

Avoiding Moving Obstacles with Stochastic Hybrid Dynamics using PEARL: PrEference Appraisal Reinforcement Learning

Aleksandra Faust^{1,2}, Hao-Tien Chiang¹, Nathanael Rackley¹, and Lydia Tapia¹

Abstract—Manual derivation of optimal robot motions for task completion is difficult, especially when a robot is required to balance its actions between opposing preferences. One solution has been proposed to automatically learn near optimal motions with Reinforcement Learning (RL). This has been successful for several tasks including swing-free UAV flight, table tennis, and autonomous driving. However, high-dimensional problems remain a challenge. We address this dimensionality constraint with PrEference Appraisal Reinforcement Learning (PEARL), which solves tasks with opposing preferences for acceleration controlled robots. PEARL projects the high-dimensional continuous robot state space to a low dimensional preference feature space resulting in efficient and adaptable planning. We demonstrate that on a dynamic obstacle avoidance robotic task, a single learning on a much simpler problem performs real-time decision-making for significantly larger, high-dimensional problems working in unbounded continuous states and actions. We trained the agent with 4 static obstacles, while the trained agent avoids up to 900 moving obstacles with complex hybrid stochastic obstacle dynamics in a highly constrained space using only limited information about the environment. We compare these tasks to traditional, often manually tuned solutions for these high-dimensional problems.

I. INTRODUCTION

There are many high-dimensional, motion-based robotic tasks, including multi-robot coordination and control of complex kinematic linkages. These complex robotic problems often require planning high-dimensional motions that complete the task in a timely manner. Motion and trajectory planning identifies a sequence of actions that move the robot in accordance to its dynamics (physical constraints) and the task objectives. Since manually accounting for all possibilities is often infeasible, sampling-based, learning-based, and other intelligent methods are the norm [18]. Reinforcement learning (RL), in particular, has been successful for robotic task learning [15] in several problems such as table tennis [21], swing-free UAV delivery [7], and a self-driving car [12]. However, traditional RL methods do not handle continuous and high-dimensional state spaces well [11].

A primary challenge in these problems is task description. The task goals and constraints that the robot must obey are often unknown or difficult to calculate. For example, consider a simple manipulation task; a robot is required to set a glass on a table without breaking it. We do not know precisely the amount of force that causes the glass to shatter, yet we can describe our preferences: low force

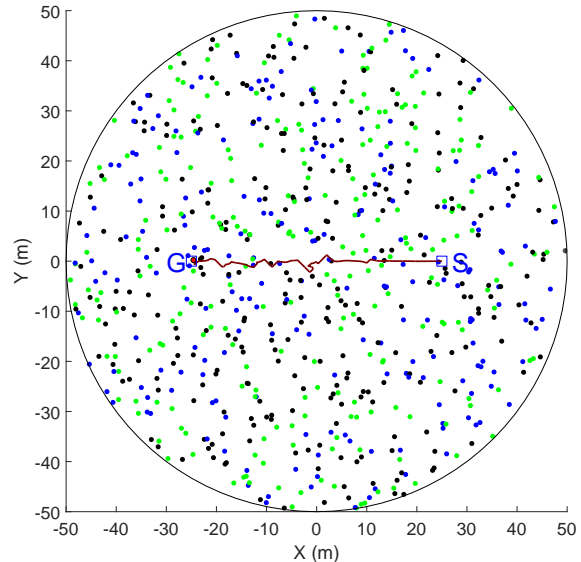


Fig. 1. Dynamic Obstacle Avoidance Task: The robot (red dot) must travel from S (start) to G (goal) without colliding with 900 obstacles moving with hybrid stochastic dynamics. The dark red curve is an example path of the robot. The black/green/blue dots are stochastically moving obstacles with swerve/linear/arc dynamics.

and fast task completion. *Preference reinforcement learning* learns and performs preference balancing tasks (PBTs) with desired qualities (preferences). These tasks cannot be easily demonstrated, but the set of preferences can be described [27].

We present PrEference Appraisal Reinforcement Learning (PEARL) for solving high-dimensional motion-based PBT problems (Figure 2). PEARL trains the planning agent on a small problem, and transfers the learned policy to be used for planning for high-dimensional problems. The key to PEARL is the feature selection method that constructs task-preference features invariant to the robot’s state space dimensionality. Because the method learns and performs the task in the feature space, such transfer is possible. Previously, we empirically showed that, using hand-crafted features, batch RL learns in small spaces and acts on larger problems [8], [7], but did not address when learning transfer is possible and how to do it for an arbitrary problem. This paper formalizes the feature selection and the conditions under which transfer is possible so that it can be applied for classes of PBTs. We include preferences that increase over time as well.

¹Computer Science Department, University of New Mexico, USA, ²Currently with: Google, Inc. {faust@google.com, lewispro@unm.edu, nprackley@gmail.com, tapia@cs.unm.edu}

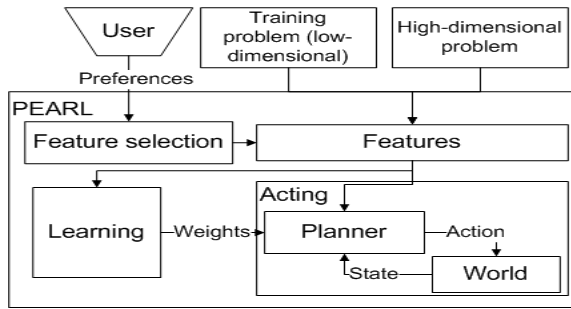


Fig. 2. PrEference Appraisal Reinforcement Learning (PEARL) framework for learning and executing PBT. The user-provided preferences are encoded into polymorphic features. The learning agent appraises preference priorities on a low-dimensional training problem. The planner takes an initial state of a high-dimensional problem and produces actions in a closed feedback loop.

We demonstrate PEARL on the dynamic obstacle avoidance problem, which is critical in several applications such as flight coordination [10], autonomous vehicles and human interacting robots [17]. The case study shows that our method is fast, easy to use, successful for high-dimensional problems, and is a useful tool in a RL-based motion planning toolbox. An agent is required to move to the goal without colliding with hundreds of moving obstacles with hybrid stochastic dynamics. PEARL learns the task with only four static obstacles, and plans trajectories in densely populated environments with up to 900 obstacles. In real-time, at 10 Hz, the planner generates shorter trajectories of higher success rates compared to a traditional Gaussian Artificial Potential Field (APF) obstacle avoidance method, which requires extensive manual tuning [17]. In addition, trajectories generated by PEARL have higher success rates and less running time per planning step compared to a state of the art velocity obstacle based collision avoidance algorithm [9]. In a previous (unpublished) IROS 2015 workshop paper, we showed results with simpler obstacle dynamics [6]. Here, we increase the complexity by considering moving obstacles with hybrid dynamics, which includes three stochastic dynamics modes and a probability to switch modes. The enclosed video submission contains examples of the simulations.

The contributions of this work are: 1) a solution for planning high-dimensional, preference-balancing, motion-based problems (PEARL), and 2) continuous RL solutions to dynamic obstacle avoidance that outperforms state-of-the-art collision avoidance algorithms. These contributions are achieved by Markov decision process formulation for PTBs, and obstacle avoidance tasks, formalizing the preference feature selection method, including preferences that increase over time, and appraising the preferences on small tasks with continuous action RL.

II. RELATED WORK

Reinforcement Learning: Function approximation RL methods typically assume user-provided features [2]. Because they map the entire robot’s state subspace to a single

point, RL is very sensitive to feature selection [2], making feature selection very difficult. Classically, two feature types are used in RL: discretization, and basis functions [28]. Discretization partitions the domain, scaling exponentially with the space dimensionality. Basis functions, such as kernels and radial basis function networks, offer more learning flexibility. These functions, however, can require manual parameter tuning, and the feature number increases MDP formulation for multi-robot systems and dynamic obstacle avoidance tasks exponentially with the state space dimensionality [28]. PEARL proposes a feature selection method that solves a particular class of motion tasks for acceleration-controlled robots. PEARL features exploit task knowledge. Their number is invariant to the problem dimensionality, and the computation time scales polynomially with the state space dimension. In related work, Voronoi decomposition solves high-dimensional manipulation problems by projecting the robot’s configuration space onto a low dimensional task space [23]. The features we propose define a basis in the preference-task space as well. However, PEARL autonomously learns the relationship between the features. Another RL approach solves manipulation tasks with hard [24] and soft [16] constraints. Our tasks, however, do not have known constraints and bounds; they are set up as preferences to guide dynamically feasible trajectory generation.

In summary, PEARL takes task preferences as objectives, and generates features. It appraises the features to come up with weights. Another approach, Weighed Real-Time Heuristic Search (RTHS) [22], uses the weights along the search during planning, and work only in finite spaces. PEARL uses RL, which intrinsically includes weights to balance short and long term gains. In addition, PEARL learns weights between features to find a good balance of priorities among opposing preferences, and works in continuous spaces.

Dynamic Obstacle Avoidance: Planning motions in dynamic environments is challenging because plans must be frequently adjusted due to moving obstacles. Sampling-based methods provide low cost solutions to high-dimensional planning problems in dynamic environments [13], [14], [20]. For example, lazy Probabilistic Roadmaps (PRM) accommodate moving obstacles by rechecking edge validity [13], [14], while [20] pre-builds a Rapidly-exploring Random Tree (RRT) that is modified whenever obstacle changes are detected. Artificial Potential Fields (APF) provide low-cost solutions to dynamic environments by using only local information near the robot [17]. Nevertheless, several parameters impact the performance, such as the relative strength between the repulsive and attractive potentials, and the size of the repulsive potential. In contrast, Velocity Obstacle (VO) based motion planning [9], a state-of-the-art multi-obstacle collision avoidance algorithm, assumes collaborative agents where each agent runs the same VO avoidance method. VO computes the velocity of other nearby agents, and plans an action that reciprocally avoids collision. Unlike VO, our approach makes no assumptions about the obstacles’ behavior, and is appropriate for environments where we do not have control over the other agents. We compare the

proposed method with a modified VO code-base that removes the agent collaboration assumption [4].

III. BACKGROUND

We consider robots as mechanical systems that can be moved using an external force, and model their motion with a special case of nonlinear systems, a *discrete-time control-affine system* [18]. Consider a robot with m degrees of freedom (DoF). If an acceleration $\mathbf{a}(n) \in \mathbb{R}^m$ is applied to the robot's center of mass at time-step n , the new position-velocity vector (state) $\mathbf{s}(n+1) \in \mathbb{R}^{2m}$ is,

$$D: \quad \mathbf{s}(n+1) = \mathbf{f}(\mathbf{s}(n)) + \mathbf{g}(\mathbf{s}(n))\mathbf{a}(n), \quad (1)$$

for some functions \mathbf{f}, \mathbf{g} . A *Markov decision process* (MDP), a tuple (S, A, D, R) with states $S \subset \mathbb{R}^{2m}$ and actions $A \subset \mathbb{R}^m$, that assigns immediate scalar rewards $R: S \rightarrow \mathbb{R}$ to states in S , formulates a task for the system (1) [2]. A solution to a MDP is a control policy $\pi: S \rightarrow A$ that maximizes cumulative discounted reward over the agent's lifetime, value, $V(\mathbf{s}(0)) = \sum_{i=0}^{\infty} \gamma^i R(\mathbf{s}(i))$, where $0 \leq \gamma \leq 1$ is the discount factor. *Approximate value iteration* (AVI) [5], finds a solution to a continuous state MDP by approximating state-value function V with a linear map

$$\hat{V}(\mathbf{s}) = \boldsymbol{\theta}^T \mathbf{F}(\mathbf{s}). \quad (2)$$

AVI takes a feature vector $\mathbf{F}(\mathbf{s})$ and learns weights $\boldsymbol{\theta}$ between the features by sampling the state-space and observing the rewards. It iteratively updates $\boldsymbol{\theta}$ in an expectation-maximization manner.

After parameter learning is completed, batch RL enters a planning phase. The planner takes the value function approximation (2), an initial condition, and generates a trajectory using the closed-loop control with a greedy policy with respect to the state-value approximation,

$$\pi^{\hat{V}}(\mathbf{s}) = \underset{\mathbf{a} \in A}{\operatorname{argmax}} \hat{V}(\mathbf{s}'); \quad (3)$$

where state \mathbf{s}' is the result of applying action \mathbf{a} to state \mathbf{s} . Action selection in continuous spaces, which calculates the greedy policy (3), is a multivariate optimization over an unknown function. Several sampling-based methods that approximate the policy efficiently exist, such as Hierarchical Optimistic Optimization applied to Trees (HOOT) [19]. HOOT uses hierarchical discretization to progressively narrow the search on the most promising areas of the input space, thus ensuring an arbitrarily small error [19]. In practice, HOOT works well for single-agent planning with value functions that have many small-scale maxima.

IV. METHODS

Our aim is to solve tasks that can be described with a set of goals (attractors) and obstacles (repellents) for acceleration-controlled robots with unknown dynamics. We require that the solution, PEARL, is efficient and adaptive. By efficient, we mean that PEARL controls agents in real-time in fully continuous and physically unbounded spaces; by adaptive, we mean that a single learning can be applied to a number of tasks.

PEARL solves PBTs in two phases, learning and acting (Figure 2), adhering to the batch RL paradigm. To start the learning phase, a user provides PEARL with basic information about the problem; the robot's DoFs, maximum accelerations, etc., and a set of objectives (preferences). The basic system information is encoded into a MDP as presented in Section IV-A. Meanwhile, given the preferences which consist of task goals and obstacles, PEARL generates the features using the methods described in Section IV-B.

With the MDP and features setup, the learning phase uses one of the AVI-based RL algorithms on a simplified problem space to discover the relative weights between the preferences (preference appraisal). The value function, approximated with Equation (2), is a linear map of preference-based features. Once the preference weights are learned, they are handed over to the planner, and the acting phase can begin.

The acting is a closed-loop feedback system, decision-making that can work online or plan trajectories offline in simulation. The planner solves problems with larger state and action domains, because the features are valid in the larger domain and capture the important elements of the task, rather than the physical space. The features enable both the efficiency, by learning on small problems, and adaptation, by allowing the policy transfer to larger problems. It is the use of the polymorphic, automatically generated features that separates PEARL from standard batch RL, and creates a virtually tuning-free task learning and completion method.

A. MDP Setup

For the **General PEARL Formulation**, we assume the robot works in continuous state and action spaces, is controlled through acceleration applied to its center of mass, and has dynamics that are not explicitly known. Let $\mathbf{s}, \dot{\mathbf{s}}, \ddot{\mathbf{s}} \in \mathbb{R}^{d_r}$ be the robot's position, velocity, and acceleration, respectively. The MDP state space is $S = \mathbb{R}^{d_s}$, where $d_s = 2d_r$, where d_r is the robot's DoFs. The state $\mathbf{s} \in S$ is joint vector $\mathbf{s} = [\mathbf{s}, \dot{\mathbf{s}}]^T$, and action $\mathbf{a} \in A = \mathbb{R}^m$ is the joint acceleration vector, $\mathbf{a} = \ddot{\mathbf{s}}$. The state transition function, which we assume is unknown, is a control-affine system (1).

We assume for training purposes the presence of a simulator or dynamics samples for the robot. The reward R is set to one when the robot achieves the goal, and zero otherwise. The tuple (S, A, D, R) defines the MDP for the robot problem.

For the **Dynamic Obstacle Avoidance Task**, the MDP setup is the joint vector of robot position and velocity, $S = \mathbb{R}^4$. The action space is the acceleration on each axis with dimension $A = \mathbb{R}^2$.

B. Feature Selection

For the **General PEARL Formulation**, we define a PBT with n_o objectives, $\mathbf{o}_1, \dots, \mathbf{o}_{n_o}$, and preferences with respect to the objectives. The objectives, points in position or velocity space, $\mathbf{o}_i \in \mathbb{R}^{d_{r_i}}$, $i = 1, \dots, n_o$, either attract or repel the one or more agents. We call preferences that attract the agent *distance-reducing*, whereas the preferences that repel

it are *intensity-reducing*; both preference types have the goal of reducing their measure to an objective.

To learn PBT with n_o objectives, $\mathbf{o}_1, \dots, \mathbf{o}_{n_o}$, we form a feature for each objective. Assuming the low dimensional task space and high-dimensional MDP space $n_o \ll d_s$, we consider *task-preference features*,

$$\mathbf{F}(\mathbf{s}, d_s) = [F_1(\mathbf{s}, d_s), \dots, F_{n_o}(\mathbf{s}, d_s)]^T.$$

Parametrized with the state space dimensionality, d_s , the features map the state space S to the preference space, and, depending on the preference type, measure either the squared intensity or distance to the objective. Let $\mathbf{p}^{o_i}(\mathbf{s})$ be a projection of the robot's state onto the minimal subspace that contains \mathbf{o}_i . For instance, when an objective \mathbf{o}_i is a point in a position space, $\mathbf{p}^{o_i}(\mathbf{s})$ is the robot's position. Similarly, when \mathbf{o}_i is a point in a velocity space, $\mathbf{p}^{o_i}(\mathbf{s})$ is the robot's velocity. Then, *distance-reducing features* are defined with

$$F_i(\mathbf{s}, d_s) = \|\mathbf{p}^{o_i}(\mathbf{s}) - \mathbf{o}_i\|^2, \quad (4)$$

and *intensity-reducing features* are defined with

$$F_i(\mathbf{s}, d_s) = (1 + \|\mathbf{p}^{o_i}(\mathbf{s}) - \mathbf{o}_i\|^2)^{-1}. \quad (5)$$

Algorithm 1 summarizes the feature selection procedure.

Algorithm 1 PEARL feature selection.

Input: $\mathbf{o}_1, \dots, \mathbf{o}_{n_o}$ objectives, pt_1, \dots, pt_{n_o} preference types

Input: MDP (S, A, D, R) ,

Output: $\mathbf{F}(\mathbf{s}, d_s) = [F_1(\mathbf{s}, d_s), \dots, F_n(\mathbf{s}, d_s)]^T$

```

1: for  $i = 1, \dots, n_o$  do
2:   if  $pt_1$  is intensity then
3:      $F_i(\mathbf{s}, d_s) = (1 + \|\mathbf{p}^{o_i}(\mathbf{s}) - \mathbf{o}_i\|^2)^{-1}$  {intensity preference}
4:   else
5:      $F_i(\mathbf{s}, d_s) = \|\mathbf{p}^{o_i}(\mathbf{s}) - \mathbf{o}_i\|^2$ , {distance preference}
6:   end if
7: end for
8: return  $\mathbf{F}(\mathbf{s}, d_s)$ 
```

For the **Dynamic Obstacle Avoidance Task**, there are two natural preferences: 1) minimize the distance to the goal, and 2) maximize the distance from obstacles. The feature vector is then formulated as the combination of the two preferences: $\mathbf{F}(\mathbf{s}) = [F_1(\mathbf{s}) \ F_2(\mathbf{s})]^T$. Providing Algorithm 1 with the goal's and obstacle's coordinates, the features are $F_1(\mathbf{s}) = \|\mathbf{p}(\mathbf{s}) - \mathbf{G}\|^2$, and $F_2(\mathbf{s}) = (\beta + d^2)^{-1}$, where \mathbf{G} is the goal's position, $\mathbf{p}(\mathbf{s})$ is the position of the robot, d is the minimum distance to the closest obstacle, and β is a constant empirically selected to be 0.01m.

V. ANALYSIS

Feature properties: Features selected in this manner have the following properties, allowing PEARL to have the potential to learn on small problems and transfer the learning to larger problems:

- *Feature domain:* The features are Lipschitz continuous and defined for the entire state space, $F_i : \mathbb{R}^m \rightarrow$

\mathbb{R} , $i = 1, \dots, d_s$, in contrast to tiling and radial basis functions [28], [2] that are active only on a section of the problem domain.

- *Projection to preference space:* Features project the state subspace into a point that measures the quality of the preferences. Thus, the state-value function approximation from (2) is an d_s -dimensional manifold in the preference space, and their number does not change as domain space change.
- *State space polymorphism:* Because they are based on the vector norm and projections, the features are polymorphic with respect to the domain dimensionality. Since learning is more computationally intensive than planning, we use lower-dimensional problems for training. The feature vector size is invariant to the number of agents, state space dimensionality, and physical space dimensions. If the agents operate in 2D space, the features consider only planar space. But, when the same agents are placed in a 3D environment, the feature set remains unchanged even though the 3D space is considered in feature calculations [7].
- *Polynomial computation time:* The feature computation time is polynomial in state space dimensionality.

Local minima analysis: For tasks with mixed objectives, such as moving obstacle avoidance, the agents follow preferences, but there are no formal completion guarantees. In fact, the value function (2) has potentially two maxima, one on each side of the obstacle.

For the purpose of this analysis, we assume that the problem is well formulated, and contains one attractor, i.e., the intersection of subspaces defined by distance reducing objectives is non-empty, and forms a connected set. Note that in order for both distance and intensity features to behave as attractor and repellent respectively, the resulting weights must be negative $\theta_i < 0$. Since a straight line is the shortest path between an agent and its attractor, we analyze the value function restricted to that line with varying obstacle distances. To simplify the analysis, we transform the value function (2). Without loss of generality, we rotate and scale the coordinate system, such that the attractor is in the origin, and the agent is on the x -axis. The two nearest obstacles lay on $(1, d)$, and $(1, -d)$. In addition, we multiply the entire function (2) by minus one, to give a rise to function $V_x(x)$. Now, we are interested in finding necessary conditions $V_x(x)$ minima, which correspond to the $V(x, y)$ maxima.

First to construct $V_x(x)$, let $c = \frac{\theta_2}{\theta_1} > 0$ be the ratio between learned weights for the obstacle and the attractor feature. The value function after the affine transformation is

$$V_x(x) = -1 * V(x, 0) = x^2 + \frac{c}{(x-1)^2 + d^2}.$$

We examine necessary conditions for $V_x(x)$'s minima based on the obstacle distance, d , and the coefficient c . Point x_0 is local minima if

$$\frac{dV_x}{dx}(x_0) = 2x_0 - \frac{2c(x_0 - 1)}{((x_0 - 1)^2 + d^2)^2} = 0, \quad (6)$$

and the second derivative is positive,

$$\frac{d^2 V_x}{dx^2}(x_0) = 2 + 2c \frac{(x_0 - 1)^2 - d^2}{((x_0 - 1)^2 + d^2)^3} > 0. \quad (7)$$

For 6 to hold, the following must be the case,

$$x_0 < -1, \text{ or } x_0 > 0, \quad (8)$$

and given (8), (7) holds when,

$$\|x_0 - 1\| > \|d\|. \quad (9)$$

We conclude that the value function (2) for dynamical obstacle avoidance has two local maxima, one to the left of the goal, and the other one to the right of the obstacles. If the agent is located between the goal and the obstacles, it will settle at the equilibrium point to the left of the goal. If the agent is on the right side of the obstacles, it will settle to the right of the obstacle, unless the obstacles move. Further, the equilibrium point tend to the goal, as the distance between the obstacles increase.

To verify the analysis, we perform an empirical study, depicted in 3. The value function either has a single maxima near an attractor (magenta, and cyan lines in Figure 3), has two maxima (blue and green lines), or has an inflection point near the obstacles (red line). Inspection of the partial derivative $\frac{\partial V}{\partial y}$ at the minima points in Figure 3 reveals that these points are saddle points.

In summary, when the obstacles are far enough apart there is only a global maximum. As the obstacles come closer together, a new region of attraction forms on the other side of the obstacle. If the agent gets into the local maximum region of attraction, gradient-descend methods will trap it. Sampling-based greedy methods such as HOOT, however, might get the agent out of the region of attraction if it is sufficiently close to the boundary.

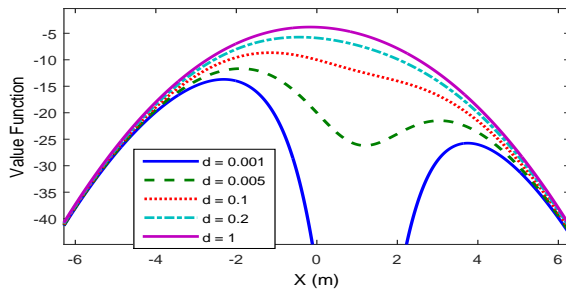


Fig. 3. Value function inflection points for $c = 100$.

VI. RESULTS

We now demonstrate PEARL for the Dynamic Obstacle Avoidance Task. PEARL was implemented in MATLAB and all experiments were run on an Intel i3-2120 at 3.3GHz with 4GB RAM.

Figure 1 illustrates the testing environment and task for the point-like holonomic robot to navigate from the starting location to the goal without colliding with circular obstacles.

The robot observes only the current position of the closest obstacle, and has no information about its velocity.

Learning Setup: We use 4 stationary obstacles placed at $[3 \text{ m}, 0 \text{ m}]$, $[0 \text{ m}, 3 \text{ m}]$, $[0 \text{ m}, -3 \text{ m}]$, $[-3 \text{ m}, 0 \text{ m}]$ to learn the weights between the two features. The goal is at the origin. The sampling space is inside a two-dimensional hypercube $[-5 \text{ m}, 5 \text{ m}]^2$. The robot has a maximum speed of 0.37 m/s , and a maximum acceleration of 3 m/s^2 . We run AVI [5] with HOOT policy approximation [19], as the learning agent in PEARL for 300 iterations to learn the feature vector weights.

Learning Result: The resulting weights are $\theta = [-0.23 - 0.1696]^T$. All simulations are done at 10 Hz . The time to learn is 123 s .

Planning Environment Setup: The planning task environment is illustrated in Figure 1. The robot must travel from the start position $[25 \text{ m}, 0 \text{ m}]$ to the goal at $[-25 \text{ m}, 0 \text{ m}]$ under the same speed and acceleration constraints as used for learning.

We maintain the constant density of moving obstacles by restricting the robot and moving obstacles to lie in a circle with radius 50 m . When an obstacle hits the boundary of the circle, it is transported to the antipodal position on the circle and continues evolving from this new position. The resulting density of moving obstacles is similar to [3].

Obstacles Setup: To demonstrate motion planning with PEARL in environments with many moving obstacles with complex stochastic dynamics, our planning environment has $N = \{300, 450, 600, 750, 900\}$ randomly placed moving obstacles with hybrid stochastic dynamics. The hybrid dynamics involves three stochastic dynamics modes (linear, arc, and swerve). An obstacle is in one of the three modes at any given moment and the initial mode is randomly determined according to the initial probability shown in Figure 4. Every $T_{\text{resample}} = 1 \text{ s}$, the obstacle has a probability of 0.1 to switch mode and the probability of switching to a particular mode is specified in Figure 4.

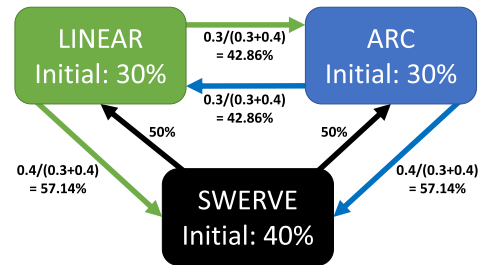


Fig. 4. Hybrid stochastic dynamics includes three stochastic dynamics modes (linear, arc and swerve) and probabilities to switch modes. The initial percentage value represents the number ratio at the beginning of the experiment.

Obstacles in linear dynamics mode has a fixed heading but the speed of travel is sampled stochastically from the set: $\{0.1, 0.2, 0.5, 0.7\} \text{ m/s}$ with probability $\{0.4, 0.1, 0.2, 0.3\}$. Obstacles in arc dynamics mode move counter clock-wise with radius 5 m at a stochastically sampled angular speed

from the set $\{0.039, 0.058, 0.088, 0.117\}$ rad/s with probability $\{0.4, 0.1, 0.2, 0.3\}$. Obstacles in swerve dynamics mode are changing heading linearly at the rate of $\pi/3$ rad/s between $[-\phi_{invert}, \phi_{invert}]$. ϕ_{invert} is sampled uniformly between $[-\pi/2, \pi/2]$. All obstacles, regardless of dynamics modes, resample stochastically with frequency $\frac{1}{T_{resample}}$.

The hybrid stochastic dynamics is highly complex and unpredictable. Figure 5 trajectories show the hybrid dynamics (magenta) generates unpredictable and diverging paths between runs, e.g., (dark and light trajectories, also numbered 1 and 2 in the legend). For comparison, trajectories with constant dynamics modes (blue, grey and green) are also shown.

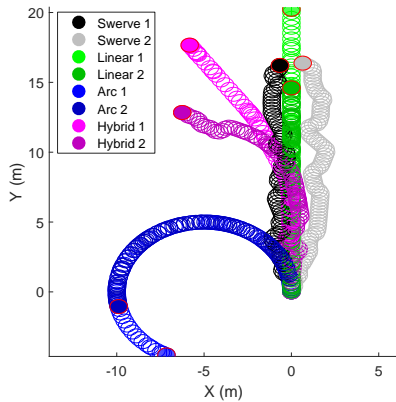


Fig. 5. Example paths of stochastically moving obstacles in various dynamics modes. All obstacles start at the origin and move for 50 s. The dark and light magenta obstacles have hybrid dynamics, which switches between all three modes probabilistically. For comparison, constant dynamics modes are also shown. The dark and light green/blue trajectories are in constant linear/arc dynamics mode while the black and grey trajectories are in constant swerve mode.

The obstacles are circles with radius $r_{obs} = 0.5$ m. The average speed of obstacles (0.37 m/s regardless of dynamics modes) is identical to the maximum speed of the robot.

Comparison Methods: We compared our method with two obstacle avoidance methods, which, similarly to PEARL, only consider the current position (and velocity) of obstacles: Gaussian APF and VO. The Gaussian APF method considers only the position of obstacles. It combines a linear attractive potential toward the goal and a repulsive potential from obstacles [17]. The obstacle potentials are Gaussians with $\sigma = 0.45$ m around obstacles, tuned empirically for this problem. The relative strength between the attractive and repulsive potential, α , has a significant impact on the success rate and needs to be manually tuned. Larger α represents a more goal-greedy robot behavior. We compared various values of α to our method. VO considers the position and velocity of obstacles [9]. A C++ implementation of ORCA [26] was downloaded and modified to support a single robot with multiple moving obstacles from [1]. This was done by disabling the VO calculation and reciprocal collision avoidance for agents acting as moving obstacles. The agent acting as the robot calculates VOs in order to avoid moving

obstacles.

Planning Result: Figure 6a and 6b show that planning with PEARL has a higher probability of successfully avoiding obstacles, and reaches the goal in less time compared to the Gaussian APF method. The success rate and task finish time of the Gaussian method depends greatly on the parameter α . This parameter has to be tweaked manually or by optimization algorithms [25] after many planning trials. PEARL balances the features (similar to finding the optimal α) in the learning phase with a simplified scenario, and is able to transfer the weights to the online plan with comparable or better performance.

Figure 6a also shows that PEARL has a higher success rate than VO. This is primarily due to VO's optimal velocity obstacle formulation, which assumes the obstacle has a fixed velocity. This results in trajectories optimal in finish time but very small clearance to obstacles. This often causes collisions with the highly unpredictable obstacles. PEARL on the other hand, by balancing the goal reaching and obstacle avoidance features, was able to generate trajectories with enough clearance to account for the unpredictability of obstacle motion. In addition, Figure 6b shows the running time per planning step for PEARL is lower than VO, even though PEARL is implemented in MATLAB and VO is in C++. PEARL scales linearly with the number of obstacles and is capable of generating higher success rate trajectories in real-time (Figure 6c). This suggests that PEARL is a viable alternative method for dynamic obstacle avoidance, even when the obstacles are moving with highly unpredictable hybrid stochastic dynamics.

VII. CONCLUSION AND FUTURE WORK

This paper presents PEARL, a solution for high-dimensional preference-balancing motion problems, that is efficient, adaptive, and controls the agent in real-time. The method uses features that work in continuous domains, scale polynomially with the problem's dimensionality, and are polymorphic with respect to the domain dimensionality. PEARL was demonstrated on a complex Dynamic Obstacle Avoidance Task where the agent has to progress toward the goal while avoiding collision with 900 moving obstacles with highly unpredictable hybrid stochastic dynamics using only limited information about the environment. Our experiments show that PEARL outperforms the traditional Gaussian APF method and VO, a state of art algorithm for the obstacle avoidance.

In future work, we plan to investigate the relationship between the difficulty of the training environment and the planning environment. In particular, we will impose stronger kinodynamic constraints to the robot and faster moving obstacles to tune the trade-off between learning and planning difficulties.

ACKNOWLEDGMENTS

The authors thank Marco Morales for the feedback on the manuscript. We thank Jur Van Den Berg for feedbacks on modifying ORCA. PEARL was developed by Faust while

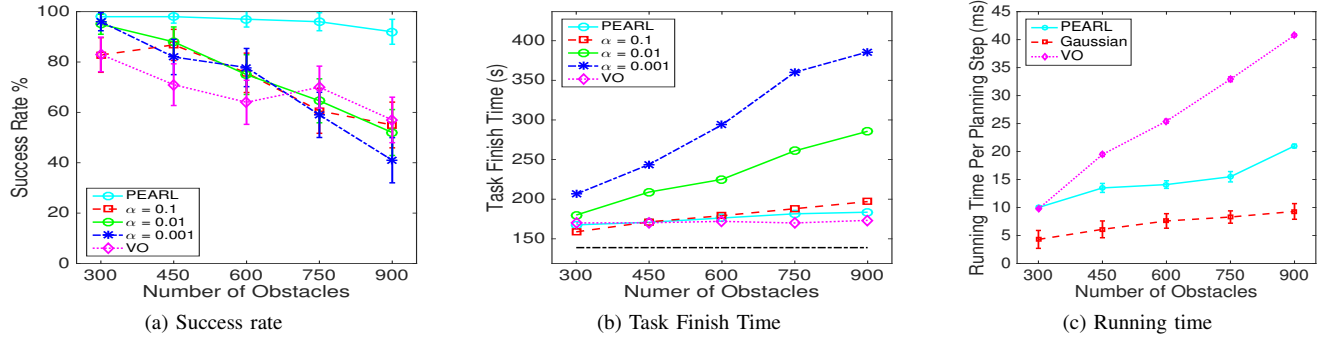


Fig. 6. Trajectory characteristics (averaged over 200 trials) for environments with varied complexity (number of obstacles). The agent planned with the same policy for all runs. A policy learned with four static obstacles. (a) Collision-free success rate of finishing the task. The error bars are computed using the 99% confidence interval formalism. (b) Amount of time for the agent to reach the goal without collision, the dotted line is the minimum finish time without obstacles. (c) Running time per planning step. PEARL has a higher success rate and less running time per planning step compared to VO.

at University of New Mexico, and partially supported with New Mexico Space Grant. This work is partially supported by National Science Foundation Grant IIS-1528047.

REFERENCES

- [1] J. V. D. Berg, S. J. Guy, J. Snape, M. C. Lin, and D. Manocha. Rvo2 library: Reciprocal collision avoidance for real-time multi-agent simulation. <http://gamma.cs.unc.edu/RVO2/>.
- [2] L. Buşoniu, R. Babuška, B. De Schutter, and D. Ernst. *Reinforcement Learning and Dynamic Programming Using Function Approximators*. CRC Press, Boca Raton, Florida, 2010.
- [3] H.-T. Chiang, N. Malone, K. Lesser, M. Oishi, and L. Tapia. Aggressive moving obstacle avoidance using a stochastic reachable set based potential field. In *Proc. Int. Workshop on Algorithmic Foundations of Robotics (WAFR)*, 2014.
- [4] H.-T. Chiang, N. Rackley, and L. Tapia. Stochastic ensemble simulation motion planning in stochastic dynamic environments. In *Proc. IEEE Int. Conf. Intel. Rob. Syst. (IROS)*, pages 3836–3843, 2015.
- [5] D. Ernst, M. Glavic, P. Geurts, and L. Wehenkel. Approximate value iteration in the reinforcement learning context. application to electrical power system control. *International Journal of Emerging Electric Power Systems*, 3(1):1066.1–1066.37, 2005.
- [6] A. Faust, H.-T. Chiang, N. Rackley, and L. Tapia. Dynamic obstacle avoidance with PEARL: Preference appraisal reinforcement learning. In *Second Annual Machine Learning in Planning and Control of Robot Motion Workshop at IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS) 2015, Hamburg, Germany*, October 2015.
- [7] A. Faust, I. Palunko, P. Cruz, R. Fierro, and L. Tapia. Automated aerial suspended cargo delivery through reinforcement learning. *Artificial Intelligence*, 2014.
- [8] A. Faust, P. Ruymgaart, M. Salman, R. Fierro, and L. Tapia. Continuous action reinforcement learning for control-affine systems with unknown dynamics. *Automatica Sinica, IEEE/CAA Journal of*, 1(3):323–336, 2014.
- [9] P. Fiorini and Z. Shiller. Motion planning in dynamic environments using velocity obstacles. *The International Journal of Robotics Research*, 17(7):760–772, 1998.
- [10] W. Glover and J. Lygeros. A stochastic hybrid model for air traffic control simulation. In *Hybrid Systems: Computation and Control*, pages 372–386. Springer, 2004.
- [11] H. Hasselt. Reinforcement learning in continuous state and action spaces. In M. Wiering and M. Otterlo, editors, *Reinforcement Learning*, volume 12 of *Adaptation, Learning, and Optimization*, pages 207–251. Springer Berlin Heidelberg, 2012.
- [12] T. Hester and P. Stone. TEXPLORE: real-time sample-efficient reinforcement learning for robots. *Machine Learning*, 90(3):385–429, 2013.
- [13] L. Jaillet and T. Simeon. A PRM-based motion planner for dynamically changing environments. In *Proc. IEEE Int. Conf. Intel. Rob. Syst. (IROS)*, 2004.
- [14] M. Kallman and M. Mataric. Motion planning using dynamic roadmaps. In *Proc. IEEE Int. Conf. Robot. Autom. (ICRA)*, pages 4399–4404, 2004.
- [15] J. Kober, D. Bagnell, and J. Peters. Reinforcement learning in robotics: A survey. *International Journal of Robotics Research*, 32(11):1236–1272, 2013.
- [16] T. Kunz and M. Stilman. Manipulation planning with soft task constraints. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 1937–1942, October 2012.
- [17] C.-P. Lam, C.-T. Chou, K.-H. Chiang, and L.-C. Fu. Human-centered robot navigation towards a harmoniously human-robot coexisting environment. *Robotics, IEEE Transactions on*, 27(1):99–112, 2011.
- [18] S. M. LaValle. *Planning Algorithms*. Cambridge University Press, Cambridge, U.K., 2006.
- [19] C. Mansley, A. Weinstein, and M. Littman. Sample-based planning for continuous action markov decision processes. In *Proc. of Int. Conference on Automated Planning and Scheduling*, 2011.
- [20] M. Otte and E. Frazzoli. RRT-X: Real-time motion planning/replanning for environments with unpredictable obstacles. In *Proc. Int. Workshop on Algorithmic Foundations of Robotics (WAFR)*, 2014.
- [21] K. Mülling, J. Kober, and J. Peters. A biomimetic approach to robot table tennis. *Adaptive Behavior*, 19(5):359–376, 2011.
- [22] N. Rivera, J. A. Baier, and C. Hernandez. Weighted real-time heuristic search. In *Proceedings of the 2013 International Conference on Autonomous Agents and Multi-agent Systems, AAMAS '13*, pages 579–586, Richland, SC, 2013. International Foundation for Autonomous Agents and Multiagent Systems.
- [23] A. C. Shkolnik and R. Tedrake. Path planning in 1000+ dimensions using a task-space voronoi bias. In *Proc. IEEE Int. Conf. Robot. Autom. (ICRA)*, pages 2061–2067. IEEE, may 2010.
- [24] M. Stilman. Global manipulation planning in robot joint space with task constraints. *IEEE/RAS Transactions on Robotics*, 26(3):576–584, 2010.
- [25] K. C. Tan and W. Ming-Liang. Evolutionary artificial potential fields and their application in real time robot path planning. In *Evolutionary Computation, 2000. Proceedings of the 2000 Congress on*, pages 256–263, 2000.
- [26] J. Van Den Berg, S. J. Guy, M. Lin, and D. Manocha. Reciprocal n-body collision avoidance. In *Robotics research*, pages 3–19. Springer, 2011.
- [27] C. Wirth and J. Fürnkranz. Preference-based reinforcement learning: a preliminary survey. In *ECML/PKDD-13 Workshop on Reinforcement Learning from Generalized Feedback: Beyond Numeric Rewards*, Sep 2013.
- [28] C. Wu. *Novel function approximation techniques for large-scale reinforcement learning*. PhD thesis, Northeastern University, Apr 2010.