sub>obj</sub> w/o categorization	Total A <sub>obj</sub> w ACR
Terrans	4	36	10
Protoss	7	63	21
Zergs	9	81	28

Table 2: Total  $A_{obj}$  for the different build order exploration phases

specified in the build order) are first identified [Ugur et al., 2011]. This exploration phase adds to the overall planning time. Thus, a reduced  $A_{obj}$  in the exploration phase would reduce overall planning time. We compare ACR and baseline (without categorization) on  $A_{obj}$ , during exploration phase of build order planning.

We explore build orders from all 3 civilizations. Table 2 shows the total number of agent-object interactions during the exploration phase, along with the number of previously unseen objects for which the object-action relations had to be learned.

As shown in the Table 2, the number of object interactions with ACR is significantly reduced compared to the baseline approach that does not use categorization. In the baseline case, every action (of the 9 actions shown in Table 1) has to be attempted on each new object in the build order to discover all of its associations which is mitigated by the use of action categories. The results obtained here highlight the benefits previously discussed in the mathematical analysis of Section 4.1. While the feedback for an invalid action in StarCraft is instantaneous and incurs no significant time cost, in other domains such as task execution with robots, there may be implicit time and cost constraints associated with each interaction. Hence, with ACR it is possible to minimize interactions with the environment by grouping actions.

### Improved Planning Performance with ACR

In this section, we combine ACR with an existing off-the-shelf PDDL Planner (Fast Forward) [Helmert, 2006] to

demonstrate how ACR reduces the action space that the planner has to contend with, thus reducing the planning time.

For this evaluation we use a combat formation problem where combat units (Dragoons) have to form a particular arrangement on a section of the battlefield. We compared the classical planning approach without action categories (baseline) to planning with ACR. We increase the number of Dragoons and demonstrate its effect on the two planning approaches. Figure 4 shows a sample initial and goal states for the Dragoon formation.

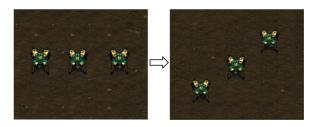


Figure 4: Combat formation problem showing initial (left) and goal (right) states of formation using 3 Dragoons

The overall pipeline for combining ACR with the planner is shown in Figure 5. Contrary to the classical approach, ACR introduces the learned action categories into the domain and problem definitions for planning, leading to computational improvements. The classical planning approach instantiates each Dragoon as a separate entity with distinct variables, while the ACR based approach instantiates all the Dragoons in terms of the action category that they are mapped to, which is  $A^{c_2}$  in this case (Table 1). The domain and problem definitions are thus automatically generated from ACR. The plan is then generated using the domain and problem definitions.

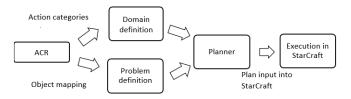
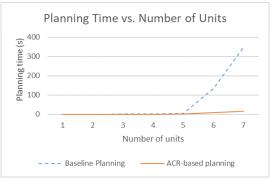
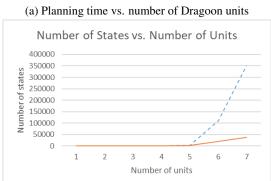


Figure 5: Pipeline for integrating ACR with PDDL planner for planning in StarCraft

As shown in Figure 6, increasing number of Dragoons from 1 to 7, exponentially increases the search time for the classical planning approach as compared to the ACR-based approach (Figure 6a). This is because the number of states increases exponentially with the number of Dragoons. For instance, on a 5x5 grid, the number of states for 3 Dragoons is  $\binom{25}{3} * 3 = 6900$  for the classical planning approach and  $\binom{25}{3} = 2300$  for the ACR-based planner since ACR instantiates all Dragoons in terms of their action category. Similarly, in the case of forward chaining planners such as [Uyanik *et al.*, 2013], ACR can help reduce the branching factor from  $n_o * |A|$  (where,  $n_o$  is the total number of objects or Dragoons in this case, and |A| is total number of actions) to  $\sum_1^n |o^{c_i}| * |A^{c_i}|$  (where,





(b) Number of search states vs. number of Dragoon units

ACR-based planning

Baseline Planning

Figure 6: Graphs showing effect of number of Dragoons on planning time (Fig 6a) and number of search states (Fig 6b) for the baseline planning and ACR-based planning approaches

 $|o^{c_i}|$  indicates number of objects mapped to action category  $A^{c_i}$ ). Thus, ACR leads to computational benefits by pruning the action space the planner has to contend with.

## 6 Case Study in Lightworld

In this section, we discuss the computational benefits of applying ACR with RL. We first discuss our domain design, inspired by the Lightworld domain used in the Options RL framework [Konidaris and Barto, 2007]. We then discuss the benefits of applying ACR with RL.

Figure 7 shows a sample domain. The game consists of a 7x8 grid of locked rooms with some doors operated by a switch and some doors operated by a key. The goal of the agent is to unlock the doors and move to the final reward. The agent has to move over the button or the key to either press or pick up the object. There are also spike pits that the agent needs to avoid while navigating the room. The agent receives a reward of +100 for reaching the goal state, and a negative reward of -10 for falling into spike pits which are terminal states. Additionally, the agent receives a negative step reward of -0.04. There are a total of 6 actions with each of the four grid directions, a pickup action and press action. The environment is deterministic and unsuccessful actions do not change the state.

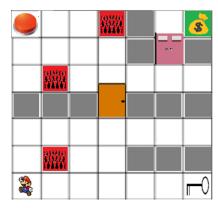


Figure 7: Example domain

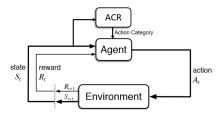


Figure 8: Pipeline for integrating ACR with RL

## 6.1 Integrating ACR with RL

Figure 8 shows the integration of ACR with general RL algorithms. ACR influences the action selection step, given an observed state. With ACR, the agent chooses an action from within an action category  $A^{c_i} \subset A$  of the known objects it can interact with, given the state. For previously unseen objects whose action categories are unknown, the agent chooses the entropy-wise selected action as described in Sec 4.1 to simultaneously infer the action category of the objects during learning. If the agent cannot interact with any objects, it chooses from the non object-related set of actions (analogous to an "agent" action category) given by,  $A - \bigcup_{i=1}^m A^{c_i}$  where m denotes the number of learned action categories. This reduces the action space for the agent.

As noted in Sec 3, we use ACR with RL to bias the initial learning rather than applying it as a hard constraint over the entire learning phase. We achieve this by allowing ACR to influence the action choice for a fixed number of episodes, denoted by  $N_{ACR}$ . This allows the agent to leverage the reduced action space while also learning the optimal policies in states where the ACR-guided policy may be suboptimal.

## 6.2 Computational Benefits of ACR with RL

We compare three RL agents: Q-learning, Q-learning with ACR and Q-learning with Human-Agent Transfer (HAT [Taylor *et al.*, 2011]). HAT uses human demonstrations to learn strategies (decision list) from demonstration summaries. We used Q-learning with  $\epsilon$ -greedy exploration with  $N_{ACR}=50$ ,  $\alpha=0.25$ ,  $\gamma=0.99$  and  $\epsilon=0.1$ .

We used 5 expert and 5 suboptimal demonstrations separately, to compare the effect of demonstration quality on ACR and HAT. Suboptimal demonstrations refer to where the

Method	Episodes to	Avg. number of actions
	convergence	taken to convergence
HAT w. expert demos	13	11.38
ACR	63	21.27
Q-learning	193	49.24
HAT w. suboptimal demos	167	45.71

Table 3: Comparison of the different approaches based on convergence episode and average number of actions taken by the agent (bold values correspond to ACR)

demonstrator either failed to complete the goal or took a suboptimal path to reach the goal state. In all experiments below, the ACR was built from a single demonstration that exposed the agent to all object-related actions necessary to complete the game. That is, the ACR encodes what the agent *can* do, rather than what the agent *should*. This imposes a more relaxed constraint on the teacher since it is easier to show the "rules" rather than the "strategy" which requires an expert. Importantly, unlike most existing approaches, the ACR built from expert or suboptimal demonstrations **do not differ** if the agent learned the same rules from either demonstration.

Additionally, to evaluate the benefit of entropy-based action selection in RL, the ACR-based approach treats keys and switches as previously unseen objects whose action-categories are unknown and must be simultaneously inferred during the course of the learning process.

We demonstrate three benefits of using ACR with RL:

- robustness to demonstration quality: we show that ACR has a higher learning rate compared to HAT and Q learning when trained on suboptimal demonstrations
- learning from few demonstrations: we show that ACR learns more efficiently than both HAT and Q learning when only a single demonstration is available
- improved performance when combining ACR and HAT: we show that best overall performance is achieved when ACR is used to improve the performance of other LfD methods, in this case Human-Agent Transfer.

### **Effect of Demonstration Quality on ACR:**

As shown in 9, ACR performs much better than Q-learning approach and HAT trained on suboptimal demonstrations. ACR also minimizes the number of attempted actions as shown in Table 3. However, it does not perform better than HAT that uses expert demonstrations. This is because, with enough expert demonstrations, the information contained within ACR can be implicitly learned in the form of rules. Since ACR does not fully leverage the capabilities of a good teacher it does not outperform HAT that uses multiple expert demonstrations.

However, HAT is quite sensitive to the optimality of the demonstration. The starting reward for HAT is dependent on the teacher performance. As shown in Figure 9, the starting reward for HAT trained on expert demonstrations is much higher when compared to HAT trained on suboptimal demonstrations.

To summarize, in cases involving non-expert users, ACR can leverage the rules of the task in order to improve learning performance over the baseline approaches.

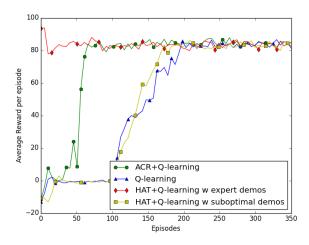


Figure 9: Comparison of Q-learning, HAT + Q-learning with 5 expert, 5 suboptimal demonstrations and ACR + Q-learning

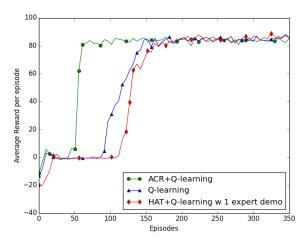


Figure 10: Comparison of Q-learning, HAT + Q-learning with single expert demonstration and ACR + Q-learning

### Effect of Number of Demonstrations on ACR:

Given only a single expert demonstration, HAT fails to accurately summarize the source policy (Figure 10). The building of the decision list in HAT requires more data depending on the complexity of the domain. However, ACR was able to perform better than the baseline approaches with one expert demonstration. Hence, this makes ACR a feasible approach when there are not enough demonstrations available to learn a good demonstration policy.

Figure 11 summarizes the effects of number and quality of demonstrations on learning performance of the different approaches. Q-learning is also shown for comparison.

### Combining ACR with HAT

ACR is an algorithm and domain independent representation; as a result, one of its strengths is that it can be easily combined with complex learning methods, including ones that in themselves influence action selection, such as Human-Agent

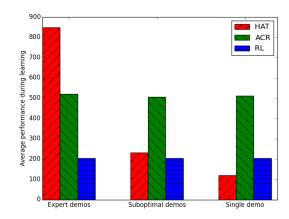


Figure 11: Summary of the average learning performances based on number and quality of demonstrations for HAT and ACR.

Q-learning (RL) also shown for comparison.

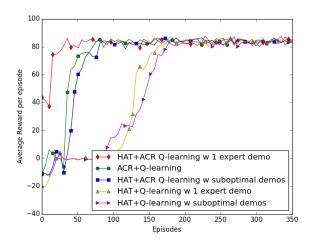


Figure 12: Performances of HAT + ACR for single expert demonstration and 5 suboptimal demonstrations compared to the baselines that consider the two separately

Transfer.

HAT consists of 3 steps: demonstration, policy summarization and independent learning. Given one or more demonstrations, the teacher's behavior is summarized in the form of a decision list and then used to bootstrap the learning. We utilize the *Extra Action* method from [Taylor *et al.*, 2011], in which the agent executes the action suggested by the decision list for a fixed number of initial episodes before running regular RL.

HAT can bootstrap the learning in states where a good policy is obtained from the demonstrations, while for the "bad" states in which the demonstrator's performance was suboptimal, ACR helps accelerate the learning of the optimal policy by reducing the action space. We combine ACR and HAT by verifying that the action suggested by the decision list conforms with the retrieved action category. It does so, by making two checks:

- 1. In states with objects: action selection is restricted to the union of all non-object actions (e.g. movement) and object-related actions within ACR, ensuring that the agent does not try an incorrect action on the object (e.g. "pick" button).
- In states without objects: action selection is restricted to non-object actions.

As in the *Extra Action* method, the above action selection method is used to bias exploration early in the learning process before continuing to classical RL using  $\epsilon$ -greedy Q-learning.

The results of this method are presented in Figure 12, showing improved performance when ACR is combined with HAT. Combining ACR with HAT trained on a single expert demonstration improves the learning performance beyond the case where either of the two approaches are considered separately. Hence, by combining ACR with HAT, it is possible to reduce the effect of number of demonstrations and their optimality on HAT, while also allowing ACR to maximally utilize the teacher demonstrations.

### 7 Conclusion and Future Work

To conclude, we presented the Action-Category Representation that allows online categorization of objects to action categories based on action codes. Our results demonstrate some of the key benefits of ACR in terms of reduced action space resulting from the action groupings, computational improvements when used with planning and RL, and reduced demonstration requirements with robustness to demonstration errors.

While the domains described here are discrete in nature, ACR is also applicable to continuous domains by discretizing the state space into states where interaction with an object is possible/not possible. For instance an object may be interacted with, if the agent is within a certain distance of it. Approaches such as [Mugan and Kuipers, 2008] have discussed discretization of continuous state spaces for RL and in this manner, ACR can also be extended to continuous domains.

In our future work, we aim to address some of the limitations of our work in its current form. Since ACR currently models single parameter actions, it limits the applicability of ACR to real-world tasks. We plan to address this by incorporating multi-parameter actions by decomposing them into single-parameter actions [Bach *et al.*, 2014]. Additionally, future work will explore the possibility of using ACR as a bias with planning (as opposed to the hard constraint on action selection), by utilizing "Ontology-Repair" [McNeill and Bundy, 2007] to update ACR and improve its flexibility. Finally, we wish to extend the application of ACR to Deep Learning for computational improvements in learning performance.

### 8 Acknowledgement

This work is supported by NSF IIS 1564080.

### References

- [Abdo et al., 2012] Nichola Abdo, Henrik Kretzschmar, and Cyrill Stachniss. From low-level trajectory demonstrations to symbolic actions for planning. In *ICAPS Workshop on Combining Task and Motion Planning for Real-World App*, 2012.
- [Agostini et al., 2015] Alejandro Agostini, Mohamad Javad Aein, Sandor Szedmak, Eren Erdal Aksoy, Justus Piater, and Florentin Würgütter. Using structural bootstrapping for object substitution in robotic executions of human-like manipulation tasks. In *Intelligent Robots and Systems (IROS)*, 2015 IEEE/RSJ International Conference on, pages 6479–6486. IEEE, 2015.
- [Bach et al., 2014] Patric Bach, Toby Nicholson, and Matthew Hudson. The affordance-matching hypothesis: how objects guide action understanding and prediction. Frontiers in human neuroscience, 8, 2014.
- [Barth-Maron *et al.*, 2014] Gabriel Barth-Maron, David Abel, James MacGlashan, and Stefanie Tellex. Affordances as transferable knowledge for planning agents. In 2014 AAAI Fall Symposium Series, 2014.
- [Brys *et al.*, 2015] Tim Brys, Anna Harutyunyan, Halit Bener Suay, Sonia Chernova, Matthew E Taylor, and Ann Nowé. Reinforcement learning from demonstration through shaping. In *IJCAI*, pages 3352–3358, 2015.
- [Chernova and Thomaz, 2014] Sonia Chernova and Andrea L Thomaz. Robot learning from human teachers. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 8(3):1–121, 2014.
- [Cobo et al., 2013] Luis C Cobo, Charles L Isbell, and Andrea L Thomaz. Object focused q-learning for autonomous agents. In Proceedings of the 2013 international conference on Autonomous agents and multi-agent systems, pages 1061–1068. International Foundation for Autonomous Agents and Multiagent Systems, 2013.
- [Cruz et al., 2014] Francisco Cruz, Sven Magg, Cornelius Weber, and Stefan Wermter. Improving reinforcement learning with interactive feedback and affordances. In *Development and Learning and Epigenetic Robotics (ICDL-Epirob), 2014 Joint IEEE International Conferences on*, pages 165–170. IEEE, 2014.
- [Ekvall and Kragic, 2008] Staffan Ekvall and Danica Kragic. Robot learning from demonstration: a task-level planning approach. *International Journal of Advanced Robotic Systems*, 5(3):33, 2008.
- [Ellis and Tucker, 2000] Rob Ellis and Mike Tucker. Microaffordance: The potentiation of components of action by seen objects. *British journal of psychology*, 91(4):451–471, 2000.
- [Geib et al., 2006] Christopher Geib, Kira Mourao, Ron Petrick, Nico Pugeault, Mark Steedman, Norbert Krueger, and Florentin Wörgötter. Object action complexes as an interface for planning and robot control. In *IEEE RAS International Conference on Humanoid Robots*, 2006.

- [Gibson, 1977] James J Gibson. Perceiving, acting, and knowing: Toward an ecological psychology. *The Theory of Affordances*, pages 67–82, 1977.
- [Gupta and Davis, 2007] Abhinav Gupta and Larry S Davis. Objects in action: An approach for combining action understanding and object perception. In *Computer Vision and Pattern Recognition*, 2007. CVPR'07. IEEE Conference on, pages 1–8. IEEE, 2007.
- [Helmert, 2006] Malte Helmert. The fast downward planning system. *Journal of Artificial Intelligence Research*, 26:191–246, 2006.
- [Kjellström et al., 2011] Hedvig Kjellström, Javier Romero, and Danica Kragić. Visual object-action recognition: Inferring object affordances from human demonstration. Computer Vision and Image Understanding, 115(1):81–90, 2011.
- [Konidaris and Barto, 2007] George Konidaris and Andrew G Barto. Building portable options: Skill transfer in reinforcement learning. In *IJCAI*, volume 7, pages 895–900, 2007.
- [Krening et al., 2016] Samantha Krening, Brent Harrison, Karen M Feigh, Charles Isbell, and Andrea Thomaz. Object-focused advice in reinforcement learning. In Proceedings of the 2016 International Conference on Autonomous Agents & Multiagent Systems, pages 1447– 1448. International Foundation for Autonomous Agents and Multiagent Systems, 2016.
- [Lopes et al., 2007] Manuel Lopes, Francisco S Melo, and Luis Montesano. Affordance-based imitation learning in robots. In *Intelligent Robots and Systems*, 2007. *IROS* 2007. *IEEE/RSJ International Conference on*, pages 1015–1021. IEEE, 2007.
- [Lörken and Hertzberg, 2008] Christopher Lörken and Joachim Hertzberg. Grounding planning operators by affordances. In *International Conference on Cognitive Systems* (CogSys), pages 79–84, 2008.
- [McGrenere and Ho, 2000] Joanna McGrenere and Wayne Ho. Affordances: Clarifying and evolving a concept. In *Graphics interface*, volume 2000, pages 179–186, 2000.
- [McNeill and Bundy, 2007] Fiona McNeill and Alan Bundy. Dynamic, automatic, first-order ontology repair by diagnosis of failed plan execution. *International Journal on Semantic Web and Information Systems (IJSWIS)*, 3(3):1–35, 2007.
- [Min et al., 2016] Huaqing Min, Chang'an Yi, Ronghua Luo, Jinhui Zhu, and Sheng Bi. Affordance research in developmental robotics: A survey. *IEEE Transactions* on Cognitive and Developmental Systems, 8(4):237–255, 2016.
- [Montesano *et al.*, 2008] Luis Montesano, Manuel Lopes, Alexandre Bernardino, and José Santos-Victor. Learning object affordances: from sensory–motor coordination to imitation. *IEEE Transactions on Robotics*, 24(1):15–26, 2008.

- [Mugan and Kuipers, 2008] Jonathan Mugan and Benjamin Kuipers. Continuous-domain reinforcement learning using a learned qualitative state representation. 2008.
- [Şahin *et al.*, 2007] Erol Şahin, Maya Çakmak, Mehmet R Doğar, Emre Uğur, and Göktürk Üçoluk. To afford or not to afford: A new formalization of affordances toward affordance-based robot control. *Adaptive Behavior*, 15(4):447–472, 2007.
- [Schubotz *et al.*, 2014] Ricarda I Schubotz, Moritz F Wurm, Marco K Wittmann, and D Yves von Cramon. Objects tell us what action we can expect: dissociating brain areas for retrieval and exploitation of action knowledge during action observation in fmri. *Frontiers in psychology*, 5, 2014.
- [Stoytchev, 2005] Alexander Stoytchev. Behavior-grounded representation of tool affordances. In *Robotics and Automation*, 2005. ICRA 2005. Proceedings of the 2005 IEEE International Conference on, pages 3060–3065. IEEE, 2005.
- [Suay et al., 2016] Halit Bener Suay, Tim Brys, Matthew E Taylor, and Sonia Chernova. Learning from demonstration for shaping through inverse reinforcement learning. In *Proceedings of the 2016 International Conference on Autonomous Agents & Multiagent Systems*, pages 429–437. International Foundation for Autonomous Agents and Multiagent Systems, 2016.
- [Sun et al., 2010] Jie Sun, Joshua L Moore, Aaron Bobick, and James M Rehg. Learning visual object categories for robot affordance prediction. *The International Journal of Robotics Research*, 29(2-3):174–197, 2010.
- [Taylor et al., 2011] Matthew Edmund Taylor, Halit Bener Suay, and Sonia Chernova. Using human demonstrations to improve reinforcement learning. In AAAI Spring Symposium: Help Me Help You: Bridging the Gaps in Human-Agent Collaboration, 2011.
- [Thomaz *et al.*, 2006] Andrea Lockerd Thomaz, Cynthia Breazeal, et al. Reinforcement learning with human teachers: Evidence of feedback and guidance with implications for learning performance. In *Aaai*, volume 6, pages 1000–1005, 2006.
- [Tucker and Ellis, 1998] Mike Tucker and Rob Ellis. On the relations between seen objects and components of potential actions. *Journal of Experimental Psychology: Human perception and performance*, 24(3):830, 1998.
- [Ugur *et al.*, 2011] Emre Ugur, Erhan Oztop, and Erol Sahin. Goal emulation and planning in perceptual space using learned affordances. *Robotics and Autonomous Systems*, 59(7):580–595, 2011.
- [Uyanik *et al.*, 2013] Kadir Firat Uyanik, Yigit Caliskan, Asil Kaan Bozcuoglu, Onur Yürüten, Sinan Kalkan, and Erol Sahin. Learning social affordances and using them for planning. In *CogSci*, 2013.
- [Wang et al., 2013] Chang Wang, Koen V Hindriks, and Robert Babuska. Robot learning and use of affordances in goal-directed tasks. In *Intelligent Robots and Systems (IROS)*, 2013 IEEE/RSJ International Conference on, pages 2288–2294. IEEE, 2013.