

Appendix: Automaton-Guided Dynamic Task Sampling for Reinforcement Learning Agents

Anonymous submission

Abstract

Reinforcement Learning (RL) has made significant strides in enabling artificial agents to learn diverse behaviors. However, learning an effective policy often requires a large number of environment interactions. To mitigate sample complexity issues, recent approaches have used high-level task specifications, such as Linear Temporal Logic (LTL_f) formulas or Reward Machines (RM), to guide the learning progress of the agent. In this work, we propose a novel approach, called Automaton-Guided Dynamic Task Sampling (*AGTS*), that learns a set of RL policies to guide an agent from an initial state to a goal state based on a high-level LTL_f task objective, while minimizing the number of environmental interactions. Unlike previous work, *AGTS* does not assume information about the environment dynamics or the Reward Machine, and dynamically samples promising tasks that lead to successful goal policies. We evaluate *AGTS* on a grid-world and show that it achieves improved time-to-threshold performance on complex sequential decision-making problems compared to state-of-the-art RM and Automaton-guided RL baselines, such as Q-Learning for Reward Machines and Compositional RL from logical Specifications (DiRL). Moreover, we demonstrate that our method outperforms RM and Automaton-guided RL baselines in terms of sample-efficiency, both in a partially observable robotic task and in a continuous control robotic manipulation task.

1 Introduction

Agents are now capable of learning optimal control behavior for a broad spectrum of tasks, ranging from Atari games (Gao and Wu 2021) to robotic manipulation tasks (Nguyen and La 2019), thanks to recent advancements in Reinforcement Learning (RL). Despite the progress made in RL, learning an optimal task policy using model-free RL techniques still suffers from sample-complexity issues because of sparse reward settings and unknown transition dynamics (Lattimore, Hutter, and Sunehag 2013). These challenges further intensify in long-horizon settings, where the agent needs to perform a series of correct sequential decisions to achieve the goal. To alleviate this issue in complicated tasks, several lines of work have explored representing the goal using an intricately shaped reward function that guides the agent toward the goal (Grzes 2017). However, generating such a reward function requires the human engineer to assign ‘importance’ weights to various aspects of the

task, which is time consuming and assumes knowledge on which aspects of the task are important. Poorly engineered reward functions can lead to local optima, where the agent learns to satisfy only a subset of goals and ignores the rest.

Recent research has investigated representing the goal using high-level specification languages, such as finite-trace Linear Temporal Logic (LTL_f) (De Giacomo and Vardi 2013) or Reward Machines (RM) (Icarte et al. 2022) that allow defining the goal using a graphical representation of sub-tasks. The high-level objective is known before commencing the task, and the graphical representation allows the agent to achieve easier sub-goals initially, and build upon them to achieve complex goals. RM approaches still require human guidance in defining the reward structure of the machine, which is dependent on knowing how much reward should be assigned for each sub-goal. In contrast, our method does not require access to the reward structure. Compositional RL from Logical Specifications (DiRL) (Jothimurugan et al. 2021) mitigates this issue by using Dijkstra’s algorithm to reach those nodes in the graph that yield the highest success rate, encouraging the agent to learn policies for satisfying *all* outgoing edge propositions from such nodes. However, this approach requires the agent to explore a sub-task for a manually specified number of interactions, which requires knowledge about the environment. DiRL ends up spending a lot of interactions learning unproductive policies as some sub-tasks can be unpromising, yet the agent has to spend the defined number of interactions learning a policy for the sub-task. Unlike DiRL, our approach finds unpromising tasks and discards them, saving costly interactions and converging to a successful policy faster. Minimizing the overall number of interactions while learning a set of successful policies is non-trivial as this problem is equivalent to finding the shortest path in a graph whose edge weights are unknown *a priori* (Szepesvári 2004). In our case, the edge weight denotes the total number of environmental interactions required by the agent to learn a successful policy for a sub-task, in which the agent must induce a visit to a state where certain properties hold true. Additionally, we can only sample interactions for a sub-task if we have a policy that can reach the edge’s source node from the start node of the graph.

To address the above challenges, we present Automaton-Guided Dynamic Task Sampling (*AGTS*). We begin with

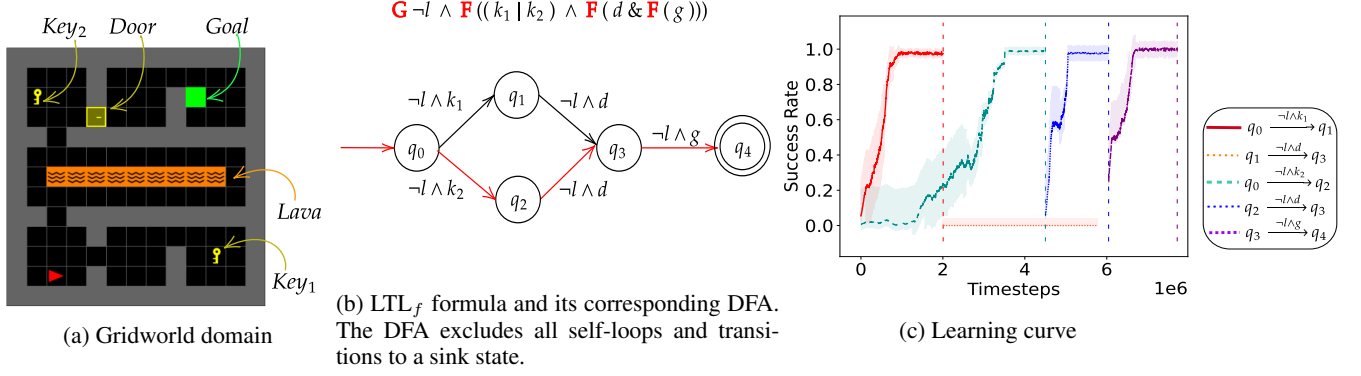


Figure 1: (a) Gridworld domain and descriptors. The agent (red triangle) needs to collect one of the keys and open the door to reach the goal; (b) The LTL_f formula for the task and its DFA. The operators **G** and **F** denote *always* and *eventually* respectively. Formulas l , k_1 , k_2 , d and g correspond to *Lava*, *Key1*, *Key2*, *Door* and *Goal* respectively; (c) Learning curves for individual sub-tasks (averaged over 10 trials) generated using AGTS. The path chosen by AGTS is highlighted in red in Fig.1(b)

a high-level objective represented using finite-trace Linear Temporal Logic (LTL_f) formulas which can equivalently be represented using Deterministic Finite Automaton (DFA). The DFA structure encodes memory, helping the agent understand what events of interest have occurred in the past, and which events must occur to reach the accepting states. Our key insight is to learn policies for sub-tasks defined using the edges of the DFA. Specifically, a transition from a state q to p occurs in a DFA when the propositional logic formula labeling the edge (q, p) evaluates to true. We use the set of propositional logic formulas labeling the outgoing edges from a given node in DFA to define sub-tasks. The trajectory induced by a successful RL policy for the sub-task enables the DFA's state to transition from the source node to the destination node of the edge defining the sub-task. We employ an adaptive Teacher-Student learning strategy, where, (1) the Teacher agent uses its high-level policy along with exploration techniques to actively sample a sub-task for the Student agent to learn. The high-level policy considers the DFA representation and the Student agent's expected performance on all the sub-tasks, aiming to satisfy the high-level objective in the fewest number of interactions, and (2) the Student agent interacts with the environment for a few steps (much fewer than the interactions required to learn a successful policy for the sub-task) while updating its low-level RL policy for the sampled sub-task. The Teacher observes the Student's performance on these interactions and updates its high-level policy. Steps (1) and (2) continue alternately until the Student agent learns a set of successful policies that guide the agent to reach a goal state.

As an example, let us look at the environment shown in Fig. 1a. The goal for the agent is to collect *any* of the two *Keys*, followed by opening the *Door* and then reaching the *Goal* while avoiding the *Lava* at all times. The task's high-level objective is represented using the LTL_f formula and its corresponding DFA representation \mathcal{U}^1 in Fig. 1b. The DFA does not contain information about the environment config-

uration, such as: the optimal number of interactions required to reach *Door* from *Key1* are much higher compared to the interactions required to reach *Door* from *Key2*, making the *Key1* to *Door* trajectory sub-optimal. Hence, it is crucial to prevent any additional interactions the agent spends in learning a policy for the sub-task defined by the edge $q_1 \xrightarrow{\neg l \wedge d} q_3$ as the path $q_0 \xrightarrow{\neg l \wedge k_1} q_1 \xrightarrow{\neg l \wedge d} q_3 \xrightarrow{\neg l \wedge g} q_4$ will always be sub-optimal. Our proposed approach, AGTS begins with the aim of learning two distinct policies π_1 for the task of visiting *Key1* and π_2 for the task of visiting *Key2*, both avoiding *Lava*. The Teacher initially samples evenly from these two sub-tasks but later biases its sampling toward the sub-task on which the Student agent shows higher learning potential. Once the agent learns a successful policy for one of the sub-tasks (let's say the learned policy π_1^* corresponding to the sub-task defined by the transition $q_0 \xrightarrow{\neg l \wedge k_1} q_1$), the Teacher does not sample that task anymore, identifies the next task(s) in the DFA representation, and appends them to the set of tasks it is currently sampling (in this case, the only next task is: $q_1 \xrightarrow{\neg l \wedge d} q_3$). Since the agent only has access to the state distribution over q_0 , it follows the trajectory given by π_1^* to reach a state that lies in the set of states where the proposition $\neg \text{Lava} \wedge \text{Key}_1$ holds true before commencing its learning for the policy (π_3) for $q_1 \xrightarrow{\neg l \wedge d} q_3$. If the agent learns the policies π_2^* for satisfying the sub-task defined by $q_0 \xrightarrow{\neg l \wedge k_2} q_2$ and π_4^* for $q_2 \xrightarrow{\neg l \wedge d} q_3$ before learning π_3 , it effectively has a set of policies to reach the node q_3 . Thus, the Teacher will now only sample the next task in the DFA representation $q_3 \xrightarrow{\neg l \wedge g} q_4$, as learning policies for paths that reach q_3 are effectively redundant. This process continues iteratively until the agent learns a set of policies that reach the goal node (q_4) from the start node (q_0). The learning curves in Fig. 1c empirically validate the running example. As evident from the learning curves, the agent learns policies for the path $q_0 \xrightarrow{\neg l \wedge k_2} q_2 \xrightarrow{\neg l \wedge d} q_3 \xrightarrow{\neg l \wedge g} q_4$ that produce trajectories to reach the goal node q_4 from the initial

¹Cycles and sink nodes are pruned for ease of understanding

node q_0 , without excessively wasting interactions on the unpromising sub-task $q_1 \xrightarrow{\neg l \wedge d} q_3$. The dashed lines in Fig. 1c signify the interactions at which a task policy converged.

The dynamic task sampling strategy promotes *AGTS* to achieve sample-efficient learning on complex temporally extended tasks by identifying unpromising tasks and discarding them, saving costly interactions. Our empirical results show that *AGTS* reduces environmental interactions by orders of magnitude compared to state-of-the-art Specifications-Guided RL Baseline DiRL, Reward Machine-based baselines QRM (Icarte et al. 2018), GSRS (Camacho et al. 2018), and curriculum learning baseline TSCL (Matiisen et al. 2020). We also evaluate *AGTS^{ct}*, a modified algorithm that further improves sample efficiency by continuing exploration on a new sub-task once a goal state for a sub-task is reached. We perform evaluation on two robotic navigation and manipulation tasks and demonstrate that *AGTS* reduces the number of interactions by orders-of-magnitude when compared to state-of-the-art automaton-guided RL baselines.

2 Related Work

Automaton-guided RL approaches utilize temporal logic-based language specifications to define tasks (Toro Icarte et al. 2018; Bozkurt et al. 2020; Xu and Topcu 2019; Alur et al. 2022; De Giacomo et al. 2019). Separating policies for task sub-goals aids in abstracting knowledge that can be utilized in novel tasks (Icarte et al. 2018), without reliance on a dense reward function. Another technique is to shape the reward in proportion to the distance from the accepting node in the automaton (Camacho et al. 2018); however, this often leads to suboptimal reward settings. Augmenting the reward function with Monte Carlo Tree Search helps mitigate this issue (Velasquez et al. 2021). This approach requires the ability to plan-ahead in the environment, which is not always feasible. Automaton-guided RL has been used to aid navigational exploration for robotic domains (Cai et al. 2023; Li, Vasile, and Belta 2017) and for multi-agent settings (Hammond et al. 2021). Generating a curriculum given the high-level objective (Shukla et al. 2023) requires access to the Object-Oriented MDP (Diuk, Cohen, and Littman 2008), which cannot be obtained if environment details are not known in advance. DiRL interleaves high-level planning with RL to learn a policy for each edge, which overcomes challenges arising from poor representations (Jothimurugan et al. 2021). This approach becomes inefficient in terms of number of interactions, as it requires the agent to act for a predetermined number of interactions, even if learning the task does not show any promise. Unlike previous works, in this paper, we propose an automaton-guided dynamic task sampling approach that does not require access to the environment dynamics or the Reward Machine, and efficiently samples tasks that show promise toward the high-level objective, saving interactions on unpromising tasks.

Teacher-Student algorithms (Matiisen et al. 2020) have been previously studied in Curriculum Learning literature (Narvekar et al. 2020; Shukla et al. 2022) and in the Intrinsic Motivation literature (Oudeyer and Kaplan 2009).

The idea is to have the Teacher propose those tasks to Student on which the Student shows most promise. This strategy helps Student learn simpler tasks first, transferring its knowledge to complex tasks. The technique reduces the overall number of interactions necessary to learn a successful policy. These approaches tend to optimize a curriculum to learn a single policy, which does not scale well to temporally-extended tasks. Instead, we propose an Automaton-guided Teacher-Student learning strategy that learns a policies for promising automaton transitions, promoting sample-efficient training compared to the baselines.

3 Theoretical Framework

Episodic MDP. An episodic labeled Markov Decision Process (MDP) M is a tuple $(\mathcal{S}, \mathcal{A}, \mathcal{V}, R, \mathcal{S}_0, \gamma, K, AP, L)$, where \mathcal{S} is the set of states, \mathcal{A} is the set of actions, $\mathcal{V}(s'|s, a)$ denotes the transition probability of reaching state $s' \in \mathcal{S}$ from $s \in \mathcal{S}$ using action $a \in \mathcal{A}$, $R : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ is the reward function, \mathcal{S}_0 is the initial state distribution, $\gamma \in [0, 1]$ is the discount factor, K is the maximum number of interactions in any episode, AP is a set of atomic propositions, and $L : \mathcal{S} \rightarrow 2^{AP}$ is a labeling function that maps a state $s \in \mathcal{S}$ to a subset of atomic propositions that are true in that state.

In every interaction, the agent observes the current state s and selects an action a according to its policy function $\pi(a|s, \theta)$ with parameters θ . The MDP transitions to a new state $s' \in \mathcal{S}$ with probability $\mathcal{V}(s' | s, a)$. The agent’s goal is to learn an *optimal policy* π^* that maximizes the discounted return $G_0 = \sum_{k=0}^K \gamma^k R(s'_k, a_k, s_k)$ until the end of the episode, which occurs after at-most K interactions.

High-level Specification. We express the high-level specification for a complex, temporally extended task finite-trace Linear Temporal Logic (LTL_f) formulas (De Giacomo and Vardi 2013). LTL_f succinctly expresses intricate temporal tasks such as safety, reachability, persistence, recurrence or a combination of these. LTL_f formulas are constructed inductively from a set of atomic propositions AP , using the operations $p, \neg\phi, \phi_1 \vee \phi_2, \mathbf{X}\phi, \mathbf{G}\phi, \mathbf{F}\phi$, and $\phi_1 \mathbf{U}\phi_2$. Here, $p \in AP$, and ϕ, ϕ_1 , and ϕ_2 are LTL_f formulas. The negation operator is represented by \neg , the disjunction operator is represented by \vee , and the operators \mathbf{X} , \mathbf{G} , \mathbf{F} , and \mathbf{U} correspond to *Next*, *Always*, *Eventually*, and *Until*, respectively.

The language of any LTL_f formula can be expressed equivalently by a DFA (De Giacomo and Vardi 2013). A DFA is a tuple $\mathcal{U} = (Q, q_0, \Sigma, \delta, F)$, where Q is the set of nodes, $q_0 \in Q$ is an initial node, $\Sigma = 2^{AP}$ is an alphabet defined over a set of atomic propositions AP , $\delta : Q \times \Sigma \rightarrow Q$ is a deterministic transition function, and $F \subseteq Q$ is the set of accepting nodes. A path $\zeta = q_0 q_1 \dots q_n$ in a DFA is a finite-length sequence of nodes such that, for any $i = 0 \dots n-1$, we have $q_{i+1} = \delta(q_i, \sigma_i)$. A path ζ is said to be accepting if its last node is accepting, i.e., $q_n \in F$.

Every path $\rho = s_0 s_1 \dots s_n$ in MDP induces a word $w = L(s_0)L(s_1) \dots L(s_n) \in \Sigma^*$. The word w induces a path $\zeta = q_0 q_1 \dots q_n$ in the DFA. A path ρ in MDP satisfies an LTL_f formula φ if and only if the last state q_n in the path ζ induced by ρ is an accepting state in the DFA corresponding to φ .

Problem Formulation. Given an MDP M with unknown

transition dynamics and an LTL_f formula φ representing the high-level objective of the agent, let \mathcal{U} be the DFA representing the language of φ . Let $\text{Paths}(q, X)$ be the set of all finite paths originating in q and terminating at a node in $X \subseteq Q$. The aim of this work is to learn a set of policies π_i^* , $i = 0, \dots, n-1$, with the following three properties: (i) Following π_0^* results in a path in MDP that induces a transition from q_0 to some state $q_1 \in Q$ in the DFA, following π_1^* results in a path in MDP that induces a transition from q_1 to some state $q_2 \in Q$ in the DFA, and so on. (ii) The resulting path $q_0 q_1 \dots q_n$ in the DFA terminates at an accepting state, i.e., $q_n \in F$, with probability greater than a given threshold, $\eta \in (0, 1)$. We intend to minimize the number of environmental interactions required to learn these policies.

4 Methodology

Given the DFA \mathcal{U} representing the language of φ , we define a set of sub-tasks based on the edges of the DFA. Intuitively, given any MDP state $s \in \mathcal{S}$ and a DFA node $q \in Q$, a sub-task defined by an edge from node q to $p \in Q$ defines a reach-avoid objective for the agent represented by the LTL_f formula,

$$\text{Task}(q, p) = \mathbf{F}(\phi_{(q,p)}) \wedge \mathbf{G} \left(\bigwedge_{r \in \text{Succ}(q), r \neq p} \neg \phi_{(q,r)} \right),$$

where $\phi_{(q,p)}$ is the propositional formula labeling the edge from q to p in the DFA and $\text{Succ}(q)$ is the set of successors of node q in DFA. For example, in Fig. 1b, the propositional logic formula labeling the edge from q_0 to q_1 is $\phi_{(q_0,q_1)} = \neg l \wedge k_1$. When $e = (q, p)$, we sometimes write $\text{Task}(e)$ instead of $\text{Task}(q, p)$ for notational convenience.

Each sub-task $\text{Task}(q, p)$ defines a problem to learn a policy $\pi_{(q,p)}^*$ such that, given any MDP state $s_0 \in \mathcal{S}$, following $\pi_{(q,p)}^*$ results in a path $s_0 s_1 \dots s_n$ in MDP that induces the path $qq \dots qp$ in the DFA. That is, the DFA state remains at q until it transitions to p . While constructing the set of sub-tasks, we omit the self-loops and the transitions that lead to a ‘sink’ state (from which final states are unreachable).

Given the MDP M with unknown transition dynamics and the LTL_f objective, φ , we first translate φ to its corresponding Deterministic Finite Automaton (DFA) representation $\mathcal{U} = (Q, q_0, \Sigma, \delta, F)$. Next, we define the set of sub-tasks. For this, we consider the edges that lie on some path from q_0 to some state in F . This is because any path that does not visit F must lead to a sink state from which the objective cannot be satisfied. The set of such edges is identified using breadth-first-search (Moore 1959).

The algorithm for AGTS is described in Algo. 1. We begin by initializing the following quantities (lines 2-4): (1) Set of: Active Tasks AT, Learned Tasks LT, Discarded Tasks DT; (2) A Dictionary \mathcal{P} that maps a sub-task $\text{Task}(e)$ of the DFA \mathcal{U} to a policy π_e ; (3) A Dictionary that represents the Teacher Q-Values Q by mapping a sub-task $\text{Task}(e)$ to a numerical value associated with that sub-task.

Firstly, we convert \mathcal{U} into an Adjacency Matrix \mathcal{X} (line 6), and find the set of all the outgoing edges $\bar{\mathcal{E}}_{q_0} \subseteq \mathcal{E}$

from the initial node q_0 (line 7). Satisfying the edge’s formula $\phi_{(q_0,q_1)} \in \bar{\mathcal{E}}_{q_0}$ represent a reachability sub-task M' where each goal state $s \in \mathcal{S}_f^{M'}$ of M' satisfy the condition $\phi_{(q_0,q_1)} \subseteq L(s)$. The agent receives a positive reward for satisfying $\phi_{(q_0,q_1)}$ and a small negative reward at all other time steps. The state space, the action space and the transition dynamics of M' are equivalent to M . To complete the sub-task, the agent must learn a policy $\pi_{(q,p)}^*$ that ensures a visit from q to p with probability greater than a predetermined threshold (η). Moreover, the policy must also avoid triggering unintended transitions in the DFA. For instance, while picking up Key_1 , the policy must not inadvertently pick up Key_2 .

We set the Teacher Q-Values for all the sub-tasks corresponding to edges in AT (i.e., $\bar{\mathcal{E}}_{q_0}$) to zero (line 8). We formalize the Teacher’s goal of choosing the most promising task as a *Partially Observable MDP* (Kaelbling, Littman, and Cassandra 1998), where the Teacher does not have access to the entire Student agent state but only to the Student agent’s performance on a task (e.g. success rate or average returns), and as an action, chooses a task $\text{Task}(e) \in \text{AT}$ the Student agent should train on next. In this POMDP setting, each Teacher action (a sub-task) has an Q-Value associated with it. Intuitively, higher Q-Values correspond to tasks on which the Student agent is more successful, and the Teacher should sample such tasks at a higher frequency to satisfy φ (reach a goal node) in fewest overall interactions.

(A) Given the Teacher Q-Values, we sample a sub-task $\text{Task}(e) \in \text{AT}$ using the ϵ -greedy exploration strategy (line 10), and (B) The Student agent trains on task $\text{Task}(e)$ using the policy $\mathcal{P}[e]$ for ‘ x ’ interactions (line 11). In one Teacher timestep, the Student trains for x environmental interactions. Here, $x \ll$ total number of environmental interactions required by the agent to learn a successful policy for $\text{Task}(e)$, since the aim is to keep switching to a task that shows highest promise. (C) The Teacher observes the Student agent’s average return g_t on these x interactions, and updates its Q-Value for $\text{Task}(e)$ (line 12):

$$Q[e] \leftarrow \alpha(g_t) + (1 - \alpha)Q[e] \quad (1)$$

where α is the Teacher learning rate, g_t is the average Student agent return on $\text{Task}(e)$ at the Teacher timestep t . As the learning advances, g_t increases as t increases, and hence we use a constant parameter α to tackle the non-stationary problem of a moving return distribution. Several other algorithms could be used for the Teacher strategy (e.g., UCB (Agrawal and Goyal 2012), Thompson Sampling (Auer, Cesa-Bianchi, and Fischer 2002)). Steps (A), (B) and (C) continue successively until the policy for *any* $\text{Task}(e) \in \text{AT}$ task converges.

We define a policy for $\text{Task}(q, p)$ to be converged (line 13) if a trajectory ω produced by the policy triggers the transition with probability $\Pr_{\omega \in \Omega}[\omega \text{ satisfies } \text{Task}(q, p)] \geq \eta$ and $\Delta(g_t, g_{t-1}) < \tau$ where η is the expected performance and τ is a small numerical value. Intuitively, a converged policy attains an average success rate $\geq \eta$ and has not improved further by maintaining constant average returns. Like all other RM and automaton-based approaches, we assume

access to the labeling function L to examine if the trajectory ω satisfies the formula $\phi(q, p)$ by checking if the final environmental state s of the trajectory satisfies the condition $\phi(q, p) \subseteq L(s)$. The sub-goal regions need not be disjoint, i.e., the same state s can satisfy propositions for multiple DFA nodes. Once a policy for the $\text{Task}(q, p)$ converges, we append $\text{Task}(q, p)$ to the set of Learned Tasks LT and remove it from the set of Active Tasks AT (line 14). In order to ensure that the learned task does not get sampled any further, we set the Teacher Q-value for this sub-task to $-\infty$ (line 15). Once we have a successful policy for the $\text{Task}(q, p)$ (the transition $q \xrightarrow{\sigma} p$), we determine the sub-tasks that can be discarded (line 16). We find the sub-tasks corresponding to edges that: (1) lie before p in a path from q_0 to any $q \in F$, and, (2) do not lie in a path to $q \in F$ that does not contain p . Intuitively, if we already have a set of policies that can generate a successful trajectory to reach the node p , we do not need to learn policies for sub-tasks that ultimately lead *only* to p . We add all such sub-tasks to the set of Discarded Tasks DT (line 17), and set the Teacher Q-values for all the discarded tasks to $-\infty$ to prevent them from being sampled for the Student learning agent (line 18). As an extension, in the limit, an optimal policy can be found by not completely discarding these sub-tasks, but rather biasing away from them so that in the limit they would still be explored.

Subsequently, we determine the next set of tasks $\bar{\mathcal{E}}_{AT}$ in the DFA to add to the AT set (line 19). This is calculated by identifying sub-tasks corresponding to all the outgoing edges from p . Since the edge $e_{q,p}$ corresponds to the transition $q \xrightarrow{\sigma} p$, we have a successful policy that can produce a trajectory that ends in a state where the propositions for p hold true, and $\bar{\mathcal{E}}_{AT}$ corresponds to $\mathcal{X}[p] \setminus \text{DT}$ (' \setminus ' refers to set-minus) i.e., sub-tasks corresponding to all the outgoing edges from p that do not lie in the DT set.

Once we identify $\bar{\mathcal{E}}_{AT}$, we set the Teacher Q-values for all $\text{Task}(\bar{e}) \in \bar{\mathcal{E}}_{AT}$ to 0 so that the Teacher will sample these tasks (line 23). We consider an episodic setting where the episode starts from a state $s \sim \mathcal{S}_0$ where the propositions for q_0 hold true, and if the current sampled sub-task is $\text{Task}(p, r)$, the agent follows a trajectory using corresponding learned policies from Π^* to reach a state where the propositions for p hold true, and then attempts learning a separate policy for $\text{Task}(p, r)$.

The above steps (lines 9-26) go on iteratively until $\bar{\mathcal{E}}_{AT}$ is an empty set. This indicates we have no further tasks to add to our sampling strategy, and we have reached a node $q \in F$. Thus, we **break** from the **while** loop (line 21) and return the set of learned policies Π^* , and edge-policy dictionary \mathcal{P} (line 27). From \mathcal{P} and Π^* , we get an ordered list of policies $\Pi^*_{list} = [\pi_{(q_1, q_2)}, \pi_{(q_2, q_3)}, \dots, \pi_{(q_{k-1}, q_k)}]$ such that sequentially following $\pi \in \Pi^*_{list}$ generates trajectories that satisfy the LTL_f objective φ^2 .

Guarantee: Given the ordered list of policies Π^*_{list} , we can generate a trajectory ω in the task M with $\Pr_{\omega \in \Omega}[\omega \text{ satisfies } \varphi] \geq \eta$ (Details in Appendix B).

²Link to code to be provided after review

Algorithm 1: *AGTS* (\mathcal{U}, M, η, x)

Output: Set of learned policies : Π^* , Edge-Policy Dictionary \mathcal{P}

```

1: Placeholder Initialization:
2: Sets of: Active Tasks ( $\text{AT}$ )  $\leftarrow \emptyset$ ;
   Learned Tasks ( $\text{LT}$ )  $\leftarrow \emptyset$ ; Discarded Tasks ( $\text{DT}$ )  $\leftarrow \emptyset$ 
3: Edge-Policy Dictionary  $\mathcal{P} : \text{Task}(e) \rightarrow \pi$ 
4: Teacher Q-Value Dictionary:  $Q : \text{Task}(e) \rightarrow -\infty$ 
5: Algorithm:
6:  $\mathcal{X} \leftarrow \text{Adjacency\_Matrix}(\mathcal{U})$ 
7:  $\text{AT} \leftarrow \text{AT} \cup \{\mathcal{X}[q_0]\}$ 
8:  $\forall \text{Task}(e) \in \text{AT} : Q[e] = 0$ 
9: while True do
10:   $e \leftarrow \text{Sample}(Q)$ 
11:   $\mathcal{P}[e], g \leftarrow \text{Learn}(M, \mathcal{U}, e, x, \mathcal{P})$ 
12:   $\text{Update\_Teacher}(Q, e, g)$ 
13:  if  $\text{Convergence}(Q, e, g, \eta)$  then
14:     $\Pi^* \leftarrow \Pi^* \cup \mathcal{P}[e]$ ;  $\text{LT} \leftarrow \text{LT} \cup \{\text{Task}(e)\}$ ;
     $\text{AT} \leftarrow \text{AT} \setminus \{\text{Task}(e)\}$ 
15:     $Q[e] = -\infty$ 
16:     $\bar{\mathcal{E}}_{DT} \leftarrow \text{Discarded\_Tasks}(\mathcal{X}, e)$ 
17:     $\text{DT} \leftarrow \text{DT} \cup \bar{\mathcal{E}}_{DT}$ 
18:     $\forall \text{Task}(\bar{e}) \in \bar{\mathcal{E}}_{DT} : Q[\bar{e}] = -\infty$ 
19:     $\bar{\mathcal{E}}_{AT} \leftarrow \text{Next\_Tasks}(\mathcal{X}, e, \text{DT})$ 
20:    if  $|\bar{\mathcal{E}}_{AT}| = 0$  then
21:      break
22:    end if
23:     $\forall \text{Task}(\bar{e}) \in \bar{\mathcal{E}}_{AT} : Q[\bar{e}] = 0$ 
24:     $\text{AT} \leftarrow \text{AT} \cup \bar{\mathcal{E}}_{AT}$ 
25:  end if
26: end while
27: return  $\Pi^*, \mathcal{P}$ 

```

5 Experimental Results

We aim to answer the following questions: (Q1) Does *AGTS* yield sample efficient learning compared to state-of-the-art baselines? (Q2) After reaching a sub-task goal state, can we sample a new sub-task to continue training and get better sample efficiency? (Q3) Does *AGTS* yield sample efficient learning for complex robotic tasks with partially observable or continuous control settings?

5.1 AGTS - Gridworld Domain

To answer (Q1), we evaluated *AGTS* on a Mini-grid (Chevalier-Boisvert, Willems, and Pal 2018) inspired domain with the LTL_f objective:

$$\varphi_f^{\text{gridworld}} := \mathbf{G} \neg \text{Lava} \wedge \mathbf{F}((\text{Key}_1 | \text{Key}_2) \wedge \mathbf{F}(\text{Door} \ \& \ \mathbf{F}(\text{Goal}))) \quad (2)$$

where $\text{Key}_1, \text{Key}_2, \text{Door}, \text{Goal}$ are the atomic propositions. The environment and the DFA representation are given in Fig. 1. Essentially, the agent needs to collect *any* of the *Keys* before heading to the *Door*. After *toggling* the *Door* open, the agent needs to visit the grid with the *Goal*. At all times, the agent needs to avoid the *Lava* object. We assume an episodic setting where an episode ends if the agent

Approach	# Interactions (Mean \pm SD)	Success Rate (Mean \pm SD)
<i>AGTS^{ct}</i>	$(5.75 \pm 0.38) \times 10^6$	0.96 ± 0.02
<i>AGTS</i>	$(6.12 \pm 0.25) \times 10^6$	0.95 ± 0.01
DiRL ^c	$(7.97 \pm 0.46) \times 10^6$	0.95 ± 0.03
DiRL	$(9.62 \pm 0.42) \times 10^6$	0.94 ± 0.01
QRM	5×10^7	0.05 ± 0.04
GSRS	5×10^7	0 ± 0
TSCL	5×10^7	0 ± 0
LFS	5×10^7	0 ± 0

Table 1: Table comparing #interactions & success rate. *AGTS* (highlighted) outperformed all baselines

touches the *Lava* object, reaches the *Goal* or exhausts the number of allocated interactions.

This is a complex sequential decision making problem as the agent needs to perform a series of correct actions to satisfy $\varphi_f^{\text{gridworld}}$. In this environment, the agent has access to three navigation actions: *move forward*, *rotate left* and *rotate right*. The agent can also perform: *pick-up* action, which adds the *Key* to the agent’s inventory if it is facing the *Key*, *drop* places the *Key* in the next grid if *Key* is present in the inventory, and, *toggle* that toggles the *Door* (closed \leftrightarrow open) only if the agent is holding the *Key*. For this environment, we assume a fully-observable setting where the environmental state is a low-level encoding of the image. For each cell in the grid, the low-level encoding returns an integer that describes the item occupying the grid, along with additional information, if any (e.g., the *Door* state can be open or closed).

For the RL pipeline, we use PPO (Schulman et al. 2017), which works for discrete and continuous action spaces. We consider a standard actor-critic architecture with 2 convolutional layers followed by 2 fully connected layers. For *AGTS*, the reward function is sparse. The agent gets a reward of $(1 - 0.9^{\frac{(\text{interactions taken})}{(\text{interactions allocated})}})$ if it achieves the goal in the sub-task, and a reward of 0 otherwise. For individual tasks, $\text{interactions allocated} = 100$. The agent does not receive any negative rewards for hitting the *Lava*.

We compare our *AGTS* method against six baseline approaches: learning from scratch (LFS), Reward Machine-based (RM) baselines: GSRS (Camacho et al. 2018), QRM (Icarte et al. 2018); and Compositional RL from Logical Specifications (DiRL) (Jothimurugan et al. 2021). All the baselines are implemented using the RL algorithm (PPO), described above. GSRS builds upon naive reward shaping by modifying the reward inversely proportional to the distance from the RM goal node and by learning policies for each RM node transition. QRM employs a separate Q-function for each node in the RM, and DiRL uses Dijkstra’s algorithm to guide the agent in choosing a path from the specification graph. The edge-costs for Dijkstra’s algorithm are calculated as the average RL policy success rate on that particular sub-task after interacting over a manually defined number of interactions. For the fifth baseline, we modify DiRL such that instead of manually specifying a limit on the

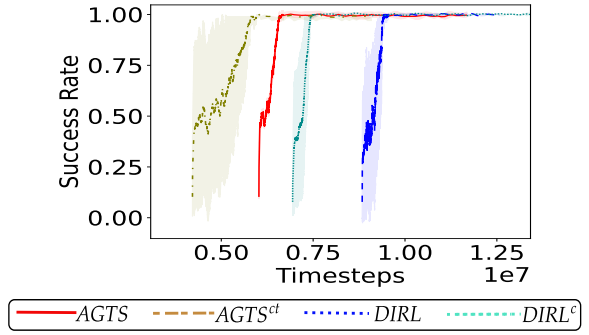


Figure 2: Averaged over 10 trials: Learning curves for approaches whose policies successfully converged.

number of interactions, which needs to be fine-tuned to suit the task, we stop learning a sub-task once it has reached the convergence criteria defined in Section 4. We call this modified baseline as DiRL^c. The sixth baseline (TSCL (Matiisen et al. 2020)) follows a curriculum learning strategy where the Teacher samples most promising task without the use of any automaton to guide the learning progress of the agent. (Baseline implementation details in Appendix C)

The results in Table 1 and Fig. 2 (averaged over 10 trials) show that *AGTS* reaches a successful policy quicker compared to the baseline approaches. *AGTS^{ct}* is a modified version of *AGTS* and is described in Sec. 5.1. The learning curves in Fig. 2 have an offset on the x-axis to account for the time spent in the initial sub-tasks before moving on to the final task in the automaton, signifying *strong transfer* (Taylor and Stone 2009). Our custom baseline, DiRL^c performs better than DiRL, and both outperform other baselines, which fail to learn a meaningful policy. A pure Teacher-Student learning strategy (TSCL) (without an automaton to guide the learning progress) fails to learn a successful policy. We performed an unpaired t-test (Kim 2015) to compare *AGTS* against the best performing baselines at the end of 10^7 training steps and we observed statistically significant results (95% confidence). Thus, *AGTS* not only achieves a better success rate, but also converges faster (more on statistical significance result details in Appendix Section D). Time-to-threshold metric is defined as the difference in number of interactions between two approaches to reach a desired performance (Narvekar et al. 2020). From Fig. 2, we see that the time-to-threshold between *AGTS* and the best-performing baseline DiRL^c is 1.85×10^6 interactions for a desired performance of $\eta = 95\%$ success rate.

AGTS^{ct} (AGTS + Cont. Training) - Gridworld Domain

In *AGTS*, while learning a policy for Task(q, p), we reinitialize the environment to a random initial environmental state $s \sim \mathcal{S}_0$ once the agent reaches a state where the propositions for p hold true. To answer the second question (Q2), instead of reinitializing the environment after reaching such a state where the propositions for p hold true, we let the Teacher agent sample a task (let’s say Task(p, r)) from the set $\mathcal{X}[p] \setminus \text{DT}$, where \mathcal{X} is the adjacency matrix for the graph, and DT is the set of Discarded Tasks, as defined in Algo. 1. This helps

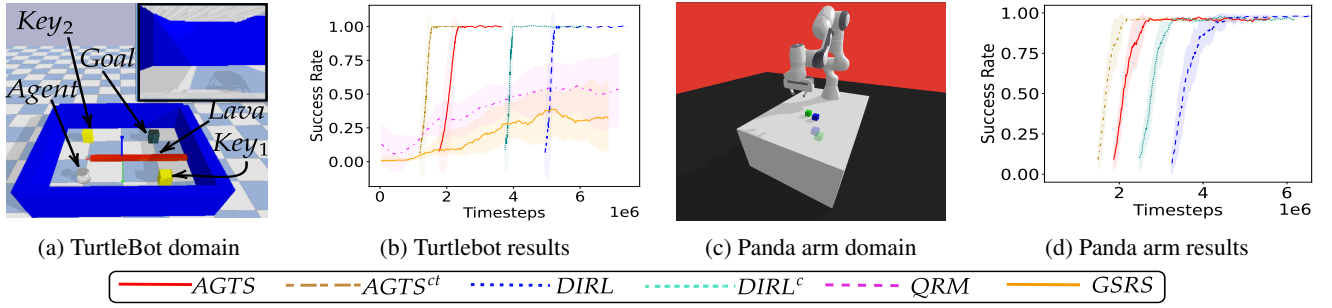


Figure 3: Learning curves (Averaged over 10 trials) for the two robotic domains.

the agent continue its training by attempting to learn a policy $\pi_{(p,r)}$ for Task(p, r) while simultaneously learning a separate policy $\pi_{(q,p)}$ for the task Task(q, p). If the agent fails to satisfy Task(p, r), we reinitialize the environment from $s \sim S_0$. Otherwise, the agent continues its training until it satisfies the high-level objective φ . We call this approach $AGTS^{ct}$ (Detailed algorithm in Appendix A). Results in Table 1 and Fig. 2 demonstrate that this approach improves sample efficiency by reducing the number of interactions required to learn a successful policy for the gridworld task, with a time-to-threshold metric of 3.7×10^5 interactions as compared to the $AGTS$ approach.

5.2 $AGTS$ and $AGTS^{ct}$ - Robotic Domains

To answer (Q3), we test $AGTS$ and $AGTS^{ct}$ on two simulated robotic environments with high interaction cost. The task in Fig. 3a has the following LTL_f objective:

$$\varphi_f^{navigation} := \mathbf{G} \neg Lava \wedge \mathbf{F}((Key_1 | Key_2) \& \mathbf{F}(Goal))$$

In this task, the agent (a simulated TurtleBot) needs to collect any of the keys (yellow blocks) present in a $[3m, 3m]$ continuous environment before reaching the goal position (gray block). At all times, the agent needs to avoid the lava object (red wall) present in the center. The move forward (backward) action causes the robot to move forward (backward) by $0.1m$ and the robot rotates by $\pi/8$ radians with each rotate action. The pick-up and drop actions have effects similar to the gridworld domain. The robotic domain is more complex as objects can be placed at continuous locations. The agent receives an ego-centric image view of the environment (top-right corner of Fig. 3a), which makes the task partially observable in nature and more complex to get a successful policy. The RL agent is described in Sec. 5.1.

The second environment (Fig. 3c) consists of a simulated robotic arm pushing two objects to their target locations (Gallou  dec et al. 2021) with the LTL_f formula:

$$\varphi_f^{manipulation} := \mathbf{F}(p_1) \& \mathbf{F}(p_2)$$

where p_1 and p_2 are the atomic propositions for ‘push-object-1’, ‘push-object-2’. The robot has continuous action parameters for moving the arm and a binary gripper action (close/open). An episode begins with the two objects randomly initialized on the table, and the robotic arm has to

push these two objects to its final location. The agent receives its current end-effector pose, positions and velocities of the two objects, and the desired goal position for the two objects. For this task, we use the Deep Deterministic Policy Gradient with Hindsight Experience Replay (DDPG-HER) (Andrychowicz et al. 2017) as our RL algorithm. DDPG-HER is implemented using the OpenAI baselines (Dhariwal et al. 2017). Both the robotic domains were modeled using PyBullet (Coumans and Bai 2021), and the reward structure for both the RL agents was sparse, similar to the one described in Sec. 5.1. The learning curves for the TurtleBot domain (Fig. 3b) and the Panda arm domain (Fig. 3d) (averaged over 10 trials) are shown in Fig. 3b and Fig. 3d respectively. For both domains, $AGTS$ outperforms all the baselines in terms of learning speed. $AGTS^{ct}$ further speeds-up learning for both the robotic domains. The time-to-threshold between $AGTS$ and the best performing baseline (our custom implementation) $DIRM^c$, is 2×10^6 for the TurtleBot domain and 5×10^5 for the Panda arm domain.

6 Conclusion

We proposed $AGTS$, a framework for dynamic task sampling for RL agents using the high-level LTL_f objective coupled with the Teacher-Student learning strategy. Through experiments, we demonstrated that $AGTS$ accelerates learning, converging to a desired success rate quicker as compared to other curriculum learning and automaton-guided RL baselines. $AGTS^{ct}$ further improves sample efficiency by continuing exploration on a new sub-task once a goal state for a sub-task is reached. We also evaluate our approach on long-horizon complex robotic tasks where the state space is large and the actions are continuous. $AGTS$ reduces training time without relying on human-guided dense reward function, accelerating learning when the high-level objective is available

Limitations & Future Work: In certain cases, the LTL_f objective can be completely novel and/or generating the labeling function can be infeasible. Our future plans involve expanding our framework to scenarios where obtaining a precise LTL_f specification is challenging. As an extension, we would like to explore biasing away from sub-tasks rather than completely discarding them once the target node is reached, so in the limit, optimal policies can be obtained. We would also like to explore complex robotic and multi-agent scenarios with elaborate LTL_f objectives.

References

- Agrawal, S.; and Goyal, N. 2012. Analysis of thompson sampling for the multi-armed bandit problem. In *Conference on learning theory*, 39–1. JMLR Workshop and Conference Proceedings.
- Alur, R.; Bansal, S.; Bastani, O.; and Jothimurugan, K. 2022. A framework for transforming specifications in reinforcement learning. In *Principles of Systems Design: Essays Dedicated to Thomas A. Henzinger on the Occasion of His 60th Birthday*, 604–624. Springer.
- Andrychowicz, M.; Wolski, F.; Ray, A.; Schneider, J.; Fong, R.; Welinder, P.; McGrew, B.; Tobin, J.; Pieter Abbeel, O.; and Zaremba, W. 2017. Hindsight experience replay. *Advances in neural information processing systems*, 30.
- Auer, P.; Cesa-Bianchi, N.; and Fischer, P. 2002. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47: 235–256.
- Bozkurt, A. K.; Wang, Y.; Zavlanos, M. M.; and Pajic, M. 2020. Control synthesis from linear temporal logic specifications using model-free reinforcement learning. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, 10349–10355. IEEE.
- Cai, M.; Aasi, E.; Belta, C.; and Vasile, C.-I. 2023. Overcoming Exploration: Deep Reinforcement Learning for Continuous Control in Cluttered Environments From Temporal Logic Specifications. *IEEE Robotics and Automation Letters*, 8(4): 2158–2165.
- Camacho, A.; Chen, O.; Sanner, S.; and McIlraith, S. A. 2018. Non-Markovian rewards expressed in LTL: Guiding search via reward shaping (extended version). In *GoalsRL, a workshop collocated with ICML/IJCAI/AAMAS*.
- Chevalier-Boisvert, M.; Willems, L.; and Pal, S. 2018. Minimalistic Gridworld Environment for Gymnasium.
- Coumans, E.; and Bai, Y. 2021. PyBullet, a Python module for physics simulation for games, robotics and machine learning. <http://pybullet.org>. Accessed: 2023-04-02.
- De Giacomo, G.; Iocchi, L.; Favorito, M.; and Patrizi, F. 2019. Foundations for restraining bolts: Reinforcement learning with LTLf/LDLf restraining specifications. In *Intl. Conf. on Automated Planning and Scheduling*, volume 29.
- De Giacomo, G.; and Vardi, M. Y. 2013. Linear temporal logic and linear dynamic logic on finite traces. In *IJCAI'13 Proc. of the Twenty-Third Intl. joint Conf. on Artificial Intelligence*, 854–860. Association for Computing Machinery.
- Dhariwal, P.; Hesse, C.; Klimov, O.; Nichol, A.; Plappert, M.; Radford, A.; Schulman, J.; Sidor, S.; Wu, Y.; and Zhokhov, P. 2017. OpenAI Baselines. <https://github.com/openai/baselines>.
- Diuk, C.; Cohen, A.; and Littman, M. L. 2008. An object-oriented representation for efficient reinforcement learning. In *25th Intl. Conf. on Machine learning*, 240–247.
- Gallouédec, Q.; Cazin, N.; Dellandréa, E.; and Chen, L. 2021. panda-gym: Open-Source Goal-Conditioned Environments for Robotic Learning. *4th Robot Learning Workshop: Self-Supervised and Lifelong Learning at NeurIPS*.
- Gao, Y.; and Wu, L. 2021. Efficiently mastering the game of nogo with deep reinforcement learning supported by domain knowledge. *Electronics*, 10(13): 1533.
- Grzes, M. 2017. Reward shaping in episodic reinforcement learning.
- Hammond, L.; Abate, A.; Gutierrez, J.; and Wooldridge, M. 2021. Multi-agent reinforcement learning with temporal logic specifications. *arXiv preprint arXiv:2102.00582*.
- Icarte, R. T.; Klassen, T.; Valenzano, R.; and McIlraith, S. 2018. Using reward machines for high-level task specification and decomposition in reinforcement learning. In *Intl. Conf. on Machine Learning*, 2107–2116.
- Icarte, R. T.; Klassen, T. Q.; Valenzano, R.; and McIlraith, S. A. 2022. Reward machines: Exploiting reward function structure in reinforcement learning. *Journal of Artificial Intelligence Research*, 73: 173–208.
- Jothimurugan, K.; Bansal, S.; Bastani, O.; and Alur, R. 2021. Compositional reinforcement learning from logical specifications. *Neural Information Processing Systems*.
- Kaelbling, L. P.; Littman, M. L.; and Cassandra, A. R. 1998. Planning and acting in partially observable stochastic domains. *Artificial Intelligence*, 101(1): 99–134.
- Kim, T. K. 2015. T test as a parametric statistic. *Korean journal of anesthesiology*, 68(6): 540–546.
- Lattimore, T.; Hutter, M.; and Sunehag, P. 2013. The sample-complexity of general reinforcement learning. In *International Conference on Machine Learning*, 28–36. PMLR.
- Li, X.; Vasile, C.-I.; and Belta, C. 2017. Reinforcement learning with temporal logic rewards. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 3834–3839. IEEE.
- Matiisen, T.; Oliver, A.; Cohen, T.; and Schulman, J. 2020. Teacher-Student Curriculum Learning. *IEEE Trans. Neural Networks Learn. Syst.*, 31(9): 3732–3740.
- Moore, E. F. 1959. The shortest path through a maze. In *Proc. of the International Symposium on the Theory of Switching*, 285–292. Harvard University Press.
- Narvekar, S.; Peng, B.; Leonetti, M.; Sinapov, J.; Taylor, M. E.; and Stone, P. 2020. Curriculum Learning for Reinforcement Learning Domains: A Framework and Survey. *JMLR*, 21: 1–50.
- Nguyen, H.; and La, H. 2019. Review of deep reinforcement learning for robot manipulation. In *2019 Third IEEE International Conference on Robotic Computing (IRC)*, 590–595. IEEE.
- Oudeyer, P.-Y.; and Kaplan, F. 2009. What is intrinsic motivation? A typology of computational approaches. *Frontiers in neurorobotics*, 6.
- Schulman, J.; Wolski, F.; Dhariwal, P.; Radford, A.; and Klimov, O. 2017. Proximal Policy Optimization Algorithms. *CoRR*.
- Shukla, Y.; Kulkarni, A.; Wright, R.; Velasquez, A.; and Sinapov, J. 2023. Automaton-Guided Curriculum Generation for Reinforcement Learning Agents. In *Proceedings of*

the 33rd International Conference on Automated Planning and Scheduling.

Shukla, Y.; Thierauf, C.; Hosseini, R.; Tatiya, G.; and Sinapov, J. 2022. ACuTE: Automatic Curriculum Transfer from Simple to Complex Environments. In *21st Intl. Conf. on Autonomous Agents and Multiagent Systems*, 1192–1200.

Szepesvári, C. 2004. Shortest path discovery problems: A framework, algorithms and experimental results. In *AAAI*, 550–555.

Taylor, M. E.; and Stone, P. 2009. Transfer learning for reinforcement learning domains: A survey. *JMLR*, 10(7).

Toro Icarte, R.; Klassen, T. Q.; Valenzano, R.; and McIlraith, S. A. 2018. Teaching multiple tasks to an RL agent using LTL. In *Autonomous Agents and MultiAgent Systems*.

Velasquez, A.; Bissey, B.; Barak, L.; Beckus, A.; Alkhouri, I.; Melcer, D.; and Atia, G. 2021. Dynamic automaton-guided reward shaping for monte carlo tree search. In *Proc. of the AAAI Conf. on Artificial Intelligence*.

Xu, Z.; and Topcu, U. 2019. Transfer of temporal logic formulas in reinforcement learning. In *IJCAI: proceedings of the conference*, volume 28, 4010. NIH Public Access.