



# Intrinsic Rewards for Maintenance, Approach, Avoidance and Achievement Goal Types

1

2 **Paresh Dhakan<sup>1\*</sup>, Kathryn Merrick<sup>2</sup>, Inaki Rano<sup>1</sup>, Nazmul Siddique<sup>1</sup>**3 <sup>1</sup>Intelligent Systems Research Centre, Ulster University4 <sup>2</sup>School of Engineering and Information Technology, University of New South Wales5 **\* Correspondence:**6 [dhakan-p@ulster.ac.uk](mailto:dhakan-p@ulster.ac.uk)

7

8 **Keywords: intrinsic reward function, goal types, open-ended learning, autonomous goal**  
9 **generation, reinforcement learning.**

10

11 **Abstract**

12 In reinforcement learning, reward is used to guide the learning process. The reward is often designed  
13 to be task-dependent, and it may require significant domain knowledge to design a good reward  
14 function. This paper proposes general reward functions for maintenance, approach, avoidance and  
15 achievement goal types. These reward functions exploit the inherent property of each type of goal  
16 and are thus task-independent. We also propose metrics to measure an agent's performance for  
17 learning each type of goal. We evaluate the intrinsic reward functions in a framework that can  
18 autonomously generate goals and learn solutions to those goals using a standard reinforcement  
19 learning algorithm. We show empirically how the proposed reward functions lead to learning in a  
20 mobile robot application. Finally, using the proposed reward functions as building blocks, we  
21 demonstrate how compound reward functions, reward functions to generate sequences of tasks, can  
22 be created that allow the mobile robot to learn more complex behaviors.

23 **1 Introduction**

24 Open-ended learning, still an open research problem in robotics, is envisaged to provide learning  
25 autonomy to robots such that they will require minimal human intervention to learn environment  
26 specific skills. Several autonomous learning frameworks exist (Santucci, Baldassarre, and Mirolli  
27 2016) (Santucci, Baldassarre, and Mirolli 2010) (Bonarini, Lazaric, and Restelli 2006) (Baranes and  
28 Oudeyer 2010a) (Baranes and Oudeyer 2010b), most of which have similar key modules that include:  
29 (a) a goal generation mechanism that discovers the goals the robot can aim to achieve; and (b) a  
30 learning algorithm that enables the robot to generate the skills required to achieve the goals. Many of  
31 the autonomous learning frameworks use reinforcement learning (RL) as the learning module  
32 (Santucci, Baldassarre, and Mirolli 2016) (Santucci, Baldassarre, and Mirolli 2010) (Bonarini,  
33 Lazaric, and Restelli 2006). In RL, an agent learns by trial and error. It is not initially instructed  
34 which action it should take in a particular state but instead must compute the most favorable action  
35 using the reward as feedback on its actions. For many dynamic environments, however, it is not

36 always possible to know upfront which tasks the agent should learn. Hence, sometimes, it is not  
37 possible to design the reward function in advance. Open-ended learning aims to build systems that  
38 autonomously learn tasks as acquired skills that can later be used to learn user-defined tasks more  
39 efficiently (Thrun and Mitchell 1995) (Weng et al. 2001) (Baldassarre and Mirolli 2013). Thus, for  
40 an open-ended learning system, autonomous reward function generation is an essential component.  
41 This paper contributes to open-ended learning by proposing an approach to reward function  
42 generation based on the building blocks of maintenance, achievement, approach and avoidance goals.

43 Existing literature reveals two common solutions to address the problem of the autonomous reward  
44 function design or at least provides a level of autonomy in designing a reward function: (1) Intrinsic  
45 motivation (Singh, Barto, and Chentanez 2004) and (2) reward shaping (Laud and DeJong 2002)  
46 (Ng, Harada, and Russell 1999). Intrinsic motivation is a concept borrowed from the field of  
47 psychology. It can be used to model reward that can lead to the emergence of task-oriented  
48 performance, without making strong assumptions about which specific tasks will be learned prior to  
49 the interaction with the environment. Reward shaping, on the other hand, provides a positive or  
50 negative bias encouraging the learning process towards certain behaviors. Intrinsic motivation,  
51 although promising, has not been validated on large-scale real-world applications and reward shaping  
52 requires a significant amount of domain knowledge thus cannot be considered as an autonomous  
53 approach. As an alternative to these solutions, we propose reward functions based on the various  
54 types of goals identified in the literature. Although the concept of creating a reward function using  
55 goals is not new, this approach is often overlooked and has not been the main focus of the RL  
56 community. In our approach, different reward functions are generated based on the type of the goal,  
57 and since the reward functions exploit the inherent property of each type of goal, these reward  
58 functions are task-independent.

59 Goals have been the subject of much research within the Beliefs, Desires, Intentions community (Rao  
60 and Georgeff 1995) and the agent community (Regev and Wegmann 2005). A goal is defined as an  
61 objective that a system should achieve (van Lamsweerde 2001), put another way, a goal is the state of  
62 affairs a plan of action is designed to achieve. Goals range in abstraction from high-level to low-  
63 level, cover functional as well as non-functional aspects and can be categorized into hard goals that  
64 can be verified in a clear-cut way to soft goals that are difficult to verify (van Lamsweerde 2001).  
65 Examples of types of goals include achievement, maintenance, avoidance, approach, optimization,  
66 test, query, and cease goals (Braubach et al. 2005). Instead of classifying goals based on types, van  
67 Riemsdijk et al. (van Riemsdijk, Dastani, and Winikoff 2008) classify them as declarative or state-  
68 based where the goal is to reach specific desired situation and procedural or action-based where the  
69 goal is to execute actions. State-based goals are then sub-classified into the query, achieve and  
70 maintain goals, and action-based goals are sub-classified into perform goal. RL is already able to  
71 solve some problems where some of these kinds of goals are present. For example, well-known  
72 benchmark problems such as the cart-pole problem are maintenance goals, while others such as maze  
73 navigation are achievement goals. Likewise, problems solved with positive reward have typically  
74 approach goal properties, while problems solved from negative reward have avoidance goal  
75 properties. The idea of generating reward signals for generic forms of these goals thus seems  
76 promising. Based on this logic we propose a domain-independent reward function for each of the  
77 goal types. This approach can be applied to the goal irrespective of its origin, i.e., whether the goal is  
78 intrinsic, extrinsic or of a social origin. In this paper though, we use the output of an existing goal  
79 generation module for a mobile robot (Merrick, Siddique, and Rano 2016) to validate the proposed  
80 reward functions. We show how the intrinsic reward functions bridge the gap between goal  
81 generation and learning by providing a task-independent reward. We further demonstrate how these  
82 primitive reward functions based on the goal types can be combined to form compound reward

83 functions that can be used to learn more complex behaviors in agents. Thus, the contributions of this  
 84 paper are: 1) A proposal for task-independent intrinsic reward functions for maintenance, approach,  
 85 avoidance and achievement goal types; 2) Metrics for the measurement of the performance of these  
 86 reward functions with respect to how effectively solutions to them can be learned; and 3) A  
 87 demonstration of how these primitive reward functions can be combined to motivate learning of more  
 88 complex behaviors.

89 The remainder of the paper is organized as follows. In Section 2, we present a background on the  
 90 design of reward functions and the solutions for task-independent reward functions found in the  
 91 literature. In Section 3, we detail the proposed reward functions based on the goal types, and the  
 92 metrics we use to measure the agents' performance using those reward functions. In Section 4, we  
 93 detail experiments to examine the performance of reward functions for maintenance, approach,  
 94 avoidance, and achievement goal types on a mobile 'e-puck' robot. In Section 5, we demonstrate  
 95 complex behaviors learned from compound reward functions constructed from the autonomously  
 96 generated primitive functions for each goal type. Finally, in Section 6, we provide concluding  
 97 remarks and discuss directions for future work.

## 98 2 Background and Related Work

99 In RL, an agent perceives the state of its environment with its sensors and takes action to change that  
 100 state. The environment may comprise variables such as the robot's position, velocity, sensor values,  
 101 etc. These parameters collectively form the state of the agent. With every action that the agent  
 102 executes in the environment, it moves to a new state. The state of the agent at time  $t$  can be expressed  
 103 as:

$$104 \quad S_t = [s_t^1, s_t^2, s_t^3, \dots, s_t^n]$$

105 where each attribute  $s_t^i$  is typically a numerical value describing some internal or external variable of  
 106 the robot, and  $n$  is the number of attributes of the state. The agent takes an action  $A_t$  to change the  
 107 state of the environment from the finite set of  $m$  actions  $\mathcal{A}$ :

$$108 \quad \mathcal{A} = \{A^1, A^2, A^3, \dots, A^m\}$$

109 This state change is denoted by event  $E_t$ , formally denoted as:

$$110 \quad E_t = [e_t^1, e_t^2, e_t^3, \dots, e_t^n]$$

111 where an event attribute  $e_t^i = s_t^i - s_{t-1}^i$ . That is,

$$112 \quad E_t = S_t - S_{t-1} = [\Delta(s_t^1 - s_{t-1}^1), \Delta(s_t^2 - s_{t-1}^2), \dots, \Delta(s_t^n - s_{t-1}^n)]$$

113 Thus, an event, which is a vector of difference variables, models the transition between the states. An  
 114 action can cause a number of different transitions, and an event is used to represent those transitions.  
 115 Since this representation does not make any task-specific assumption about the values of the event  
 116 attributes, it can be used to represent the transition in a task-independent manner (Merrick 2007).

117 Finally, the experience of the agent includes the states  $S_t$  it has encountered, the events  $E_t$  that have  
 118 occurred and the actions  $A_t$  that it has performed. Thus, the experience  $X$  is a trajectory denoted as the  
 119 following, and it provides the data from which the goals can be constructed.

$$X = \{S_0, A_0, S_1, E_1, A_1, S_2, E_2, A_2, S_3, E_3, \dots\}$$

## 121 2.1 Design of Reward Functions

122 In RL, the reward is used to direct the learning process. A simple example of a reward function is a  
123 pre-defined value assignment for known states or transitions. For example:

$$r(S_t) = \begin{cases} 1 & \text{if a particular state } S_t \text{ is reached} \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

124 A more specific, task-dependent example can be seen from the canonical cart-pole domain in which a  
125 pole is attached to a cart that moves along a frictionless track. The aim of the agent is to maintain the  
126 pole balanced on the cart by moving the cart to the right or left. The reward, in this case, depends on  
127 the attributes specific to the task:

$$r(S_t) = -c2 * (G^1 - s_t^1)^2 - c3 * (G^2 - s_t^2)^2 \quad (2)$$

128 where  $s_t^1$  is the position of the cart and  $s_t^2$  is the angle of the pole with respect to the cart,  $G$  (with  
129 attributes  $G^1$  – desired position and  $G^2$  – desired angle) is the goal state, and  $c2$  and  $c3$  are constants.

130 For an even more complex task like ball paddling, where a table-tennis ball is attached to a paddle by  
131 an elastic string with the goal to bounce the ball above the paddle, it is quite difficult to design a  
132 reward function. Should the agent be rewarded for bouncing the ball a maximum number of times?  
133 Should the agent be rewarded for keeping the ball above the paddle? As detailed in (Amodei et al.  
134 2016), the agent might find ways to ‘hack the reward’ resulting in unpredictable or unexpected  
135 behavior.

136 For some complex domains, it is only feasible to design ‘sparse reward signals’ which assign non-  
137 zero reward in only a small proportion of circumstances. This makes learning difficult as the agent  
138 gets very little information about what actions resulted in the correct solution. Proposed alternatives  
139 for such environments include ‘hallucinating’ positive rewards (Andrychowicz et al. 2017) or  
140 bootstrap with self-supervised learning to build a good world model. Also, imitation learning and  
141 inverse RL have shown reward functions can be implicitly defined by human demonstrations, so they  
142 do not allow a fully autonomous development of the agent.

143 ‘Reward engineering’ is another area that has attracted the attention of the RL community, which is  
144 concerned with the principles of constructing reward signals that enable efficient learning (Dewey  
145 2014). Dewey (2014) concluded that as artificial intelligence becomes more general and autonomous,  
146 the design of reward mechanisms that result in desired behaviors are becoming more complex. Early  
147 artificial intelligence research tended to ignore reward design altogether and focused on the problem  
148 of efficient learning of an arbitrary given goal. However, it is now acknowledged that reward design  
149 can enable or limit autonomy, and there is a need for reward functions that can motivate more open-  
150 ended learning beyond a single, fixed task. The following sub-sections review work that focus in this  
151 area.

## 152 2.2 Intrinsic Motivation

153 Reward modeled as intrinsic motivation is an example of an engineered reward leading to open-  
154 ended learning (Baldassarre and Mirolli 2013). It may be computed online as a function of

155 experienced states, actions or events and is independent of *a priori* knowledge of task-specific factors  
156 that will be present in the environment. The signal may serve to drive acquisition of knowledge or a  
157 skill that is not immediately useful but could be useful later on (Singh, Barto, and Chentanez 2004).  
158 This signal may be generated by an agent because a task is inherently ‘interesting’, leading to further  
159 exploration of its environment, manipulation/play or learning of the skill.

160 Intrinsic motivation can be used to model reward that can lead to the emergence of task-oriented  
161 performance, without making strong assumptions about which specific tasks will be learned prior to  
162 the interaction with the environment. The motivation signal may be used in addition to a task-specific  
163 reward signal, aggregated based on a predefined formula, to achieve more adaptive and multitask  
164 learning. It can also be used in the absence of a task-specific reward signal to reduce the handcrafting  
165 and tuning of the task-specific reward thus moving a step closer to creating a true task independent  
166 learner (Merrick and Maher 2009). Oudeyer and Kaplan (Oudeyer and Kaplan 2007) proposed the  
167 following categories for a computational model of motivation: knowledge-based, and competence-  
168 based. In knowledge-based motivation, the motivation signal is based on an internal prediction error  
169 between the agent’s prediction of what is supposed to happen and what actually happens when the  
170 agent executes a particular action. In competence-based motivation, the motivation signal is  
171 generated based on the appropriate level of learning challenge. This competency motivation depends  
172 on the task or the goal to accomplish. The activity at a correct level of learnability given the agent’s  
173 current level of mastery of that skill generates maximum motivation signal. Barto et al. (Barto,  
174 Mirolli, and Baldassarre 2013) further differentiated between surprise (prediction error) and novelty  
175 based motivation. Novelty motivation signal is computed based on the experience of an event that  
176 was not experienced before (Neto and Nehmzow 2004) ([Nehmzow et al. 2013](#)).

### 177 **2.3 Intrinsically Motivated Reinforcement Learning**

178 Frameworks that combine intrinsic motivation with RL are capable of autonomous learning, and they  
179 are commonly termed intrinsically motivated reinforcement learning frameworks. Singh et al. (Singh,  
180 Barto, and Chentanez 2004), and Oudeyer et al. (Oudeyer, Kaplan, and Hafner 2007) state that  
181 intrinsic motivation is essential to create machines capable of lifelong learning in a task-independent  
182 manner as it favors the development of competence and reduces reliance on externally directed goals  
183 driving learning. When intrinsic motivation is combined with RL, it creates a mechanism whereby  
184 the system designer is no longer required to program a task-specific reward (Singh, Barto, and  
185 Chentanez 2004). An intrinsically motivated reinforcement learning agent can autonomously select a  
186 task to learn and interact with the environment to learn the task. It results in the development of an  
187 autonomous entity capable of resolving a wide variety of activities, as compared to an agent capable  
188 of resolving only a specific activity for which a task-specific reward is provided.

189 Like in RL, in an intrinsically motivated reinforcement learning framework, the agent senses the  
190 states, takes actions and receives an external reward from the environment, however as an additional  
191 element, the agent internally generates a motivation signal that forms the basis for its actions. This  
192 internal signal is independent of task-specific factors in the environment. Incorporating intrinsic  
193 motivation with RL enables agents to select which skills they will learn and to shift their attention to  
194 learn different skills as required (Merrick 2012). Broadly speaking, intrinsically motivated  
195 reinforcement learning introduces a meta-learning layer in which a motivation function provides the  
196 learning algorithm with a motivation signal to focus the learning (Singh, Barto, and Chentanez 2004).

### 197 **2.4 Role of Goals to Direct the Learning**

198 Where early work focused on generating reward directly from environmental stimuli, more recent  
 199 works have acknowledged the advantages of using the intermediate concept of a goal to motivate  
 200 complexity and diversity of behavior (Santucci, Baldassarre, and Mirolli 2016) (Merrick, Siddique,  
 201 and Rano 2016). It has been shown by Santucci et al. (Santucci, Baldassarre, and Mirolli 2012) that  
 202 using intrinsic motivation (generated by prediction error) directly for skill acquisition can be  
 203 problematic and a possible solution to that is to instead generate goals using the intrinsic motivation  
 204 which in turn can be used to direct the learning. Further, it has been argued by Mirolli and  
 205 Baldassarre (Mirolli and Baldassarre 2013) that a cumulative acquisition of skills requires a  
 206 hierarchical structure, in which multiple ‘expert’ sub-structures focus on acquiring different skills and  
 207 a ‘selector’ sub-structure decides which expert to select. The expert substructure can be implemented  
 208 using knowledge-based intrinsic motivation that decides what to learn (by forming goals), and the  
 209 selector sub-structure can be implemented using competence-based intrinsic motivation that can be  
 210 used to decide which skill to focus on. Goal-directed learning is also shown to be a promising  
 211 direction for learning motor skills. Rolf et al. (Rolf, Steil, and Gienger 2010) show how their system  
 212 auto-generates goals using inconsistencies during exploration to learn inverse kinematics and that the  
 213 approach can scale for a high dimension problem.

214 Recently, using goals to direct the learning has even attracted the attention of the deep learning  
 215 community. Andrychowicz et al. (Andrychowicz et al. 2017) have proposed using auto-generated  
 216 interim goals to make learning possible even when the rewards are sparse. These interim goals are  
 217 used to train the deep learning network using experience replay. It is shown that the RL agent is able  
 218 to learn to achieve the end goal even if it has never been observed during the training of the network.  
 219 Similarly, in a framework proposed by Held et al. (Held et al. 2017), they auto-generate interim  
 220 tasks/goals at an appropriate level of difficulty. This curriculum of tasks then directs the learning  
 221 enabling the agent to learn a wide set of skills without any prior knowledge of its environment.

222 Regardless of whether the goals are intrinsic, extrinsic, of social origin, whether they are created to  
 223 direct the learning or generated by an autonomous learning framework, the approach of using goal-  
 224 based reward functions detailed in the next section can be applied to them.

### 225 **3 Primitive Goal-based Motivated Reward Functions**

226 The basis of our approach in this paper is a generic view of the function in Equation (1) as follows:

$$227 \quad r(S_t) = \begin{cases} 1 & \text{if the goal is reached} \\ 1 - \epsilon & \text{otherwise} \end{cases} \quad (3)$$

228 where  $\epsilon$  is a non-negative constant. The remainder of this section defines different representations of  
 229 ‘goal’ in Equation (3) and representation of the meaning of ‘reached’.

#### 229 **3.1 Reward Function for the Maintenance Goal Type**

230 A maintenance goal monitors the environment for some desired world state and motivates the agent  
 231 to actively try to re-establish that state if the distance between the desired state and the current state  
 232 goes beyond a set limit. For a maintenance goal, an agent’s action selection should consider both  
 233 triggering conditions as well as the constraining nature of the goal (Hindriks and Van Riemsdijk  
 234 2007). The act of maintaining a goal can be never-ending thus making the process continuous or non-  
 235 episodic.



$$M_4 = \max_{j=1\dots J}(\text{length of maintenance period } j)$$

### 3.2 Reward Function for the Approach Goal Type

An approach goal represents the agent's act of attempting to get closer to the desired world state. The main difference between an approach and maintenance goal lies in the condition of fulfillment. An approach goal is fulfilled when the agent is getting closer to the desired state whereas a maintenance state is fulfilled when the desired state is maintained and not violated. An approach attempt leads to a behavior that functions to shorten the distance, either physically or psychologically between the agent and the desired outcome (Elliot 2008).

The reward function for the approach goal can be expressed as:

$$r(S_t) = \begin{cases} \sigma & \text{if } d(S_t, G) < d(S_{t-1}, G) \text{ and } d(S_t, G) > \rho \\ \varphi & \text{otherwise} \end{cases} \quad (5)$$

where  $d(\cdot)$ , the distance function is used to check the approach attempt by comparing the distance between the current state  $S_t$  and the desired goal state  $G$  with the distance between the previous state  $S_{t-1}$  and  $G$ . The second condition of the equation ensures that the distance is more than the defined distance threshold  $\rho$  so that 'reached' means an approach attempt and not "approach and achieve". Same as in Equation (4), the reward for when the goal is not reached is  $\varphi$  with  $\varphi < \sigma$  in order to incentivize the agent to find a shorter path to the goal state.

The following metrics may thus be used to evaluate this reward function for the approach goal type. Each metric is again assumed to be measured over a fixed period  $T$  of the agent's life. Since the approach and avoidance functions (detailed in section 3.3) reward the approach and the avoidance attempts irrespective of the distance between the current and the goal state, the cumulative reward for the agent is very high. In order to get a better sense of the proportion of the reward gained per trial, we use percentage in the following metrics.

- **Number of steps the goal is approached as a percentage of  $T$  ( $M_5$ ).** This metric indicates the approachability of the goal, i.e., how easy is it to approach the goal state?

$$M_5 = \frac{M_2 \times 100}{T}$$

- **Number of approach attempts as a percentage of  $T$  ( $M_6$ ).** The agent is considered to have made an approach attempt if it receives a positive reward for two or more consecutive steps, i.e., signifying that the agent attempted to approach the goal state.

$$M_6 = \frac{M_1 \times 100}{T}$$

### 3.3 Reward Function for the Avoidance Goal Type

An avoidance goal type is the opposite of the approach goal type. Avoidance is a behavior where an agent stays away or moves away from an undesirable stimulus, object or event (Elliot 2008). An avoidance goal is considered fulfilled as long as the agent is away from the state that it wants to avoid, and it increases the distance from the state that it wants to avoid. Considering those





336 labeled in a clockwise direction as *Front-Right*, *Right*, *Rear-Right*, *Rear-Left*, *Left*, and *Front-Left*.  
 337 The red lines in Figure 1(a) show the direction in which the sensors detect an obstacle. A high sensor  
 338 reading indicates that an object is close to that sensor. Figure 1(b) shows a 5×5 meter square flat  
 339 walled arena that we use for our experimentation with primitive goal-based reward functions.

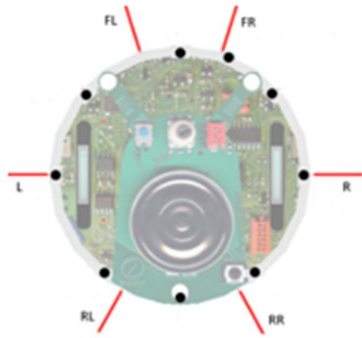


Figure 1(a): e-puck proximity sensors (shown by the red directional lines)

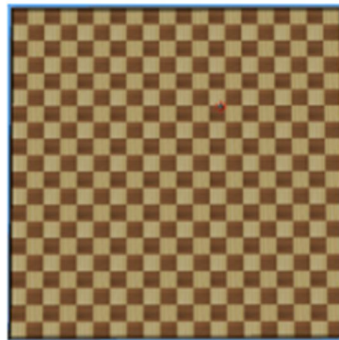


Figure 1(b): A simple walled arena

State:

$$[\omega^R \ \omega^L \ \theta \ s^L \ s^R \ s^{FL} \ s^{FR} \ s^{RL} \ s^{RR}]$$

Actions:

- 1 – Left\_Wheel\_Speed +  $\delta$
- 2 – Right\_Wheel\_Speed +  $\delta$
- 3 – Left\_Wheel\_Speed -  $\delta$
- 4 – Right\_Wheel\_Speed -  $\delta$
- 5 – No change to wheel speeds

340 The arena, the state, and the action space of the robot are the same as detailed by Merrick et al.  
 341 (Merrick, Siddique, and Rano 2016). The state of the mobile robot comprises nine parameters: left  
 342 wheel speed, right wheel speed, orientation, left sensor value, right sensor value, front-left sensor  
 343 value, front-right sensor value, rear-left sensor value and rear right sensor value, i.e., the state vector  
 344 is  $[\omega^R \ \omega^L \ \theta \ s^L \ s^R \ s^{FL} \ s^{FR} \ s^{RL} \ s^{RR}]$ .  $\omega^R$  and  $\omega^L$  are the rotational velocities of the right and the left  
 345 wheels. Their range is  $-\pi$  to  $\pi$  radians per second.  $\theta$  is the orientation angle of the mobile robot. Its  
 346 value ranges from  $-\pi$  to  $\pi$ . For our experiments, we use binary values for the proximity sensors with 0  
 347 indicating that there is no object in the proximity of the sensor, and 1 indicates that the object is near.  
 348 The rotational velocities and orientation are discretized into nine values making the state space quite  
 349 large.

350 The action space comprises five actions: 1 – increase the left wheel speed by  $\delta$ , 2- increase the right  
 351 wheel speed by  $\delta$ , 3 – decrease the left wheel speed by  $\delta$ , 4 – decrease the right wheel speed by  $\delta$  and  
 352 5 - no change to any of the wheel speeds. A fixed value of  $\pi/2$  was used as  $\delta$ .

353 In this paper, we use the goals generated for the mobile robot based experiment by Merrick et al.  
 354 (Merrick, Siddique, and Rano 2016). The main concept of the experience based goal generation  
 355 detailed in (Merrick, Siddique, and Rano 2016) is that the agent must explore its environment and  
 356 determine if the experience is novel enough to be termed a potential goal. Goal generation phase is  
 357 divided into two stages: experience gathering stage and the goal clustering stage. In the experience  
 358 gathering stage, the mobile robot moves around randomly in its environment. The states experienced  
 359 by the robot are recorded. These recorded states form an input to the goal clustering stage which uses  
 360 simplified adaptive resonance theory (SART) network (Baraldi 1998). SART is a neural network  
 361 based clustering technique. It is capable of handling a continuous stream of data thus solving the  
 362 stability-plasticity dilemma. The network layer takes a vector input and identifies its best match in  
 363 the network. Initially, the network starts with no clusters. As the data is read, its similarity is checked  
 364 with any existing clusters. If there is close enough match, it is clustered together else a new cluster is  
 365 created. As the clusters are created, they are connected to the input nodes (i.e., the recorded  
 366 experience). The number of clusters created will depend on the vigilance parameter of the SART

367 network. Higher vigilance produces many fine-grained clusters whereas a low vigilance parameter  
 368 produces a coarser level of clusters. The goals generated by this phase form input for the goal  
 369 learning phase.

370 In the learning phase, the robot learns the skills to accomplish the goals. For the goal learning, we use  
 371 an RL algorithm called Dyna-Q. Dyna-Q (Sutton and Barto 1998) is a combination of Dyna  
 372 architecture with RL’s Q Learning algorithm. With Dyna-Q, the Q-Learning is augmented with  
 373 model learning, thus combining both model-based and model-free learning. The RL agent improves  
 374 its Q value function using both the real experiences with its environment and imaginary experiences  
 375 (also called planning process) generated by the model of the environment. During the planning  
 376 process, that is typically run several times for every real interaction with the environment; the  
 377 algorithm randomly selects the samples from the model (continuously updated using the real  
 378 experiences) and updates the Q value function. This reduces the number of interactions required with  
 379 the environment which are typically expensive, and especially for the robotic applications. The model  
 380 of the environment for our experiments keeps track of the state  $s'$  that the mobile robot lands in  
 381 when it takes a particular action  $a$  in the current state  $s$ . The model also keeps track of the reward that  
 382 the robot receives during that transition. The state transitions for our experiments are deterministic in  
 383 nature, i.e., when the robot takes action  $a$  in state  $s$ , it will always land in a state  $s'$ . The number of  
 384 iterations for model learning can be varied as required. We set this parameter to 25, i.e., the algorithm  
 385 will attempt 25 actions for model learning (using imaginary experiences) before attempting one  
 386 action with the real environment.

#### 387 4.1 Maintenance Goal Learning Results

388 ~~Table 1 shows the results of the experiments for the maintenance goals. The goal ID, goal attribute~~  
 389 ~~and the meaning of the goal, are the maintenance goals generated by the SART based clustering as~~  
 390 ~~detailed by Merrick et al. (Merrick, Siddique, and Rano 2016) used SART based clustering to~~  
 391 ~~generate two sets of goals, namely, maintenance and achievement goals. Table 1 lists the set of~~  
 392 ~~maintenance goals described by the ID, goal attributes and the meaning of the goal as detailed by~~  
 393 ~~Merrick et al. (Merrick, Siddique, and Rano 2016). These goals are the actual states experienced by~~  
 394 ~~the mobile robot. This same set of goals are used in section 4.2 and 4.3 treated as approach and~~  
 395 ~~avoidance type respectively. Table 4 in section 4.4, lists the set of achievement goals generated by~~  
 396 ~~Merrick et al. (Merrick, Siddique, and Rano 2016).~~

397 Table 1 also shows the results of the experiments for these goals treated as maintenance goals. The  
 398 columns  $M_1$ ,  $M_2$ ,  $M_3$ , and  $M_4$  are the metrics detailed in section 3.1. The goals are states experienced  
 399 by the mobile robot treated as maintenance goal for these experiments, i.e., the aim of the robot is to  
 400 maintain these goal states. The e-puck mobile robot simulation was run for ten trials each of 25,000  
 401 steps for each of the 12 goals. Results were averaged over ten trials, and the standard deviation is also  
 402 shown in the table. Values of the parameters of Equation (4) were as follows:  $\rho$  was 0.9,  $\sigma$  was 1,  $\phi$   
 403 was -1 and  $d$  was the Euclidian distance. The RL exploration parameter epsilon was set to 0.15, and  
 404 the decay schedule was linear. When a trial ended, the end position and orientation of the e-puck  
 405 mobile robot became the start position and orientation for the next trial. However, the RL Q table  
 406 was reset after each trial, so no learning was carried forward between the trials.

407 Table 1: Experiments and results for maintenance goals

ID	Goal Attributes	Meaning of the Goal	$M_1$	$M_2$	$M_3$	$M_4$	Is Goal Valid?
G <sup>1</sup>	(2.5, 2.5, 1.8, 0, 0, 0, 0, 0)	Move forward at high speed	37 ±8	493 ±91	14 ±4	154 ±7	Yes

$G^2$	(0.4, 0.4, 1.2, 0, 0, 0, 0, 0, 0)	Move forward at low speed	121 ±25	568 ±124	4 ±1	88 ±0	Yes
$G^3$	(-2.4, -2.4, 1.4, 0, 0, 0, 0, 0, 0)	Move backward at high speed	88 ±8	888 ±179	10 ±2	188 ±9	Yes
$G^4$	(-0.4, -0.4, -1.3, 0, 0, 0, 0, 0, 0)	Move backward at low speed	192 ±28	866 ±110	4 ±0	71 ±0	Yes
$G^5$	(0.0, 0.0, -2.8, 0, 1, 0, 0, 0, 0)	Stop for obstacle in front	1 ±1	3 ±3	1 ±0	5 ±0	Yes
$G^6$	(-0.4, -0.4, 2.9, 0, 0, 0, 0, 0, 0)	Move backward at low speed	142 ±24	601 ±106	4 ±0	37 ±1	Yes
$G^7$	(-0.8, -0.8, 1.6, 0, 0, 0, 0, 0, 0)	Move backward at moderate speed	157 ±26	848 ±127	5 ±0	53 ±2	Yes
$G^8$	(0.2, 0.0, 2.4, 1, 0, 0, 0, 0, 1)	Stop for obstacle behind	0 ±0	0 ±0	0 ±0	0 ±0	Yes
$G^9$	(0.0, -0.3, 2.1, 1, 0, 0, 0, 1, 0)	Stop <del>in free space for</del> obstacle at left and back	0 ±0	0 ±0	0 ±0	2 ±0	Yes
$G^{10}$	(-1.9, -1.9, -2.2, 0, 0, 0, 0, 0, 0)	Move backward at moderate speed	162 ±23	763 ±105	4 ±0	52 ±2	Yes
$G^{11}$	(0.0, 0.0, 3.0, 0, 1, 1, 0, 0, 0)	Stop for obstacle in front	0 ±0	0 ±0	0 ±0	0 ±0	No
$G^{12}$	(1.2, 1.2, -2.7, 0, 0, 0, 0, 0, 0)	Move forward at moderate speed	100 ±18	427 ±85	4 ±0	36 ±1	Yes

408 Once the robot reaches the goal state, it maintains it until it comes across adverse conditions, i.e., for  
409  $G^l$  (move forward at high speed), once the goal state is reached, the robot will maintain that state  
410 while it is in the open space. However, once it reaches a wall, it is not able to maintain the state. We  
411 consider that the robot has learnt to attain the goal if the robot is able to reach the goal state over and  
412 over again and remain in that state for two time-steps or more. This is indicated by the column for  
413  $M_1$ . This measure is high for  $G^1$ ,  $G^2$ ,  $G^3$ ,  $G^4$ ,  $G^6$ ,  $G^7$ ,  $G^{10}$  and  $G^{12}$  indicating that the robot is able to  
414 maintain those goals. However, that measure is very low for goal  $G^5$  and zero for  $G^8$  which shows  
415 that the robot is not able to learn to maintain those goal states. This is due to the lack of opportunity,  
416 i.e., the robot has to be in a specific situation to be able to learn to maintain those goals. Those goals  
417 require the robot to be close to a wall, the likelihood of which is small because of the size of the  
418 arena.

419 ~~The  $M_1$  measure is zero for goal  $G^9$ , which is a valid goal, although the column ‘meaning of the goal’~~  
420 ~~does not seem correct. Meaning should be “Stop for obstacle at left and back”. The measure  $M_1$  for~~  
421 ~~goal  $G^9$ , which is a valid goal, is zero.~~ The mobile robot was not able to achieve that goal because of  
422 the lack of opportunity. The required situation to learn that goal would be that the robot should find  
423 itself in the bottom left corner at a particular orientation. The measure  $M_1$  is zero for  $G^{11}$  as well. The  
424 reason for that is because goal  $G^{11}$  is an unreasonable goal. According to that state, the wall is close  
425 to the Right and Front Left sensors but not Front Right. It is hard to imagine a position of the mobile  
426 robot that represents ~~that such~~ state. The goals created by SART are the cluster centers. It appears  
427 that this is an example of the clustering algorithm creating a hybrid, unreasonable goal which could  
428 be either because the granularity of the clusters is coarser than it should be, resulting in the cluster  
429 centroid not being a correct representative of the cluster or that invalid states experienced by the  
430 robot due to noise, ~~resulted in an invalid event ( $e_t = s_t - s_{t-1}$ ).~~ The column ‘Is Goal Valid?’ is  
431 marked ‘No’ in this case.

432 Figure 2(a) shows a sample trajectory of the mobile robot for  $G^l$ . ~~The trajectory is a two-dimensional~~  
433 ~~plot of the path followed by the mobile robot in the arena during the trial.~~ The goal is attained by  
434 maintaining a high speed at a particular orientation. The robot receives a positive reward for the time  
435 steps that it maintains the goal. It is only possible for the robot to attain  $G^l$  when it is in the open area  
436 of the arena. When it reaches the wall, it is no longer able to maintain goal  $G^l$ . The robot has to learn  
437 to turn around and attain the goal again. This is evident in figure 2(a) that shows multiple straight  
438 stretches where the robot attains  $G^l$ , reaches the wall, tries to turn around and attains the goal again.

439 Figure 2(b) shows the trajectory of the mobile robot for  $G^3$  (move backward at high speed) ~~and~~  
 440 ~~Figure 3(a) shows the trajectory for goal  $G^{12}$  (move forward at moderate speed).~~

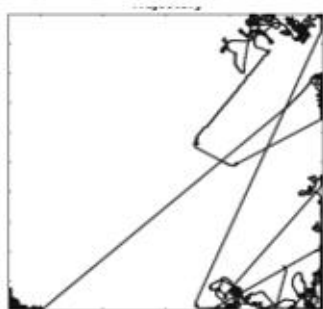


Figure 2(a): Mobile robot trajectory for  $G^1$

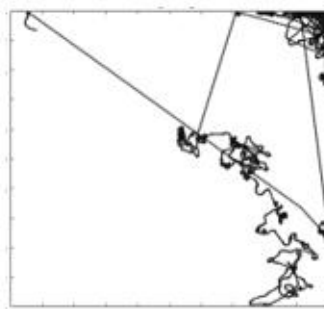


Figure 2(b): Mobile robot trajectory for  $G^3$

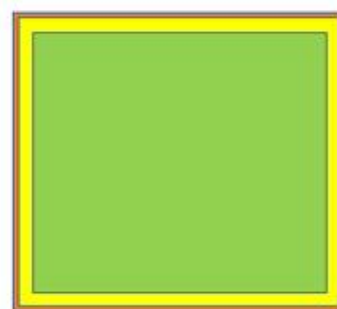


Figure 2(c): Likelihood of the reward for  $G^1$ ,  $G^3$  and  $G^{12}$

441 For goals  $G^1$ , ~~and~~  $G^3$ , ~~and~~  $G^{12}$  the robot is only able to attain the goals when it is in the open area of  
 442 the arena. Figure 2(c) shows the likelihood diagram with the wall shown in orange. In the open area of  
 443 the arena shown in green, the robot is more likely to attain the goal, i.e., to receive a positive  
 444 reward. In the area close to the wall (shown in yellow) the likelihood reduces. The probability of the  
 445 mobile robot to be in the green zone can be calculated as follows for the environment with the size of  
 446 the board  $5\text{m} \times 5\text{m}$  and sensor range of e-puck  $0.06\text{m}$ . If we were to discretize the environment into  
 447 squares of  $0.06\text{m}$ , then there would be  $83 \times 83$ , i.e., 6889 squares in the grid. Green zone for  $G^1$ , ~~and~~  
 448  $G^3$  ~~and~~  $G^{12}$  will consist of  $81 \times 81$ , i.e., 6561 squares. If we were to randomly select a square in the  
 449 green zone, the probability would be  $(81 \times 81) / (83 \times 83) = 95.23\%$ . The orientation and wheel speeds are  
 450 divided into nine buckets each. Hence the probability of the robot to be in a particular square with  
 451 particular wheel speed and orientation will be  $(81 \times 81) / (83 \times 83 \times 9 \times 9 \times 9) = 0.13\%$ .

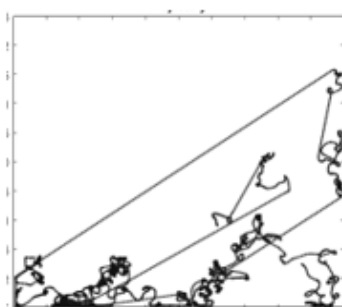


Figure 3(a): Trajectory for  $G^{12}$

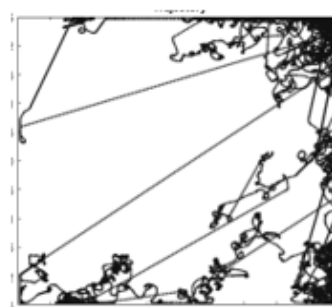
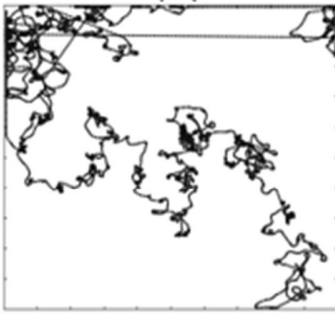
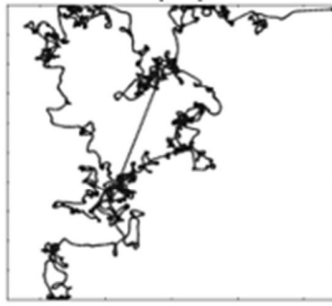
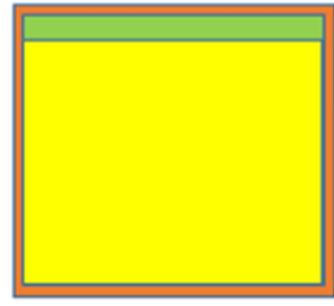


Figure 3(b): Simulation for a  $G^{12}$  run for 100,000 steps

452 ~~Figures 3(a) show the trajectories of goals  $G^{12}$  (move forward at moderate speed). The robot can~~  
 453 ~~learn to attain the goal.~~ For  $G^{12}$  we let the simulation for one of the trials continue for 100,000 steps,  
 454 the trajectory of which is shown in Figure 3(b). The straight-line trajectory shows that the robot is  
 455 maintaining the goal of moving forward at a moderate speed, i.e., it is in the region of opportunity  
 456 (Figure 3(c)). When the robot reaches the wall, it experiences states that it may not have  
 457 experienced in the past. However, it eventually learns to attain the goal of moving forward at a  
 458 moderate speed.

Figure 4(a): Trajectory for  $G^5$ Figure 4(b): Simulation for a  $G^8$ Figure 4(c): Likelihood of the reward for  $G^5$  and  $G^8$ 

459 Figure 4(a) and 4(b) shows the trajectory for goal  $G^5$  (stop for an obstacle in front) and  $G^8$  (stop for  
 460 obstacle behind) respectively. The robot does not learn to attain these goals. The obstacles in the  
 461 arena are the four walls hence the likelihood of the reward are the areas closer to the wall.  
 462 Considering the orientation for goals  $G^5$  and  $G^8$ , the mobile robot has to be beside the top wall as  
 463 shown in green in figure 4(c). The probability of the mobile robot to be in a particular square with the  
 464 orientation required for  $G^5$  or  $G^8$  is  $(81)/(83 \times 83 \times 9 \times 9 \times 9) = 0.002\%$ . This lack of opportunity is the  
 465 reason why the robot does not learn  $G^5$  and  $G^8$  goals. In order to confirm this hypothesis, we  
 466 continued the experiments with these two goals with the reduced arena size. The size of the arena  
 467 was reduced to  $0.25\text{m} \times 0.25\text{m}$  to increase the opportunity for the mobile robot to be near a wall. In  
 468 that arena, the probability of the mobile robot to find itself in the required situation is increased by  
 469 the factor of 400 ( $20 \times 20$ ) to  $0.65\%$ , thus increasing its ability to attain  $G^5$  and  $G^8$  goals. ~~In this  
 470 smaller arena, the mobile robot learnt to attain  $G^5$  and  $G^8$  goals.~~

## 471 4.2 Approach Goal Results

472 Table 2 shows the results of the experiments for the approach goals. The ~~twelve goals and their~~  
 473 ~~corresponding~~ goal IDs, goal attributes and the meaning of the goal, are the same as the ~~maintenance~~  
 474 goals detailed in Table 1. The goals for these set of experiments will be treated as approach goals,  
 475 i.e., the aim of the robot is to approach those goal states. ~~Values of the parameters of Equation (5)~~  
 476 ~~and the method in which experiments were conducted for the approach goals were the same as~~  
 477 ~~detailed in section 4.1. Similarly, in the experiments detailed in section 4.1, the e-puck mobile robot~~  
 478 ~~simulation was run for ten trials for each goal with 25,000 steps in each trial. Values of the~~  
 479 ~~parameters of Equation (5) were as follows:  $\rho$  was 0.9,  $\sigma$  was 1,  $\phi$  was -1 and  $d$  was the Euclidian~~  
 480 ~~distance. The RL exploration parameter epsilon was set to 0.15 with a linear decay schedule, and the~~  
 481 ~~Q table was reset after each trial thus there was no learning carried forward between the trials.~~

482 Table 2: Experiments and results for approach goals

ID	Goal Attributes	Meaning of Goal	$M_5$	$M_6$
$G^1$	(2.5, 2.5, 1.8, 0, 0, 0, 0, 0)	Move forward at high speed	32.49% $\pm$ 0.64	7.56% $\pm$ 0.16
$G^2$	(0.4, 0.4, 1.2, 0, 0, 0, 0, 0)	Move forward at low speed	34.66% $\pm$ 0.62	8.00% $\pm$ 0.21
$G^3$	(-2.4, -2.4, 1.4, 0, 0, 0, 0, 0)	Move backward at high speed	36.58% $\pm$ 0.41	8.39% $\pm$ 0.14
$G^4$	(-0.4, -0.4, -1.3, 0, 0, 0, 0, 0)	Move backward at low speed	35.88% $\pm$ 0.43	8.52% $\pm$ 0.11
$G^5$	(0.0, 0.0, -2.8, 0.1, 0, 0, 0, 0)	Stop for obstacle in front	37.27% $\pm$ 0.88	8.84% $\pm$ 0.34
$G^6$	(-0.4, -0.4, 2.9, 0, 0, 0, 0, 0)	Move backward at low speed	37.25% $\pm$ 0.38	8.74% $\pm$ 0.19
$G^7$	(-0.8, -0.8, 1.6, 0, 0, 0, 0, 0)	Move backward at moderate speed	36.77% $\pm$ 0.57	8.76% $\pm$ 0.22
$G^8$	(0.2, 0.0, 2.4, 1, 0, 0, 0, 1)	Stop for obstacle behind	37.15% $\pm$ 0.64	8.73% $\pm$ 0.22

$G^9$	(0.0, -0.3, 2.1, 1, 0, 0, 0, 1, 0)	Stop in free space	36.71% ±0.98	8.60% ±0.26
$G^{10}$	(-1.9, -1.9, -2.2, 0, 0, 0, 0, 0, 0)	Move backward at moderate speed	36.12% ±0.60	8.24% ±0.23
$G^{11}$	(0.0, 0.0, 3.0, 0, 1, 1, 0, 0, 0)	Stop for obstacle in front	36.89% ±0.86	8.74% ±0.26
$G^{12}$	(-1.2, 1.2, -2.7, 0, 0, 0, 0, 0, 0)	Move forward at moderate speed	33.58% ±0.58	7.40% ±0.17

483 The design of the reward function for the approach goal type is such that it rewards an approach  
484 attempt. Hence if the agent is getting closer to the goal, it receives a positive reward. Goals, when  
485 treated as approach goals, are relatively straightforward to attain as seen in the  $M_5$  column in Table 2  
486 (average number of steps positive reward received as a percentage). In the case of the goal  $G^l$ , for  
487 instance, the agent receives a positive reward for 32.49% of the time steps. This is because the  
488 attempt to approach the goal is rewarded irrespective of the distance between the current state and the  
489 goal state. Results also show that all the goals, when treated as approach type, are attainable (even the  
490 invalid goals) indicating that it is possible to approach the goal states of each of the 12 goals.

### 491 4.3 Avoidance Goal Results

492 Table 3 shows the results of the experiments for the avoidance goals. The twelve goals and their  
493 corresponding goal IDs, goal attributes and the meaning of the goal, are the same as the maintenance  
494 goals detailed in Table 1. The goal states for these experiments are treated as avoidance goals, i.e.,  
495 the aim of the robot is to avoid those goal states. Values of the parameters of Equation (6) and the  
496 method in which experiments were conducted for the avoidance goals were the same as detailed in  
497 section 4.1. Same as the experiments in section 4.1 and 4.2, the e-puck mobile robot simulation was  
498 run for ten trials for each goal with 25,000 steps in each trial. Values of the parameters of Equation  
499 (6) were as follows:  $\rho$  was 0.9,  $\sigma$  was 1,  $\phi$  was -1 and  $d$  was the Euclidian distance. The RL  
500 exploration parameter epsilon was set to 0.15 with a linear decay schedule. Also, the Q table was  
501 reset after each trial thus there was no learning carried forward between the trials.

502

Table 3: Experiments and results for avoidance goals

ID	Goal Attributes	Meaning of Goal	$M_5$	$M_6$	$M_7$
$G^1$	(2.5, 2.5, 1.8, 0, 0, 0, 0, 0, 0)	Move forward at high speed	36.67% ±0.32	8.63% ±0.14	45
$G^2$	(0.4, 0.4, 1.2, 0, 0, 0, 0, 0, 0)	Move forward at low speed	34.88% ±0.67	8.05% ±0.25	14
$G^3$	(-2.4, -2.4, 1.4, 0, 0, 0, 0, 0, 0)	Move backward at high speed	32.61% ±0.41	7.53% ±0.16	12
$G^4$	(-0.4, -0.4, -1.3, 0, 0, 0, 0, 0, 0)	Move backward at low speed	33.16% ±0.53	7.62% ±0.14	12
$G^5$	(0.0, 0.0, -2.8, 0, 1, 0, 0, 0, 0)	Stop for obstacle in front	35.60% ±1.01	8.21% ±0.30	1
$G^6$	(-0.4, -0.4, 2.9, 0, 0, 0, 0, 0, 0)	Move backward at low speed	34.22% ±0.94	7.95% ±0.25	16
$G^7$	(-0.8, -0.8, 1.6, 0, 0, 0, 0, 0, 0)	Move backward at moderate speed	33.46% ±0.55	7.75% ±0.22	13
$G^8$	(0.2, 0.0, 2.4, 1, 0, 0, 0, 0, 1)	Stop for obstacle behind	34.90% ±0.84	8.11% ±0.18	0
$G^9$	(0.0, -0.3, 2.1, 1, 0, 0, 0, 1, 0)	Stop in free space	35.54% ±0.64	8.31% ±0.17	0
$G^{10}$	(-1.9, -1.9, -2.2, 0, 0, 0, 0, 0, 0)	Move backward at moderate speed	32.74% ±0.75	7.52% ±0.16	6
$G^{11}$	(0.0, 0.0, 3.0, 0, 1, 1, 0, 0, 0)	Stop for obstacle in front	35.46% ±0.97	8.26% ±0.33	0
$G^{12}$	(-1.2, 1.2, -2.7, 0, 0, 0, 0, 0, 0)	Move forward at moderate speed	37.00% ±0.77	8.56% ±0.20	7

503 The reward function for the avoidance goal type rewards the attempt to avoid the goal, i.e., the agent  
504 is moving away from the goal state. As it can be seen in the table, the goals, when treated as  
505 avoidance goals, are relatively easy to attain. This is because the attempt to avoid the desired goal  
506 state is rewarded irrespective of the distance between the current state and the goal state. Based on  
507 the  $M_7$  column (average number of times the goal state was not avoided), it can be concluded-said  
508 that even the goals that are difficult to attain due to lack of opportunity, when treated as maintenance  
509 goals (for example,  $G^5$ ,  $G^8$ , and  $G^9$ ), are easier to avoid when treated as avoidance goals.

### 510 4.4 Achievement Goals

511 Table 4 lists the set of achievement goals generated by Merrick et al. (Merrick, Siddique, and Rano  
512 2016). Table 4 shows the results of the experiments (with 95% confidence interval) for the  
513 achievement goals. Here too ~~†~~The goal ID, goal attributes, and the meaning of the goal are the output  
514 of the SART based clustering as detailed by Merrick et al. (Merrick, Siddique, and Rano 2016). The  
515 goal states ~~are~~ is not the actual state experienced by the mobile robot but is an ~~the~~ events as described  
516 by  $e_t^i = s_t^i - s_{t-1}^i$ . Thus, for an achievement goal type, the aim of the mobile robot is to learn to  
517 achieve the transition described by that event, for example, to learn to achieve goal  $aG^5$  listed in  
518 Table 4, which is to increase speed of both wheels, the robot must learn to increase its right wheel  
519 speed by 0.9 and left wheel speed by 0.6 in a single transition of state. The goal is considered  
520 achieved when the transition  $e_t^i$  is reached regardless of what the state  $s_{t-1}^i$  is.

521 Table 4 also shows the results of the experiments (with a 95% confidence interval) for the  
522 achievement goals. The e-puck mobile robot simulation was run for 10 trials for each goal with  
523 25,000 steps in each trial. Parameters of Equation (7) were same as in the above experiments, i.e.,  $\rho$   
524 was set to 0.9,  $\sigma$  set to 1,  $\phi$  set to -1 and  $d$  was the Euclidian distance. Also, the RL exploration  
525 parameter epsilon was set to 0.15 with a linear decay schedule. For achievement goals too, when a  
526 trial was finished the next trial started at the same position and orientation of the e-puck mobile robot  
527 at which the previous trial ended. The Q table, however, was reset after each trial thus there was no  
528 learning carried forward between the trials.

529

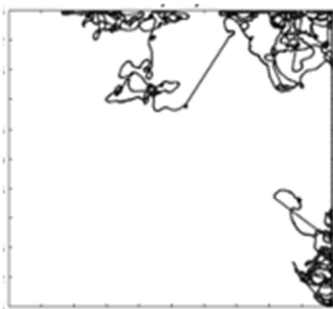
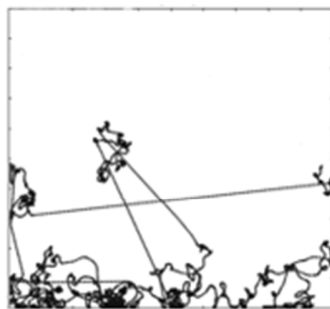
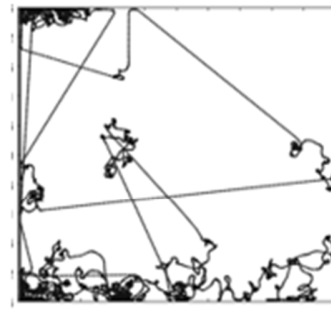
Table 4: Experiments and results for achievement goals

ID	Goal Attributes	Meaning of Goal	$M_2$	Is Goal Valid?
$aG^1$	(0.0, 0.0, 0.0, 0, 0, 0, 0, 0)	Achieve no change	25000 $\pm$ 0	Yes
$aG^2$	(0.0, 0.0, 0.0, 0, 0, 1, 0, 0)	Detect obstacle in front	43 $\pm$ 21	Yes
$aG^3$	(-0.1, 0.0, 0.0, 0, 0, -1, 0, 0)	Turn left to avoid obstacle on the right	0 $\pm$ 0	Yes
$aG^4$	(-0.6, 0.0, -0.1, 0, 0, 0, -1, 0)	Turn left to avoid obstacle on the right	0 $\pm$ 0	Yes
$aG^5$	(0.9, 0.6, 0.0, 0, 0, 0, 0, 0)	Increase speed of both wheels	6521 $\pm$ 268	Yes
$aG^6$	(-0.1, 0.1, 0.1, 0, 0, 0, 0, 0)	Turn left	0 $\pm$ 0	Yes
$aG^7$	(0.1, 0.0, -0.1, 0, 0, 0, 0, 0)	Turn right	0 $\pm$ 0	Yes
$aG^8$	(0.1, -0.4, 0.0, 0, 0, 0, 0, -1)	Turn right to avoid obstacle behind	54 $\pm$ 17	Yes
$aG^9$	(-0.3, 0.4, -0.3, 0, 0, -1, -1, 0)	Turn left to avoid obstacle on the right	0 $\pm$ 0	Yes
$aG^{10}$	(0.0, 0.5, 0.2, 0, 0, 1, 0, 0)	Turn left to detect obstacle on the right	29 $\pm$ 16	Yes
$aG^{11}$	(-0.6, -0.8, -0.2, 0, 0, -1, 0, 0)	Turn right to avoid obstacle	10 $\pm$ 4	Yes
$aG^{12}$	(0.0, 0.7, 0.3, 0, -1, 1, 0, 0)	Turn left to sense obstacle on right	0 $\pm$ 0	No
$aG^{13}$	(0.2, -0.8, -0.4, 0, 0, 0, 0, 1)	Turn right to sense obstacle on left	12 $\pm$ 4	Yes
$aG^{14}$	(0.0, 0.6, 0.1, 0, 0, 0, 0, 1)	Turn to detect obstacle behind	0 $\pm$ 0	Yes
$aG^{15}$	(0.0, -0.1, 0.0, 0, 1, 1, 0, 0)	Turn right to sense obstacle in front	0 $\pm$ 0	Yes
$aG^{16}$	(1.0, 0.5, 0.1, 0, 1, 0, 0, 0)	Turn right to sense obstacle on left	0 $\pm$ 0	NoYes
$aG^{17}$	(0.7, 0.9, 0.3, 0.0, -1, 0, 0, 0)	Turn left to sense obstacle on left	18 $\pm$ 3	Yes
$aG^{18}$	(1.2, 0.5, -0.1, 0, -1, 0, 0, 0)	Turn to avoid obstacle on left	0 $\pm$ 0	No
$aG^{19}$	(0.2, 2.7, -0.2, 0, -1, 0, 0, 0)	Turn to avoid obstacle on left	0 $\pm$ 0	No
$aG^{20}$	(-1.7, -0.5, 0.1, 0, 1, 0, 0, 0)	Turn to detect obstacle on right	0 $\pm$ 0	No
$aG^{21}$	(-0.7, -1.2, -0.3, 0, 1, 0, 0, 0)	Turn to detect obstacle on left	0 $\pm$ 0	NoYes
$aG^{22}$	(1.4, 2.0, 0.2, 0, 0, 0, 0, 0)	Turn left	0 $\pm$ 0	No

530 While the robot easily achieved goals  $aG^1$  and  $aG^5$ , it could ~~not either~~ achieve other valid goals only  
531 a few times or not able to achieve them at all. ~~most of the other goals.~~ Goals  $aG^2$ ,  $aG^8$ ,  $aG^{10}$ ,  $aG^{11}$ ,  
532  $aG^{13}$ , and  $aG^{17}$  could be achieved only a few times whereas goals  $aG^4$ ,  $aG^9$ ,  $aG^{14}$ ,  $aG^{16}$ , and  $aG^{21}$   
533 could not be achieved at all. The reason for that is due to the lack of opportunity. For example, the  
534 mobile robot ~~has to~~ must be near a wall for the event of detecting an obstacle at the front or turning  
535 right to avoid an obstacle behind. The argument made in section 4.1 regarding reducing the size of  
536 the arena to increase the opportunity for learning is valid here too.



537 Goals  $aG^2$ ,  $aG^3$ ,  $aG^6$ ,  $aG^7$ , and  $aG^{15}$  could not be achieved due to the granularity of discretization.  
 538 For the experiments in this paper, the wheel speed and orientation are discretized into nine values  
 539 ranging from  $-\pi$  to  $\pi$ . The wheel speed difference for the events for those goals was too small hence  
 540 when discretized; the values returned are 0 resulting in no change to the wheel speed, i.e., the event  
 541 of the robot turning ~~left~~, or right is not detected. For example, consider  $aG^7$  where the goal is to  
 542 turn right by increasing the right wheel speed by 0.1 (also achieving the change in orientation of -  
 543 0.1). Discretization of the range of  $2\pi$  radians into 9 buckets gives the granularity of 0.7 radians, thus  
 544 making the change of 0.1 radians difficult to detect. This, however, does not mean that the goal is  
 545 invalid. It is a valid goal, just that, for the robot to be able to learn a goal of such precise transition  
 546 would require experiments to be run with lower granularity values of wheel speed and orientation,  
 547 which in turn increases the state space and the size of the Q table and drastically increases the time to  
 548 learn to achieve those goals.

Figure 5(a): Trajectory for  $aG^5$ Figure 5(b): Trajectory for  $aG^{22}$  (run for 25,000 steps)Figure 5(c): Simulation for  $aG^{22}$  run for 100,000 steps

549 Figure 5(a) shows the trajectory for  $aG^5$  (increase speed of both wheels) for one of the trials. The  
 550 robot learns to attain this goal. In effect, this goal means that the robot has to keep increasing the  
 551 speed of its wheels. Attaining the maximum speed for both wheels results in the robot not able to  
 552 achieve the goal anymore and thus receives a negative reward. The robot, however, is again able to  
 553 attain the goal. This continues until the end of the trial.

554 Figure 5(b) shows the trajectory for  $aG^{22}$  (turn left) for 25,000 steps. The robot is not able to learn to  
 555 achieve that goal. The trajectory, however, is surprising, showing long stretches of straight line. We  
 556 let that trial continue for 100,000 steps, the trajectory for which is shown in Figure 5(c). The robot  
 557 still does not learn to achieve the goal. This is because the change in the wheel speed-difference, due  
 558 to the event (2.0 radians per second for the left wheel speed), is too large for one-time step. In a  
 559 single step, -the maximum change can only be  $\pi/2$  radians as per the design of the action set. Hence,  
 560 the goal and, as such, appears to be unreasonable. The goals  $aG^{19}$  and  $aG^{20}$  too appear to be  
 561 unreasonable for the same reason, and as can be seen from Table 4, they too could not be achieved.  
 562  $aG^{12}$  is unreasonable because goal attributes are showing transition for Right and Front-Left sensors  
 563 without any transition for Front-Right. It is hard to imagine the location of the mobile robot in the  
 564 arena that will result in such an event.  $aG^{18}$  too appears unreasonable because considering the change  
 565 to the wheel speeds (1.2 and 0.5 radians per second), the transition in the orientation (-0.1 radians) is  
 566 too small.

567 Either such unreasonable events were to be experienced by the robot during the experience gathering  
 568 stage in the experiments run by Merrick et al. (Merrick, Siddique, and Rano 2016) could be due to  
 569 noise, delay in sensing or that the mobile robot might have got stuck and then unstuck to the wall

570 resulting in an invalid event ( $e_t = s_t - s_{t-1}$ ) or that the unreasonable events were due to an error in  
 571 clustering, resulting in cluster centroid not being a correct representative of the cluster. If latter was  
 572 the case, then it requires reanalysis of the generated clusters. Possible solutions to rectify the  
 573 incorrect representation of the cluster centroid could be to place a minimum threshold on the cluster  
 574 size or to shift the cluster centroids to the nearest valid attribute value. In any case, those~~These~~ goals  
 575 appear unreasonable and are marked as invalid in the table. Based on the findings of the above  
 576 experiments, for the experiments in the next section, we have removed the orientation attribute from  
 577 the RL state vector, reduced the size of the arenas and, not used any of the invalid goals.

## 578 5 Demonstration of how Primitive Goal-based Reward Functions can be Combined

579 Not all tasks can be represented as a single goal type. Consider an example detailed in (Dastani and  
 580 Winikoff 2011), if the task for a personal assistant agent that manages a user’s calendar is to book a  
 581 meeting, it can be represented as an achievement goal, however people’s schedules change and hence  
 582 to ensure that the meeting invite remains in the calendar of all the participants, the task is better  
 583 modeled by a combination of goal types. The goal can be represented as “achieve then maintain”  
 584 where the aim is to achieve the goal and then maintain it. As another example, consider a wall  
 585 following mobile robot. The robot has to first approach a wall and then maintain a set distance from  
 586 the wall either to its left or to its right side. This goal can be represented as “approach then maintain”  
 587 where the aim of the mobile robot is to first approach the goal state (i.e., a wall to its left or right) and  
 588 then maintain it. We term this as a compound goal-based reward function, as it can be built from  
 589 multiple primitive goal-based reward functions.

590 In this section, we demonstrate compound goal-based reward functions constructed using if-then  
 591 rules to trigger different primitive reward functions in different states. In this paper, the if-then rules  
 592 are hand-crafted as we aim to demonstrate that primitive reward functions can be combined to  
 593 motivate learning of complex behaviors. The question of how to do this autonomously is discussed as  
 594 an avenue for future work in Section 6.1 and 6.2.

### 595 5.1 Experimental Setup

596 To demonstrate compound goal-based reward functions, we use the e-puck robot in three new  
 597 environments. The environments are as shown in Figures 6(a), 6(b) and 6(c). The maze environment,  
 598 shown in Figure 6(a), has walls to form a simple maze. In this environment, the goal of the robot is to  
 599 follow a wall. That goal is actually a compound goal. In order to achieve the goalgoal, the robot has  
 600 to learn primitive goals detailed in Table 1, Table 2, Table 3 and Table 4. The compound function 1  
 601 details the if-then rules to achieve this goal. The environment with obstacles, shown in Figure 6(b),  
 602 has cylindrical and cuboid objects that act as obstacles. The goal of the robot is to learn to avoid  
 603 obstacles. The compound function 2 details the if-then rules to achieve that goal. The third  
 604 environment is shown in Figure 6(c) is a circular arena with tracks. The goal of the robot is to learn  
 605 to follow a track which is detailed by compound function 3. Experiments were run for the following  
 606 goals expressed using compound goal-based reward functions. The primitive reward functions shown  
 607 in the if-then rules (Function 1, Function 2 and Function 3) are the same as in Table 1, Table 2, Table  
 608 3 and Table 4.

609 Function 1) Wall following goal in the maze arena

---

if obstacle on the left  
      $aG^{17}$  – achieve turning left  
 elseif obstacle close on the left  
      $G^1$  – maintain moving forward

---

---

```

elseif obstacle on the right
    aG11 - achieve turning right
elseif obstacle close on the right
    G1 – maintain moving forward
elseif obstacle at the front and left /*i.e. corner on the left */
    achieve turning right
elseif obstacle at the front and right /* i.e. corner on the right */
    achieve turning left
elseif obstacle at the front
    aG11 - achieve turning right
elseif no obstacle nearby
    G1 – maintain moving forward
end

```

---

610  
611

### Function 2) Obstacle avoidance goal in the arena with obstacles

---

```

if obstacle on the left
    aG13 – achieve turning right
elseif obstacle on the right
    aG4 - achieve turning left
elseif obstacle at the front and/or side
    aG11 - achieve turning right
elseif obstacle at the back
    G1 – maintain moving forward
elseif no obstacle anywhere nearby
    G1 – maintain moving forward
end

```

---

612  
613

### Function 3) Track following goal in the circular arena with tracks

---

```

if the obstacle anywhere nearby
    aG11 - achieve turning right
elseif track to the left
    achieve turning left
elseif track to the right
    achieve turning right
elseif on the track
    G1 – maintain moving forward
end

```

---

614

615 We use the same Dyna-Q algorithm that is detailed in Section 4. Action selection was using the  
616 epsilon-greedy method with epsilon parameter set to 0.1 throughout the learning process. 10 trials  
617 were run for each of the experiment with each trial consisting of 25000 steps.

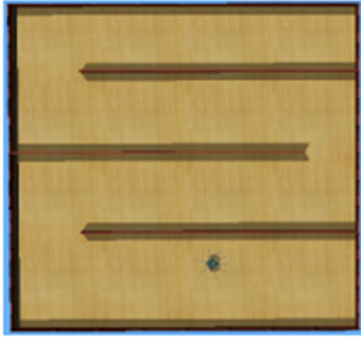


Figure 6(a): Maze arena

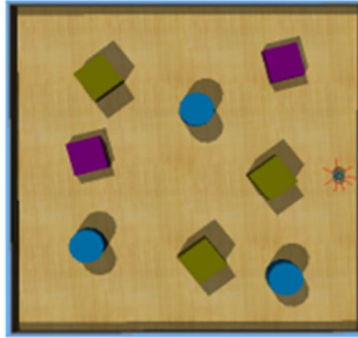


Figure 6(b): Arena with obstacles



Figure 6(c): A circular arena with tracks

618 The state space for this robot is different from that in Section 4. In addition to the six distance sensors  
 619 as detailed in the experiments in Section 4, we also use the ground sensors for these experiments. We  
 620 label the three ground sensors as *Ground-Left*, *Ground-Centre*, *Ground-Right*. The state of the  
 621 mobile robot comprises of following parameters: left wheel direction, right wheel direction, left  
 622 sensor value, right sensor value, front-left sensor value, front-right sensor value, rear-left sensor  
 623 value, rear right sensor value, ground left sensor value, ground center sensor value and ground right  
 624 sensor value. The state is a vector represented by  $[\omega^R \ \omega^L \ s^L \ s^R \ s^{FL} \ s^{FR} \ s^{RL} \ s^{RR} \ s^{GL} \ s^{GC} \ s^{GR}]$ .  $\omega^R$   
 625 and  $\omega^L$  are the rotational velocities of the right and the left wheels that are discretized to binary  
 626 values with 1 indicating that the wheel is moving forward and 0 indicating that it is moving  
 627 backwards. For the proximity sensors, we use binary values with 0 indicating that there is no object  
 628 in the proximity of the sensor and 1 indicates that the object is near. For ground sensors as [wellwell](#),  
 629 we use binary values with 0 indicating that the sensor is detecting light color and 1 indicating that it  
 630 is indicating dark color.

631 The action space comprises of three values: 1 – turn left, 2 – move forward and 3 – turn right.

## 632 5.2 Results

633 Table 5 shows the results of the wall following, obstacle avoidance, and track following goals.  
 634 Results were averaged over 10 trials, and its standard deviation is shown. The metrics used to  
 635 measure agent’s performance are the same as the ones defined in section 3 however here the metrics  
 636  $M_1$ ,  $M_2$ ,  $M_3$  and  $M_4$  measure cumulative reward gained by the agent for all the primitive goals  
 637 combined, i.e., the measurement for the compound goal-based reward.

Table 5: Results for compound goals

ID	Goal Description	$M_1$	$M_2$	$M_3$	$M_4$
G <sup>1</sup>	Wall following	1373 ±29	16833 ±115	10 ±0	78 ±6
G <sup>2</sup>	Avoiding obstacles	747 ±24	13613 ±109	11 ±0	81 ±8
G <sup>3</sup>	Following a track	992 ±24	14634 ±127	9 ±0	74 ±8

639

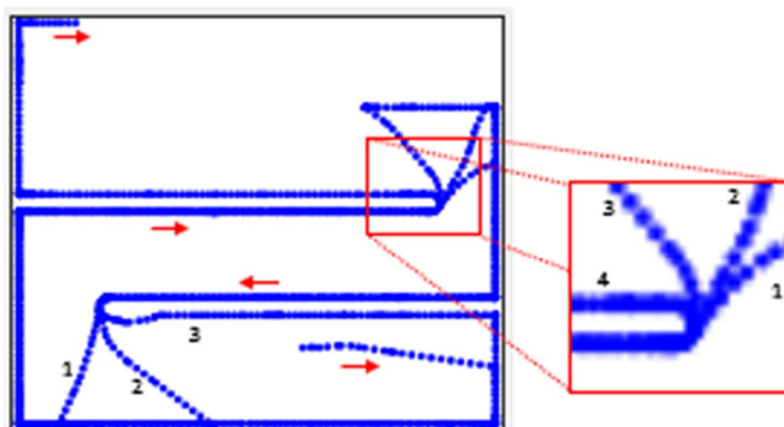


Figure 7: Trajectory for wall following goal in the maze arena

640 Figure 7 shows the trajectory for one of the trials of the mobile robot learning to follow the wall  
 641 using compound goal-based reward function (Function 1). The function comprises of a combination  
 642 of achievement and maintenance goal types each of which are triggered in a specific situation. When  
 643 there is no wall in the proximity, the robot is learning to move forward. Once it is near the wall  
 644 (either to the left or the right), it learns to follow the wall on that side. When it reaches the edge of the  
 645 wall, it is not able to follow it around for the initial two or three attempts however eventually learns  
 646 to follow the wall around and continues to follow the wall as shown in the zoomed-in section of  
 647 Figure 7. Trajectory labeled 4 in the zoomed-in section of Figure 7 is the one where the agent follows  
 648 the wall all the way around.



Figure 8(a): Trajectory for obstacle avoidance goal in the arena with obstacles



Figure 8(b): Trajectory for track following goal in the arena with tracks

649 Figure 8(a) shows the trajectory for one of the trials of the mobile robot learning to avoid obstacles  
 650 using the compound goal-based reward function (Function 2). This function too comprises a  
 651 combination of achievement and maintenance goal types each of which are triggered in a specific  
 652 situation. When there is no obstacle nearby, the robot has to learn to move forward. When it is close  
 653 to an obstacle, it has to learn to turn right and when it has the obstacle at its back it has to learn  
 654 to move forward, thus moving away from the obstacle. Figure 8(b) shows the trajectory for one of the  
 655 trials of the mobile robot learning to follow a track using the compound goal-based reward function  
 656 (Function 3). When the robot has a wall in its proximity, it has to learn to turn right. When near the

657 track, it has to learn to turn towards the track such that it is entirely on the track. Once on the track, it  
658 has to learn to move forward.

## 659 **6 Conclusion and Future Work**

660 This paper proposed reward functions for reinforcement learning based on the type of goal as  
661 categorized by the Belief Desire Intension community. The reward functions for the maintenance,  
662 approach, avoidance, and achievement goal types exploit the inherent property of its type, making  
663 them task-independent. Using simulated e-puck mobile robot experiments, we show how these  
664 intrinsic reward functions bridge the gap between autonomous goal generation and goal learning thus  
665 endowing the robot with the capability to learn in an autonomous and open-ended manner.

666 We present metrics to measure the agent's performance. The measurements show that using the  
667 proposed reward functions; all the valid goals will be learnt, some slower than the others due to the  
668 lack of opportunity. The goals that are not learnt are either very difficult to learn, unreasonable or  
669 invalid. The results also highlight the importance of attributes used in the design of the state vector as  
670 it can severely limit the learning opportunity, for example, usage of orientation attribute in the state  
671 vector. Although, this paper does not make any claim whether for or against any goal generation  
672 techniques, in the future work, the findings from this paper could be used to tune the goal generation  
673 technique used by Merrick et al. (Merrick, Siddique, and Rano 2016). We also show that the  
674 maintenance goals are easier to learn than the achievement goals. Approach and avoidance goals are  
675 even easier due to their inherent nature. This is because, for the maintenance goal, the agent is  
676 rewarded only when it can maintain the distance below a certain threshold, whereas, for approach and  
677 avoidance goals, the agent is rewarded for the approach or the avoidance attempt irrespective of its  
678 distance from the goal.

679 We further show how rather than treating the goal of a single type, the agent can decide whether it  
680 wants to maintain, approach, avoid or achieve the goal based on the situation it is experiencing. This  
681 situation specific goal type usage means the agent now knows what it has to learn in a specific  
682 situation thus directing the learning. A compound goal-based reward function can be designed by  
683 chaining any number of primitive reward functions. This raises following directions for future work.

### 684 **6.1 Autonomous Generation of Compound Reward Functions**

685 This paper demonstrated that primitive goal-based reward functions could be combined using if-then  
686 rules to create learnable compound reward functions. However, this raises a question whether it is  
687 possible for an agent to self-generate such rules or some other means of combining the primitive  
688 reward functions. One potential solution could be for the agent to autonomously determine the  
689 structure or regions in its state space each of which relates to a primitive goal. (Merrick, Siddique,  
690 and Rano 2016) have shown how the history of experienced states can be used to generate the goals.  
691 In a similar fashion, a coarse level clustering can be done on the experienced states to form these  
692 regions in the state space. Once those regions are known, one can then map the regions (primitive  
693 goal) with the goal state (compound goal) to enable the generation of the if-then rules. A formal  
694 framework is required for identifying complementary or conflicting goals so that complementary  
695 goals can be formed into compound reward functions and conflicting goals avoided.

### 696 **6.2 Conditions for Goal Accomplishment**

697 We also saw in this work that the agents learn solutions to some goals more effectively when they are  
698 in certain situations where the conditions support learning of that particular goal. This suggests that

709 there is a role for concepts such as opportunistic learning (Graham, Starzyk, and Jachyra 2012) to  
 700 maximize the efficiency of learning such that the agent only attempts goals that are feasible in a  
 701 given situation.

## 702 **7 Conflict of Interest**

703 *Paresh Dhakan, Kathryn Merrick, Inaki Rano and Nazmul Siddique declare that the research was*  
 704 *conducted in the absence of any commercial or financial relationships that could be construed as a*  
 705 *potential conflict of interest.*

## 706 **8 Author Contributions**

707 PD and KM conceived of the presented concept and planned the experiments. PD carried out the  
 708 experiments under the supervision of KM and IR. PD wrote the manuscript in consultation with KM,  
 709 IR and NS. All authors discussed the results, provided critical feedback and contributed to the final  
 710 version of the manuscript.

## 711 **9 References**

- 712 Amodei, Dario, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. 2016. "Concrete Problems  
 713 in AI Safety," 1–29. doi:1606.06565.
- 714 Andrychowicz, Marcin, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh  
 715 Tobin, Pieter Abbeel, and Wojciech Zaremba. 2017. "Hindsight Experience Replay," no. Nips.
- 716 Baldassarre, Gianluca, and Marco Mirolli. 2013. *Intrinsically Motivated Learning in Natural and Artificial Systems*.  
 717 Edited by Gianluca Baldassarre and Marco Mirolli. *Intrinsically Motivated Learning in Natural and Artificial*  
 718 *Systems*. Springer Heidelberg. doi:10.1007/978-3-642-32375-1.
- 719 Baraldi, Andrea. 1998. "Simplified ART: A New Class of ART Algorithms." *International Computer Science Institute*.
- 720 Baranes, Adrien, and Pierre-Yves Oudeyer. 2010a. "Intrinsically Motivated Goal Exploration for Active Motor Learning  
 721 in Robots: A Case Study." *IEEE/RSJ 2010 International Conference on Intelligent Robots and Systems, IROS 2010*  
 722 *- Conference Proceedings*, 1766–73. doi:10.1109/IROS.2010.5651385.
- 723 Baranes, Adrien, and Pierre-Yves Oudeyer. 2010b. "Maturationally-Constrained Competence-Based Intrinsically  
 724 Motivated Learning." *2010 IEEE 9th International Conference on Development and Learning*, August. Ieee, 197–  
 725 203. doi:10.1109/DEVLRN.2010.5578842.
- 726 Barto, Andrew G., Marco Mirolli, and Gianluca Baldassarre. 2013. "Novelty or Surprise?" *Frontiers in Psychology* 4  
 727 (DEC): 1–15. doi:10.3389/fpsyg.2013.00907.
- 728 Bonarini, Andrea, Alessandro Lazaric, and Marcello Restelli. 2006. "Incremental Skill Acquisition for Self-Motivated  
 729 Learning Animats." *Proceedings of the Ninth International Conference on Simulation of Adaptive Behavior (SAB-*  
 730 *06)* 4095: 357–68.
- 731 Braubach, Lars, Alexander Pokahr, Daniel Moldt, and Winfried Lamersdorf. 2005. "Goal Representation for BDI Agent  
 732 Systems." *Second International Workshop on Programming Multiagent Systems: Languages and Tools*, 9–20.  
 733 doi:10.1007/978-3-540-32260-3\_3.
- 734 Dastani, Mehdi, and Michael Winikoff. 2011. "Rich Goal Types in Agent Programming." In *In The 10th International*  
 735 *Conference on Autonomous Agents and Multiagent Systems*, 405–12.
- 736 Dewey, Daniel. 2014. "Reinforcement Learning and the Reward Engineering Principle." *AAAI Spring Symposium Series*,  
 737 1–8.

- 738 Duff, Simon, James Harland, and John Thangarajah. 2006. "On Proactivity and Maintenance Goals." *Proceedings of the*  
739 *Fifth International Joint Conference on Autonomous Agents and Multiagent Systems - AAMAS '06*, 1033.  
740 doi:10.1145/1160633.1160817.
- 741 Elliot, Andrew J. 2008. *Handbook of Approach and Avoidance Motivation*. Psychology Press.  
742 doi:10.1017/CBO9781107415324.004.
- 743 Graham, James T., Janusz A. Starzyk, and Daniel Jachyra. 2012. "Opportunistic Motivated Learning Agents" 7268  
744 (April). doi:10.1007/978-3-642-29350-4.
- 745 Held, David, Xinyang Geng, Carlos Florensa, and Pieter Abbeel. 2017. "Automatic Goal Generation for Reinforcement  
746 Learning Agents." <http://arxiv.org/abs/1705.06366>.
- 747 Hindriks, Koen V., and M. Birna Van Riemsdijk. 2007. "Satisfying Maintenance Goals." In *In International Workshop*  
748 *on Declarative Agent Languages and Technologies*, 4897 LNAI:86–103. doi:10.1007/978-3-540-77564-5\_6.
- 749 Laud, Adam, and Gerald DeJong. 2002. "Reinforcement Learning and Shaping: Encouraging Intended Behaviors." In  
750 *Proceedings of International Conference on Machine Learning*.
- 751 Merrick, Kathryn E. 2007. "Modelling Motivation For Experience-Based Attention Focus In Reinforcement Learning."  
752 School of Information Technologies, University of Sydney.
- 753 Merrick, Kathryn E. 2012. "Intrinsic Motivation and Introspection in Reinforcement Learning." *IEEE Transactions on*  
754 *Autonomous Mental Development* 4: 315–29. doi:10.1109/TAMD.2012.2208457.
- 755 Merrick, Kathryn E., and Mary Lou Maher. 2009. "Motivated Reinforcement Learning: Curious Characters for Multiuser  
756 Games." In *Motivated Reinforcement Learning: Curious Characters for Multiuser Games*, 1–206. Berlin,  
757 Heidelberg: Springer Berlin Heidelberg. doi:10.1007/978-3-540-89187-1.
- 758 Merrick, Kathryn E., Nazmul Siddique, and Inaki Rano. 2016. "Experience-Based Generation of Maintenance and  
759 Achievement Goals on a Mobile Robot." *Paladyn, Journal of Behavioral Robotics*, 67–84. doi:10.1515/pjbr-2016-  
760 0006.
- 761 Mirolli, Marco, and Gianluca Baldassarre. 2013. "Functions and Mechanisms of Intrinsic Motivations. The Knowledge  
762 Versus Competence Distinction." In *Intrinsically Motivated Learning in Natural and Artificial Systems*, 49–72.  
763 doi:10.1007/978-3-642-32375-1.
- 764 Neto, Hugo Vieira, and Ulrich Nehmzow. 2004. "Visual Novelty Detection for Inspection Tasks Using Mobile Robots."  
765 *In Towards Autonomous Robotic Systems: Proceedings of the 5th British Conference on Mobile Robotics*  
766 *(TAROS'04)*.
- 767 Nehmzow, U., Gatsoulis, Y., Kerr, E., Condell, J., Siddique, N. and McGuinnity, T.M., 2013. Novelty detection as an  
768 intrinsic motivation for cumulative learning robots. In *Intrinsically Motivated Learning in Natural and Artificial*  
769 *Systems* (pp. 185-207). Springer, Berlin, Heidelberg.
- 770 Ng, Andrew Y., Daishi Harada, and Stuart Russell. 1999. "Policy Invariance under Reward Transformations : Theory and  
771 Application to Reward Shaping." *Sixteenth International Conference on Machine Learning* 3: 278–87.
- 772 Oudeyer, Pierre-Yves, and Frederic Kaplan. 2007. "What Is Intrinsic Motivation? A Typology of Computational  
773 Approaches." *Frontiers in Neurorobotics* 1 (November): 6. doi:10.3389/neuro.12.006.2007.
- 774 Oudeyer, Pierre-Yves, Frederic Kaplan, and Verena V Hafner. 2007. "Intrinsic Motivation Systems for Autonomous  
775 Mental Development." *IEEE Transactions On Evolutionary Computation* 2 (2): 265–86.
- 776 Rao, Anand S, and Michael P Georgeff. 1995. "BDI Agents: From Theory to Practice." *Icmas* 95: 312–19.  
777 doi:10.1.1.51.9247.
- 778 Regev, Gil, and Alain Wegmann. 2005. "Where Do Goals Come from : The Underlying Principles of Goal-Oriented



- 779 Requirements Engineering.” In *International Conference on Requirements Engineering*, 353–62.  
780 doi:10.1109/RE.2005.80.
- 781 Rolf, Matthias, Jochen J. Steil, and Michael Gienger. 2010. “Bootstrapping Inverse Kinematics with Goal Babbling.”  
782 *2010 IEEE 9th International Conference on Development and Learning, ICDL-2010 - Conference Program*, no.  
783 May 2014: 147–54. doi:10.1109/DEVLRN.2010.5578850.
- 784 Santucci, Vieri Giuliano, Gianluca Baldassarre, and Marco Mirolli. 2010. “Biological Cumulative Learning through  
785 Intrinsic Motivations: A Simulated Robotic Study on the Development of Visually-Guided Reaching.” *Proceedings*  
786 *of the Tenth International Conference on Epigenetic Robotics: Modeling Cognitive Development in Robotic*  
787 *Systems* 0: 121–28. [http://www.im-clever.eu/publications/publications/pdfs/incollecionreference.2012-01-](http://www.im-clever.eu/publications/publications/pdfs/incollecionreference.2012-01-09.4513668372.pdf)  
788 [09.4513668372.pdf](http://www.im-clever.eu/publications/publications/pdfs/incollecionreference.2012-01-09.4513668372.pdf).
- 789 Santucci, Vieri Giuliano, Gianluca Baldassarre, and Marco Mirolli. 2010. 2012. “Intrinsic Motivation Mechanisms for  
790 Competence Acquisition.” In *IEEE International Conference on Development and Learning*, 1–6.  
791 doi:10.1109/DevLrn.2012.6400835.
- 792 Santucci, Vieri Giuliano, Gianluca Baldassarre, and Marco Mirolli. 2010. 2016. “GRAIL: A Goal-Discovering Robotic  
793 Architecture for Intrinsically-Motivated Learning.” *IEEE Transactions on Cognitive and Developmental Systems* 8  
794 (3): 214–31. doi:10.1109/TCDS.2016.2538961.
- 795 Singh, Satinder, Andrew G. Barto, and Nuttapon Chentanez. 2004. “Intrinsically Motivated Reinforcement Learning.”  
796 In *18th Annual Conference on Neural Information Processing Systems (NIPS)*, 17:1281–1288.  
797 [http://machinelearning.wustl.edu/mlpapers/paper\\_files/NIPS2005\\_724.pdf](http://machinelearning.wustl.edu/mlpapers/paper_files/NIPS2005_724.pdf).
- 798 Sutton, Richard S., and Andrew G. Barto. 1998. *Reinforcement Learning: An Introduction*. MIT Press Cambridge.  
799 doi:10.1.1.32.7692.
- 800 Thrun, Sebastian B, and Tom M Mitchell. 1995. “Lifelong Robot Learning.” *Robotics and Autonomous Systems* 15  
801 (March 1993): 25–46. doi:10.1016/0921-8890(95)00004-Y.
- 802 van Lamsweerde, Axel. 2001. “Goal-Oriented Requirements Engineering: A Guided Tour.” In *Proceedings Fifth IEEE*  
803 *International Symposium on Requirements Engineering*, 249–62. Toronto. doi:10.1109/ISRE.2001.948567.
- 804 van Riemsdijk, M. Birna, Mehdi Dastani, and Michael Winikoff. 2008. “Goals in Agent Systems: A Unifying  
805 Framework.” In *Proc. of 7th Int. Conf. on Autonomous Agents and Multiagent Syst.*, 713–20.
- 806 Weng, Juyang, James McClelland, Alex Pentland, Olaf Sporns, Ida Stockman, Mriganka Sur, and Esther Thelen. 2001.  
807 “Autonomous Mental Development by Robots and Animals.” *Science* 291 (5504): 599–600.  
808 doi:10.1126/science.291.5504.599.
- 809