

Task-Motion Planning with Reinforcement Learning for Adaptable Mobile Service Robots

Yuqian Jiang¹, Fangkai Yang², Shiqi Zhang³, and Peter Stone¹

Abstract—Task-motion planning (TMP) addresses the problem of efficiently generating executable and low-cost task plans in a discrete space such that the (initially unknown) action costs are determined by motion plans in a corresponding continuous space. A task-motion plan for a mobile service robot that behaves in a highly dynamic domain can be sensitive to domain uncertainty and changes, leading to suboptimal behaviors or execution failures. In this paper, we propose a novel framework, *TMP-RL*, which is an integration of TMP and reinforcement learning (RL), to solve the problem of robust TMP in dynamic and uncertain domains. The robot first generates a low-cost, feasible task-motion plan by iteratively planning in the discrete space and updating relevant action costs evaluated by the motion planner in continuous space. During execution, the robot learns via model-free RL to further improve its task-motion plans. RL enables adaptability to the current domain, but can be costly with regards to experience; using TMP, which does not rely on experience, can jump-start the learning process before executing in the real world. *TMP-RL* is evaluated in a mobile service robot domain where the robot navigates in an office area, showing significantly improved adaptability to unseen domain dynamics over TMP and task planning (TP)-RL methods.

I. INTRODUCTION

Building mobile robots that behave intelligently in real environments is one of the central problems of robotics and artificial intelligence. Future service robots are expected to take general requests such as “deliver coffee to Alice”. To achieve a goal like this, the integration between high-level task planning and low-level motion planning, also known as *task-motion planning* (TMP), has been widely studied for robot manipulators [1], [2], [3] and navigation tasks in service robot [4] or self-driving cars [5]. TMP algorithms typically consist of a task planner that generates high-level task sequences in an abstract discrete space, possibly using an AI planning approach [6], and a motion planner that expands each task using algorithms such as Probabilistic Random Map [7] or Rapidly-exploring Random Trees [8] to generate a collision-free trajectory based on the current status of the environment. This hierarchical approach reduces the complexity of long-horizon motion planning by improving plan feasibility, quality, and scalability.

Despite the progress made on generating feasible and quality *offline plans*, during *execution*, a robot can still face domain uncertainties and changes that are not available at modeling time. This challenge is particularly pervasive for

mobile service robots that cohabit with human and serve human requests [9], [10]. Mobile service robots usually have to navigate in building-wide areas, whose environmental dynamics involve, among many things, crowds of people, changing lighting conditions, and changed furniture layout, which are not practical for the motion planner to accurately model. Such dynamic changes may invalidate task-motion plans, leading to suboptimal behaviors and execution failures. Continually learning from execution experience and adapting to the changing domain is therefore crucial for mobile service robots to achieve long-term autonomy. To this end, reinforcement learning (RL) has been used to build highly-adaptive autonomous agents [11] and improve symbolic plan robustness and adaptability [12], [13], [14] in various simulation domains, becoming an attractive approach to enable learning adaptive task-motion plans for mobile service robots.

Aiming to *improve adaptability of task-motion plans for mobile service robots*, in this paper, we propose to integrate TMP with RL such that the robot can constantly generate feasible, high-quality task-motion plans and rapidly learn from execution experience to adapt to domain changes. Inspired by PETLON, a recent task-level-optimal TMP algorithm [4], and PEORL, a state-of-the-art task planning-RL framework [12], our approach features *two nested planning-reinforcement learning loops*:

- The *inner loop* is a complete TMP algorithm, where a symbolic plan is generated and each symbolic action is evaluated by the motion planner. Iterative learning and plan improvement is performed on rewards derived from motion plan costs.
- The *outer loop* is for learning to generate an optimal task-motion plan to accommodate domain uncertainty, change, and extra reward information. Each task-motion plan generated by the inner loop is executed in the outer loop to learn from environmental rewards. The inner loop then uses the learned values to generate an improved plan in the next episode. When the outer loop terminates, the robot has learned a task-motion plan that has adapted to the observed domain changes.

In the framework above, the inner loop generates a high quality task-motion plan based on its own discrete and continuous models, leading to a jump-start of plan quality. The outer-loop helps the task-motion planner to fine-tune the plans by learning from the environment, improving the adaptability of TMP facing domain uncertainty and change. The duality between the inner and outer loop allows a

¹Department of Computer Science, University of Texas at Austin, Austin, TX, USA

² NVIDIA Corporation, Redmond, WA, USA

³Department of Computer Science, SUNY-Binghamton, Binghamton, NY, USA

seamless integration of TMP with RL such that motion planning in a continuous model and reinforcement learning from the real execution experience can jointly contribute to improving TMP.

Our approach is generic in the sense that a variety of task planning, motion planning, and reinforcement learning approaches can be used. In this paper, we instantiate our approach using the same planning and learning technologies used in PEORL [12]: action language \mathcal{BC} [15] and answer set solver CLINGO for task planning, and R-learning [16] for reinforcement learning. We have evaluated the approach in an office environment using a real robot and the Gazebo simulator [17]. Compared to PETLON and PEORL, TMP-RL demonstrates superior adaptability to environmental uncertainties; it achieves better task performance than PETLON, and faster convergence than PEORL.

II. RELATED WORK

Existing research on TMP algorithms have various foci, such as ensuring symbolic actions’ feasibility via motion planning [2], integrated symbolic planning under uncertainty and motion planning [18], and leveraging symbolic search heuristics in motion planning space [5], [3]. Recently proposed PETLON [4], which uses sampling-based probabilistic motion planning methods to evaluate costs of task-level actions is most similar to the inner loop of our work, but in our work, we use RL to learn rewards derived from real action costs, whereas PETLON is purely a planning method. While generating feasible, low-cost task-motion plans is the major focus of existing work on TMP, to the best of our knowledge, mixing task-motion planning and learning from execution to accommodate domain uncertainty and change for long range navigation tasks in mobile robots has not been investigated before.

The integration of symbolic planning with reinforcement learning has been studied in a variety of approaches [19], [20], [21]. Recent approaches such as PEORL [12] and SDRL [13] utilize closed-loop communication between planning and learning: an optimal symbolic plan is obtained from a mutually beneficial, iterative process of planning and learning. These approaches are mainly evaluated in simulation domains, but an integrated robot system is usually equipped with motion planners that can be used to evaluate the outcomes of task plans before execution, calling for an integration of TMP with RL. To the best of our knowledge, our work is the first to apply reinforcement learning to improve adaptability of task-motion plans for service robots, where task planning, motion planning, execution, and learning are handled in a unified framework.

III. PRELIMINARIES

Symbolic Planning. An *action description* \mathbb{D} in the language \mathcal{BC} [15] includes *fluent constants* that represent the properties of the world and *action constants*. A fluent atom is an expression of the form $f = v$, where f is a fluent constant and v is an element of its domain. An action description is a finite set of *causal laws* that describe

how fluent atoms are related with each other in a single time step, or how their values are changed from one step to another, possibly by executing actions. For instance, $(A \text{ if } A_1, \dots, A_m)$ is a *static law* that states that at a time step, if A_1, \dots, A_m holds then A is true. Another static law (**default** $f = v$) states that by default, the value of f equals v at any time step. $(a \text{ causes } A_0 \text{ if } A_1, \dots, A_m)$ is a *dynamic law*, stating that at any time step, if A_1, \dots, A_m holds, by executing action a , A_0 holds in the next step. (**nonexecutable** $a \text{ if } A_1, \dots, A_m$) states that at any step, if A_1, \dots, A_m holds, action a is not executable.

A *state* s is a complete set of fluent atoms, and a transition is a tuple $\langle s_1, a, s_2 \rangle$ where s_1, s_2 are states and a is a (possibly empty) set of actions. Let \mathbb{I} be the initial state and \mathbb{G} be goal state. The triple $(\mathbb{I}, \mathbb{G}, \mathbb{D})$ is called a planning problem. A plan can be computed using answer set solver such as CLINGO. Throughout the paper, we use Π to denote both the plan and the transition path by following the plan. Automated planning can be achieved by an answer set solver.

Motion Planning. A configuration space includes a set of all possible, potentially high-dimensional, configurations, where a configuration describes a possible pose of the robot. In this work, we consider a mobile robot that moves in 2D spaces, where we directly search in the 2D workspace (instead of higher-dimensional configuration space). A motion planning problem can be specified by an initial position x^{init} and a goal set X^{goal} . The 2D space is represented as a region in Cartesian space such that the position and orientation of the robot can be uniquely represented as a *pose* (x, θ) . Some parts of the space are designated as free space, and the rest is designated as obstacle.

The motion planning problem is solved by the motion planner \mathcal{P}^{geo} to compute a collision-free trajectory ξ^* (connecting \mathbf{x}^{init} and a pose $\mathbf{x}^{goal} \in \mathbf{X}^{goal}$ taking into account any motion constraints on the part of the robot) with minimal trajectory length $Len(\xi) = L$. We use Ξ to represent the trajectory set that includes all satisfactory trajectories. The *optimal* trajectory is $\xi^* = \operatorname{argmin}_{\xi \in \Xi} Len(\xi)$, where $\xi(0) = \mathbf{x}^{init}$ and $\xi(L) = \mathbf{x}^{goal} \in \mathbf{X}^{goal}$. In particular, we use `global_planner`, an off-the-shelf package from the Robot Operating System (ROS) [22] community for motion planning, which generates trajectories using gradient descent together with standard A^* and Dijkstra’s search.

R-learning for Finite Horizon Problems. A Markov Decision Process (MDP) is defined as a tuple $(\mathcal{S}, \mathcal{A}, P_{ss'}^a, r, \gamma)$, where \mathcal{S} and \mathcal{A} are the sets of symbols denoting states and actions, the transition kernel $P_{ss'}^a$ specifies the probability of transition from state $s \in \mathcal{S}$ to state $s' \in \mathcal{S}$ by taking action $a \in \mathcal{A}$, $r(s, a) : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$ is a reward function bounded by r_{\max} , and $0 \leq \gamma < 1$ is a discount factor. A solution to an MDP is a policy $\pi : \mathcal{S} \mapsto \mathcal{A}$ that maps a state to an action. Model-free RL concerns on learning a near-optimal policy by executing actions and observing transitions and rewards.

To evaluate a policy π , R-learning [16] applies to the expected un-discounted sum of reward for finite hori-

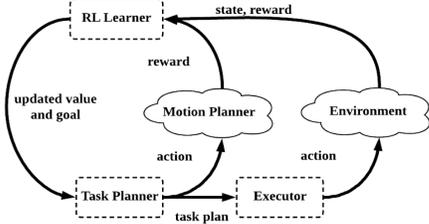


Fig. 1: An illustration of our TMP-RL framework

zon problems. Define $J_{\text{avg}}^{\pi}(s) = \mathbb{E}[\sum_{t=0}^T r_t | s_0 = s]$, and the *gain reward* $\rho^{\pi}(s)$ reaped by policy π from s as $\rho^{\pi}(s) = \lim_{T \rightarrow \infty} \frac{J_{\text{avg}}^{\pi}(s)}{T} = \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}[\sum_{t=0}^T r_t]$. R-learning is a model-free value iteration algorithm that can be used to find the optimal policy for the average reward criterion. At the t -th iteration (s_t, a_t, r_t, s_{t+1}) , the following update is performed:

$$\begin{aligned} R_{t+1}(s_t, a_t) &\stackrel{\leftarrow \alpha_t}{=} r_t - \rho_t(s_t) + \max_a R_t(s_{t+1}, a) \\ \rho_{t+1}(s_t) &\stackrel{\leftarrow \beta_t}{=} r_t + \max_a R_t(s_{t+1}, a) - \max_a R_t(s_t, a), \end{aligned} \quad (1)$$

where α_t, β_t are the learning rates, and $a_{t+1} \stackrel{\leftarrow \alpha}{=} b$ denotes the soft update rule $a_{t+1} = (1 - \alpha)a_t + \alpha b$.

IV. TMP-RL FRAMEWORK

The TMP-RL framework proposed in this paper is shown in Fig. 1. The inner loop consists of a task planner, a motion planner and a reinforcement learner that iteratively performs planning and learning to generate *a feasible and low-cost task-motion plan*. Once the inner loop returns a task-motion plan, it is sent to execution in the outer loop, where the reinforcement learner performs value iteration on the reward derived from execution experience. When the inner loop runs again, it generates an improved task-motion plan based on the learned values. The framework is explained in detail below.

A. Optimal Task Planning Conditioned on Motion Planning

A task planning problem defines the objective of generating a satisfactory plan Π^{τ} , i.e., a sequence of actions given a planning problem $(\mathbb{I}^{\tau}, \mathbb{G}^{\tau}, \mathbb{D}^{\tau})$, where \mathbb{D}^{τ} is a domain independent symbolic formulation given by human experts, \mathbb{I}^{τ} is an initial state and \mathbb{G}^{τ} is a goal state. As in PEORL, \mathbb{D}^{τ} consists of causal laws that formulates preconditions and effects of actions, such as *approach* door D_1 causes the robot beside D_1 if currently the robot is beside D_2 and D_1 is accessible from D_2 :

approach(D_1) **causes** *beside*(D_1) **if** *beside*(D_2), *acc*(D_2, D_1)

and static relationship on fluents, such as symmetry of accessible relationship: *acc*(D_1, D_2) **if** *acc*(D_2, D_1).

A motion planning problem concerns on generating a collision free trajectory $\xi(\mathbb{I}^m, \mathbb{G}^m)$ given a motion planning problem $(\mathbb{D}^m, \mathbb{I}^m, \mathbb{G}^m)$ where \mathbb{D}^m is a motion planning

domain, \mathbb{I}^m is an initial position and \mathbb{G}^m is the goal position, such that the position \mathbb{I}^m is connected with position \mathbb{G}^m .

We use a mapping function $f : X = f(s)$ that maps a symbolic state s into a set of feasible poses X in continuous space, for the motion planning algorithm to sample from. We assume the availability of at least one pose $x \in X$ in each state s , such that the robot is in the free space of \mathbb{D}^m . If it is not the case, the state s is declared infeasible. Given a motion planning domain \mathbb{D}^m and a task plan Π^{τ} for task planning problem $(\mathbb{D}^{\tau}, \mathbb{I}^{\tau}, \mathbb{G}^{\tau})$, a plan refinement of Π^{τ} w.r.t motion planner, denoted as Π^m , is a sequence of collision free trajectories obtained by perform motion planning on each navigation actions, i.e., $\Pi^m = \bigcup_{\langle s, a, s' \rangle \in \Pi^{\tau}} \xi(x, x')$, where $x \in f(s)$, $x' \in f(s')$. The cumulative cost of a task plan Π^{τ} is obtained by the cumulative length of its motion planning refinement Π^m , i.e., $Cost(\Pi^{\tau}) = \sum_{\xi \in \Pi^m} Len(\xi)$. An *optimal task plan conditioned on motion plan* is defined as the task plan Π_o^{τ} such that Π_o^{τ} has the minimal length among all task plans.

B. TMP with Reinforcement Learning

1) *Reward*: Given a symbolic transition $\langle s, a, s' \rangle$ where a can be refined by motion planner, we define a reward function r that is negative and inversely proportional to a distance metric of the motion plan that refines the navigation action a , mapped by a function $R : \mathbb{R}^+ \mapsto \mathbb{R}^-$:

$$r(s, a) = R(Len(\xi(x, x'))) \propto \frac{1}{Len(\xi(x, x'))},$$

where $x \in f(s), x' \in f(s')$. One way to instantiate R is

$$r(s, a) = R(Len(\xi(x, x'))) = -Len(\xi(x, x')).$$

If a motion plan fails for transition $\langle s, a, s' \rangle$, we define $r(s, a) = -\infty$.

2) *Domain Formulation*: We enrich the domain formulation \mathbb{D}^{τ} with the following causal laws formulating the effect of actions on cumulative plan quality:

a causes quality = $C + Z$ **if** $s, \rho(s, a) = Z, quality = C$

where s is a state. The ρ -values are initialized optimistically to the upper-bound of gain reward, which is the reward derived from the L_p metric in the configuration space:

$$\text{default } \rho(s, a) = \max_{x, x'} R(\|x - x'\|_p)$$

where $x \in f(s), x' \in f(s'), x \neq x'$, for $\langle s, a, s' \rangle$ in $T(\mathbb{D}^{\tau})$, $p \in \mathbb{R}^+$, and L_p metric stands for $\|x - x'\|_p = \left(\sum_{i=1}^n |x_i - x'_i|^p \right)^{-p}$.

3) *Planning Goal*: At any episode t , planning goal \mathbb{G}_t^{τ} contains a regular logical constraint describing the goal condition plus a linear constraint of the form

$$quality \geq quality(\Pi_t^{\tau}) \quad (2)$$

where $quality(\Pi_t^{\tau}) = \sum_{\langle s, a, s' \rangle \in \Pi_t^{\tau}} \rho(s, a)$ for some task plan Π_t^{τ} . In the planning - learning loop, the linear

constraint guides the planner to generate a plan with cumulative quality higher than a previous one, measured by learned ρ -values, leading to the iterative process of plan improvement based on reinforcement learning.

4) *Algorithm for TMP*: Algorithm 1 describes our inner loop of task-motion planning. The input to the algorithm includes a motion planning domain and a task planning problem. q_0 is initialized to be $-\infty$, and $P_0 = \emptyset$. The algorithm first generates a task plan (Line 4). Then it iterates on each symbolic transition in the plan, and for each navigation action, it obtains the initial and goal poses in 2D domain (Line 12), generates motion plan (Line 13) and returns reward (Line 14). Value iteration of R-learning is performed with the reward (Line 15). At the end of this process, plan quality is computed using the learned ρ values (Line 17), and it is used as the new constraint in the planning goal (Line 18), setting a baseline for the planner in the next iteration. The learned ρ values are also updated in the symbolic formulation (Line 19). When the algorithm terminates, it outputs a task plan that cannot be further improved w.r.t the motion planner.

Algorithm 1 Task-Motion Planning

Require: $(\mathbb{I}^\tau, \mathbb{G}^\tau, \mathbb{D}^\tau, f, \mathbb{D}^m, q_0, P_0)$ where $quality > q_0 \in G^\tau$, and an exploration probability ϵ

- 1: $t \leftarrow 0$
- 2: **while** $t < +\infty$ **do**
- 3: $\Pi^* \leftarrow \Pi_t^\tau$
- 4: obtain a plan $\Pi_t^\tau \leftarrow Plan(\mathbb{I}^\tau, \mathbb{G}^\tau, \mathbb{D}^\tau \cup P_t)$
- 5: **if** $\Pi_t^\tau = \emptyset$ **then**
- 6: **return** Π^*
- 7: **end if**
- 8: **for** symbolic transition $\langle s, a, s' \rangle \in \Pi_t^\tau$ **do**
- 9: **if** a cannot be refined by motion planner **then**
- 10: continue
- 11: **end if**
- 12: obtain initial pose $x = f(s)$ and goal pose $x' = f(s')$
- 13: generate motion plan $\xi(x, x')$
- 14: calculate reward $r(s, a) = R(Len(\xi(x, x')))$
- 15: update $R(s, a)$ and $\rho^a(s)$ using (1).
- 16: **end for**
- 17: calculate quality of Π_t^τ by (2).
- 18: update planning goal $G \leftarrow (quality > quality_t(\Pi_t^\tau))$.
- 19: update facts $P_t \leftarrow \{\rho(s, a) = z : \langle s, a, s' \rangle \in \Pi_t^\tau, \rho_t^a(s) = z\}$
- 20: $t \leftarrow t + 1$
- 21: **end while**

C. TMP Execution and Learning

Once a task-motion plan is generated, it is sent for execution, which goes to the outer loop of learning from real execution experience, where Algorithm 1 becomes the planning subroutine (Line 4) in Algorithm 2. In Algorithm 2, each action is executed in the environment, and the true reward is obtained (Line 9). The value iteration performed on the true reward received during execution further rewrites the value learned through motion planner and feed back into the TMP algorithm (Line 14) to iteratively generate a task-motion plan that is adaptable to domain change.

The difference between Algorithm 2 and Algorithm 1, is in Line 4: Algorithm 1 makes a task planning call and

Algorithm 2 Task-Motion Planning and Learning

Require: $(\mathbb{I}^\tau, \mathbb{G}^\tau, \mathbb{D}^\tau, f, \mathbb{D}^m)$ where $quality > 0 \in G^\tau$, and an exploration probability ϵ

- 1: $P_0 \leftarrow \emptyset, \Pi_0^\tau \leftarrow \emptyset, q_0 = -\infty, t = 0$
- 2: **while** $t < +\infty$ **do**
- 3: $\Pi^* \leftarrow \Pi_t^\tau$
- 4: obtain a task-motion plan by calling Algorithm 1 $\Pi_t^\tau \leftarrow TMP(\mathbb{I}^\tau, \mathbb{G}^\tau, \mathbb{D}^\tau, f, \mathbb{D}^m, q_t, P_t)$.
- 5: **if** $\Pi_t^\tau = \emptyset$ **then**
- 6: **return** Π^*
- 7: **end if**
- 8: **for** symbolic transition $\langle s, a, s' \rangle \in \Pi_t^\tau$ **do**
- 9: execute a and obtain true reward $r(s, a)$.
- 10: update $R(s, a)$ and $\rho^a(s)$ using (1).
- 11: **end for**
- 12: calculate quality of Π_t^τ by (2).
- 13: update plan quality $q_t \leftarrow quality_t(\Pi_t^\tau)$.
- 14: update facts $P_t \leftarrow \{\rho(s, a) = z : \langle s, a, s' \rangle \in \Pi_t^\tau, \rho_t^a(s) = z\}$
- 15: $t \leftarrow t + 1$
- 16: **end while**

Algorithm 2 makes a TMP call. Such duality brings a unification of refining task plans through motion planner and through learning from the environment: the quality of task plans are learned in the same framework and the learned values are propagated back so that motion planners and execution experience can jointly improve the task plan.

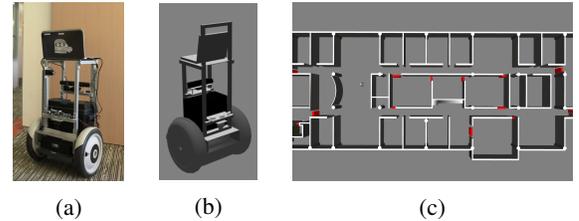


Fig. 2: (a) a BWIBot (b) a simulated BWIBot (c) simulation environment.

V. EXPERIMENTS

The TMP-RL framework is fully implemented and tested on a real service robot, BWIBot [10], as shown in the accompanying video¹. Additional experiments are conducted using a simulation closely matching the real platform and environment (Fig. 2). We compare the performance of the proposed TMP-RL algorithm (Algorithm 2) with PETLON [4], a TMP algorithm, and PEORL [12], a TP-RL approach. We measure the actual reward the robot receives in each episode by executing the plan generated by each algorithm, and compare the learning curves. Answer set solver Clingo 5.3 is used for task (symbolic) planning. Path planning is implemented using the navigation stack of Robot Operating System (ROS) [22]. In TMP and TMP-RL implementations, the global path planner is called to generate a trajectory for each navigation action, and the motion costs are estimated by the sum of distances between waypoints on the trajectory.

¹<https://youtu.be/EyoqrpO3Qkk>

Plan	Task Plan Cost	Motion Plan Cost	Average Execution Time
(1)	3	60.8	80.6
(2)	9	45.5	126.9
(3)	3	53.1	116.6

TABLE I: Plan costs at different levels of abstraction.

The plan quality constraints are implemented for T(M)P-RL algorithms. Learning of the ρ values is implemented using Equation (1) with learning rates $\alpha = 0.1, \beta = 0.5$. The default ρ -values of *approach* actions are implemented as the Euclidean distance between the target location and the landmark closest to the robot. The default reward of an *open_door* action is -3. A state-action pair without an evaluated or default ρ value is assumed to have reward -1.

In this experiment, the robot needs to go in a room where its service is requested. The robot starts near a landmark in an open space and the robot’s end position can be anywhere in the target room. The reward is defined as the negative of the execution time. The room has three initially closed doors that are available for entrance. The task planner determines which entrance the robot will use. Fig. 3a shows the experiment set-up and three competitive task plans. With 30 regions and 12 doors in the domain, many other feasible plans may be generated by task planner, and some plans involve significant detour. Table I shows the task plan length, motion plan length, and average execution time of the competitive plans. Among the three plans, plan (2) features the shortest navigation distance, but it takes 9 actions and requires crossing 3 doors. Plan (3) has 3 actions and the second shortest navigation distance. The duration of executing an *open_door* action is sampled from a normal distribution with a standard deviation of (10 seconds). Opening the bottom door takes 60 seconds on average, while the mean open time is 20 seconds for other doors. Therefore, plan (1) using the top door has the lowest expected execution time. This example shows one situation where all three levels of capability are required to efficiently find the optimal real-world plan.

Evaluation of TMP, TP-RL, TMP-RL. In this evaluation, we use PETLON (TMP) and PEORL (TP-RL) for comparison. For every approach we conducted 50 runs with 40 episodes in each run. The variability among the trials are caused by noisy action costs of navigation and opening doors. For RL-based methods, they can generate different plans depending on experiences in previous episodes.

Fig. 3b plots the learning curves for reward received in 40 episodes, averaged over 50 runs with the shaded regions representing one standard deviation from the mean. Fig. 3c shows for each approach, the average number of episodes that the three competitive plans and other feasible plans are executed. Equipped with the reinforcement learner to refine their plans, TMP-RL and TP-RL converge to the practically optimal plan, but TMP-RL converges significantly faster. Using motion plan costs in task plan evaluation ensures that TMP-RL makes steady improvements after the first two episodes, while TP-RL learns everything from executing the plans in the environment. TP-RL has low variances in the

first two episodes, because the task planner first selects the plans with the smallest number of actions (plans (1) and (3) in Fig. 3a), but much higher variances afterwards. As shown in Fig. 3c, TP-RL had to execute many task plans that are logically valid but significantly worse in quality directly in the environment, which is expensive and potentially dangerous for real robots. TMP executes the plan with the shortest navigation distance in every episode (plan (2)), without learning any information from execution.

Evaluation of TMP-RL in Multiple Tasks. In long-term deployments, the robot can be asked to achieve the same end goal from different starting positions. Since the initial states are different, the task planner and motion planner have to solve them as different problems, but TMP-RL can leverage the learned ρ -values to speed up exploration in later tasks. In order to demonstrate TMP-RL’s ability to generalize learned values to different scenarios, we extend the previous experiment with two more tasks, each with a different starting position of the robot (shown in Fig. 4a).

In this scenario, we compared continuously running TMP-RL for all three tasks against using TMP-RL to learn the second and the third tasks from scratch. In the former setting, the robot explored the first task for 15 episodes, and then switched to the second position and the third position while keeping the learned values. In the latter setting, the robot started at episode 15 and performed the second task, or started at episode 30 and performed the third task. Fig. 4b presents the learning curves averaged over 40 runs in these three settings, showing that learning the first task leads to faster, lower variance learning in the later tasks, in comparison with learning from scratch, indicating that the learned values can be transferred to accelerate learning other tasks.

VI. CONCLUSION

We introduce a novel TMP-RL framework integrating task-motion planning (TMP) and reinforcement learning (RL) for adaptable robot sequential decision making. The framework mixes task planning, motion planning, and reinforcement learning in a closed loop with iterative improvements on plan quality over the course of execution. Experiments show that TMP-RL combines the strengths of high quality offline plan generation of TMP and adaptability of RL, achieving practical learning time in a real service robot domain. Future work includes using TMP-RL to improve the long-term performance of service robots in a variety of tasks, as well as extending the framework to multi-robot scenarios.

ACKNOWLEDGEMENT

This work has taken place partly in the Learning Agents Research Group (LARG) at UT Austin. LARG research is supported in part by NSF (IIS-1637736, CPS-1739964, IIS-1724157), ONR (N00014-18-2243), FLI (RFP2-000), ARL, DARPA, and Lockheed Martin. Peter Stone serves on the Board of Directors of Cogitai, Inc. The terms of this arrangement have been reviewed and approved by the University of Texas at Austin in accordance with its policy on objectivity in research.

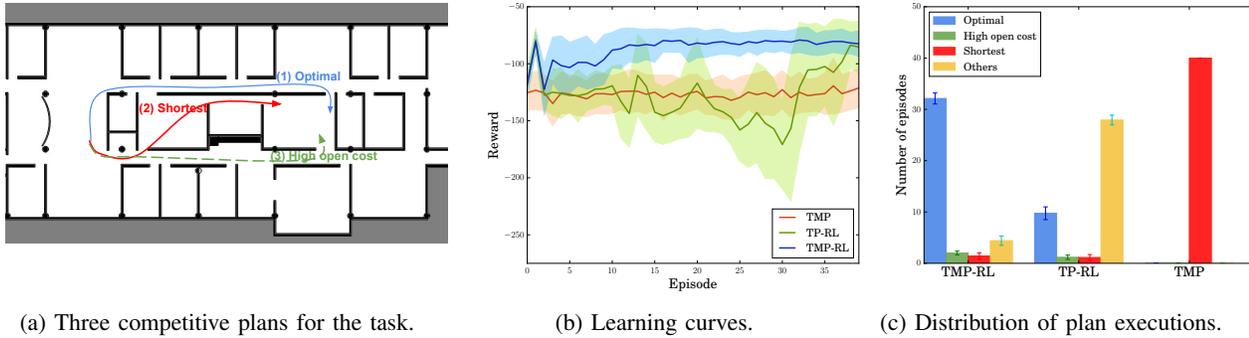
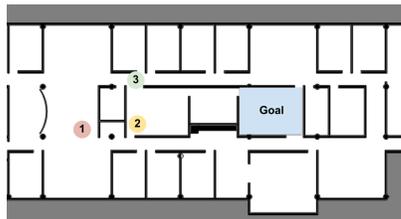
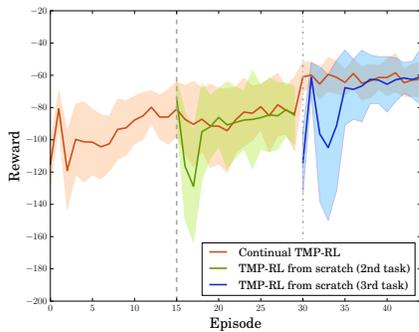


Fig. 3: Comparisons between TMP-RL and baselines of TMP and TP-RL in the task of room-to-room navigation.



(a) Extended experiment with three different initial states.



(b) TMP-RL with continuous transfer vs. learning from scratch.

Fig. 4: Evaluation of TMP-RL in multiple tasks

REFERENCES

- [1] E. Erdem, K. Haspalmutgil, C. Palaz, V. Patoglu, and T. Uras, "Combining high-level causal reasoning with low-level geometric reasoning and motion planning for robotic manipulation," in *Robotics and Automation (ICRA), 2011 IEEE International Conference on*. IEEE, 2011, pp. 4575–4581.
- [2] S. Srivastava, E. Fang, L. Riano, R. Chitnis, S. Russell, and P. Abbeel, "Combined task and motion planning through an extensible planner-independent interface layer," in *Robotics and Automation (ICRA), 2014 IEEE International Conference on*. IEEE, 2014, pp. 639–646.
- [3] C. R. Garrett, T. Lozano-Perez, and L. P. Kaelbling, "Ffrob: Leveraging symbolic planning for efficient task and motion planning," *The International Journal of Robotics Research*, vol. 37, no. 1, pp. 104–136, 2018.
- [4] S.-Y. Lo, S. Zhang, and P. Stone, "Petlon: Planning efficiently for task-level-optimal navigation," in *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*. International Foundation for Autonomous Agents and Multiagent Systems, 2018, pp. 220–228.
- [5] C. Chen, A. Gaschler, M. Rickert, and A. Knoll, "Task planning for highly automated driving," in *Intelligent Vehicles Symposium*, 2015, pp. 940–945.
- [6] A. Cimatti, M. Pistore, and P. Traverso, "Automated planning," in *Handbook of Knowledge Representation*, F. van Harmelen, V. Lifschitz, and B. Porter, Eds. Elsevier, 2008.
- [7] L. Kayraki, P. Svestka, J. Latombe, and M. Overmars, "Probabilistic roadmaps for path planning in high-dimensional configurations space," *Proc IEEE Trans Robot Autom*, vol. 12, no. 4, pp. 566–80, 1996.
- [8] S. M. LaValle, "Rapidly-exploring random trees: A new tool for path planning," 1998.
- [9] M. Veloso, J. Biswas, B. Coltin, and S. Rosenthal, "CoBots: Robust Symbiotic Autonomous Mobile Service Robots," in *Proceedings of IJCAI'15, the International Joint Conference on Artificial Intelligence*, Buenos Aires, Argentina, July 2015.
- [10] P. Khandelwal, S. Zhang, J. Sinapov, M. Leonetti, J. Thomason, F. Yang, I. Gori, M. Svetlik, P. Khante, V. Lifschitz, and P. Stone, "Bwibots: A platform for bridging the gap between ai and human-robot interaction research," *The International Journal of Robotics Research*, vol. 36, no. 5-7, pp. 635–659, 2017.
- [11] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, et al., "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, p. 529, 2015.
- [12] F. Yang, D. Lyu, B. Liu, and S. Gustafson, "Peorl: Integrating symbolic planning and hierarchical reinforcement learning for robust decision-making," in *International Joint Conference of Artificial Intelligence (IJCAI)*, 2018.
- [13] D. Lyu, F. Yang, B. Liu, and S. Gustafson, "Sdrl: Interpretable and data-efficient deep reinforcement learning leveraging symbolic planning," in *AAAI*, 2019.
- [14] —, "A human-centered data-driven planner-actor-critic architecture via logic programming," in *35th International Conference on Logic Programming (ICLP'19)*, 2019.
- [15] J. Lee, V. Lifschitz, and F. Yang, "Action Language BC: A Preliminary Report," in *International Joint Conference on Artificial Intelligence (IJCAI)*, 2013.
- [16] S. Mahadevan, "Average reward reinforcement learning: Foundations, algorithms, and empirical results," *Machine Learning*, vol. 22, pp. 159–195, 1996.
- [17] N. Koenig and A. Howard, "Design and use paradigms for Gazebo, an open-source multi-robot simulator," in *International Conference on Intelligent Robots and Systems (IROS)*, 2004.
- [18] L. P. Kaelbling and T. Lozano-Pérez, "Integrated task and motion planning in belief space," *The International Journal of Robotics Research*, vol. 32, no. 9-10, pp. 1194–1227, 2013.
- [19] R. Parr and S. J. Russell, "Reinforcement learning with hierarchies of machines," in *Advances in neural information processing systems*, 1998, pp. 1043–1049.
- [20] C. Hogg, U. Kuter, and H. Munoz-Avila, "Learning methods to generate good plans: Integrating htn learning and reinforcement learning," in *AAAI*, 2010.
- [21] M. Leonetti, L. Iocchi, and P. Stone, "A synthesis of automated planning and reinforcement learning for efficient, robust decision-making," *Artificial Intelligence*, vol. 241, pp. 103–130, 2016.
- [22] M. Quigley, K. Conley, B. Gerkey, J. Faust, T. Foote, J. Leibs, R. Wheeler, and A. Y. Ng, "Ros: an open-source robot operating system," in *ICRA workshop on open source software*, vol. 3, no. 3.2. Kobe, Japan, 2009, p. 5.