# Risk Conditioned Neural Motion Planning

Xin Huang[1], Meng Feng[1], Ashkan Jasour[1], Guy Rosman[2], and Brian Williams[1]

*Abstract* — **Risk-bounded motion planning is an important yet difficult problem for safety-critical tasks. While existing mathematical programming methods offer theoretical guarantees in the context of constrained Markov decision processes, they either lack scalability in solving larger problems or produce conservative plans. Recent advances in deep reinforcement learning improve scalability by learning policy networks as function approximators. In this paper, we propose an extension of soft actor critic model to estimate the execution risk of a plan through a risk critic and produce risk-bounded policies efficiently by adding an extra risk term in the loss function of the policy network. We define the execution risk in an accurate form, as opposed to approximating it through a summation of immediate risks at each time step that leads to conservative plans. Our proposed model is conditioned on a continuous spectrum of risk bounds, allowing the user to adjust the risk-averse level of the agent on the fly. Through a set of experiments, we show the advantage of our model in terms of both computational time and plan quality, compared to a state-of-the-art mathematical programming baseline, and validate its performance in more complicated scenarios, including nonlinear dynamics and larger state space.**

## I. INTRODUCTION

Motion planning is an important task in many robotics applications, such as rescue robots and autonomous vehicles. One of the most popular approaches for motion planning is reinforcement learning, which learns an optimal policy that minimizes the cost, through exploration and exploitation. Recently, deep reinforcement learning has been proposed to approximate the policy function or the cost function by deep neural networks, and has achieved great success in applications where the state and action space is high-dimensional, such as robotics manipulation [1] and Go game [2].

Despite the success of reinforcement learning in the context of Markov decision processes (MDPs), its objective of naively minimizing a cost function is not sufficient in safety-critical tasks. For instance, in planning with obstacles for rescue robots (see Fig. 1) or autonomous vehicles, the objective is not only to minimize a certain cost function of the policy, such as time to a designated goal or fuel consumption, but also to satisfy the safety constraint of not colliding with obstacles, which may sacrifice the cost performance. Although in many cases, we can add the constraint as an infinite term in the cost function, this may lead to infeasible solutions, when the constraint is unavoidable (i.e. if we want

[1]Computer Science and Artificial Intelligence Laboratory (CSAIL), Massachusetts Institute of Technology, Cambridge, MA 01239, USA xhuang@csail.mit.edu

[2]Toyota Research Institute, Cambridge, MA 01239, USA
This article solely reflects the opinions and conclusions of its authors and not TRI or any other Toyota entity.
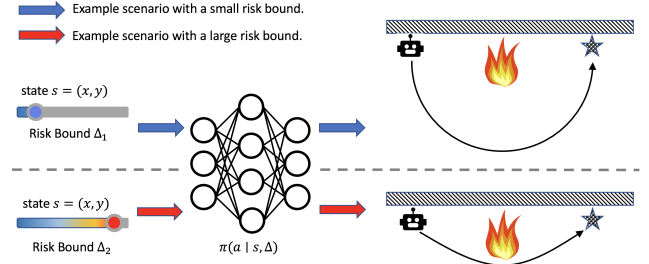


Fig. 1: Overview of the proposed risk conditioned policy network. Given the current state and risk bound, the network generates an optimal action for the agent, subject to the risk bound. The policy is conditioned on an arbitrary upper bound on the probability of violating safety constraints, which has important applications in safety-critical domains such as rescue robots and autonomous vehicles. Top and bottom depict planning scenarios with a small risk bound and a large risk bound, respectively, using the same model.

to avoid an obstacle with unbounded uncertainties) and we can never guarantee safety. Therefore, it is more desirable to bound the probability of violating the constraint by a certain level, which is defined as a chance constraint in [3].

The requirement of satisfying additional constraint motivates the formulation of constrained Markov decision process (CMDP) [4], in which the agent needs to satisfy the constraint over an auxiliary cost while minimizing the cost function. Although solving finite CMDP has been well studied with methods such as linear programming [4] or mixed integer linear programming (MILP) [5], it remains a challenge to solve for high-dimensional or large CMDP instances. In cases where the constraint is probabilistic and the objective is to bound the probability of violating a constraint, [6] introduces a special form of CMDP, named chance constraint MDP (CC-MDP), and proposes a heuristic forward search method in the discrete action space to find the desired policy efficiently, yet the performance heavily relies on the quality of heuristics.

Inspired by the recent success of deep reinforcement learning, we leverage deep neural networks to learn function approximators for the policy function and the risk (or chance constraint) function, similar to [7]–[9]. The learned chance constraint is used in the policy optimization step to produce risk-bounded policies. While existing methods assume a fixed upper bound on the chance constraint, our model is trained to generate policies conditioned on a range of continuous upper bound values. This allows us to control the

agent with different risk tolerance levels on the fly, without the need to retrain the model. One motivating example is that in rescue tasks, as illustrated in Fig. 1, we may want to assign different risk bounds to the rescue robot based on the hard time limit, so that it can rescue the target in time. This requires the planner to generate policies in real-time conditioned on different risk bounds. Although it is possible to pre-generate a library of policies under different bound levels, it is less time efficient and space efficient, especially when the discretization of the bound requires a high resolution. An overview of our approach is presented in Fig. 1.

Our contributions are as follows: i) We learn a risk estimator as a deep neural network to approximate the execution risk of a given policy, which provides task-specific information in addition to standard cost measures in MDP. ii) We leverage the learned risk estimator to learn a risk-bounded policy network, by adding the exceeded risk as an additional penalty in the policy loss, which balances performance and safety in the resulting plan. iii) Our policy network is conditioned on a continuous spectrum of upper risk bound, providing the flexibility to change the agent behavior on the fly at test time. iv) We provide a detailed comparison between a state-of-the-art MILP-based method and our proposed method, and show our method achieves better performances in terms of solution quality and computational time, without violating the risk bound.

## II. RELATED WORK

### A. Risk-Bounded Motion Planning

Risk-bounded motion planning problems have been studied in the forms of CMDPs and CC-MDPs in many works. Early methods show that CMDPs can be solved with linear programs (LPs) [4], [10] and MILPs [11]. Most of these methods work with deterministic constraints. In [5], the constraint is modeled as the probability of failure, and an approach named Iterative Risk Allocation (IRA), is proposed to generate constrained policies by iteratively assigning the global probability constraint into individual probability budgets at each time step through a union bound, and solve for a MILP instance given the individual budgets. Aside from MILP-based methods, [12] proposes a sampling-based approach to generate constrained motion plans with probabilistic guarantees. A common assumption of these works is that the global constraint can be expressed as the sum of sub-constraints over each time step, which may not hold if we define the constraint as the probability of collision over the entire path, since the joint disjunctive probability is not equivalent to the summation of probabilities of individual events. In fact, the summation is usually an upper bound of the joint probability, and optimizing for the summed constraints will likely lead to conservative policies.

The conservatism is resolved in [6], which defines the risk as the probability of collisions in an exact form, and proposes RAO* to find optimal policies through heuristic forward search, which is demonstrated to work well in vehicle planning [13], [14] and aircraft routing [15] domains.

However, RAO* is limited to discrete action space due to its tree-based search approach.

Overall, due to the NP-hard nature of CMDP problems [16], both (MI)LP and search-based approaches suffer from scalability, especially in high-dimensional and continuous state space and action space.

### B. Safe Reinforcement Learning

Safety is a popular topic in reinforcement learning as many algorithms have been applied to real world applications, which requires the agent to explore the policy space while being safe to a certain extent. In this paper, we focus on RL methods that model risk as an explicit constraint in the overall objective function. In [7], [17], the constrained RL problem is converted to an unconstrained RL problem using Lagrangian multiplier, and then solved by standard RL algorithms such as actor critic. Instead of solving the dual problem using the Lagrangian, [9] approximates the objective function and the constraint function through quadratic surrogate functions, and solves for convex quadratically constrained quadratic program directly. As an extension to [7], [8] combines a constrained policy optimization algorithm with imitation learning to bootstrap training time and achieves comparable performance. These methods define the constraint as a summed auxiliary cost from each time step, which leads to conservative policies if we model the constraint as the probability of collisions, as discussed in Sec. II-A. Furthermore, the upper bound of the constraint is usually predefined as a fixed value, which would require the user to retrain the model if a different bound is needed. In our paper, we extend the state-of-the-art deep RL models to handle probabilistic constraints, in an accurate form, and flexible risk bounds, in the task of risk-bounded motion planning.

In addition to constrained MDP-based approaches, there exist works that improve safety by generating more samples in the risky region to bootstrap performance in critical scenarios [18]; by using a safety layer at the end of a deep neural network to verify the safety of the resulting policy and replacing with a backup safe action if needed [19]; by proposing a reachability-based trajectory safe guard to ensure the safety of a policy [20], etc. In this paper, we focus on generating risk-bounded policies directly by modeling the risk as an explicit constraint in the objective function.

### C. Conditioned Reinforcement Learning

Conditioned reinforcement learning provides great flexibility and improves generalizability when applying RL algorithms. A notable example is goal-conditioned RL [21]–[23], which produces goal-conditioned policies without assuming a fixed goal. This provides more capability to the agent as it can navigate to any goal locations in the environment without retraining. In our work, we are inspired by goal-conditioned policy and propose risk conditioned policy that generates policies conditioned on a continuous spectrum of upper risk bound.

## III. PROBLEM FORMULATION

In this work, we consider the risk-bounded motion planning problem in the absence of the agent models. More precisely, we aim at solving the chance constrained Markov decision process (CC-MDP) problem defined as follows:

**Definition 1. Chance-Constrained Markov Decision Process:** A Chance-Constrained Markov Decision Process (CC-MDP) is defined a tuple $\langle \mathcal{S}, \mathcal{A}, T, R, s_0, h, \mathcal{C}, \Delta \rangle$ as follows [6]:

- $\mathcal{S}$ is a set of continuous states.
- $\mathcal{A}$ is a set of actions for the agent model.
- $T : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to \mathbb{R}$ is a state transition function between the states.
- $R : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ is a reward function.
- $s_0$ is the initial state.
- $h$ is the finite execution horizon.
- $\mathcal{C}$ is a set of safety constraints defined over $\mathcal{S}$.
- $\Delta$ is the upper risk bound.

Given a definition of CC-MDP, the objective is to find a policy that maximizes the expected cumulative reward function:

$$\pi^* = \arg\max_{\pi} \ \mathbb{E}_{(s_t, a_t) \sim \pi} \left[ \sum_{t=0}^{h} R(s_t, a_t) \right], \qquad (1)$$

while satisfying the chance constraint with respect to the safety constraints $\mathcal{C}$:

$$er(s_0, \mathcal{C}|\pi) \leq \Delta, \qquad (2)$$

where $er$ is the execution risk defined as follows:

$$er^{\pi}(s_t) = er(s_t, \mathcal{C}|\pi) = 1 - \Pr\left( \bigwedge_{i=t}^{h} Sa_i = 1 \middle| s_t, \pi \right). \quad (3)$$

where $Sa_i$ is a Bernoulli random variable with value 1, when the state has not violated any constraint in $\mathcal{C}$, defined over $\mathcal{S}$, at time $i$.

The execution risk can be computed in a recursive way as follows [6]:

$$er^{\pi}(s_t) = r_b(s_t) + (1 - r_b(s_t))\mathbb{E}[er^{\pi}(s_{t+1})], \qquad (4)$$

where $r_b(s_t)$ is the immediate risk at time $t$. An example of risk is the probability of collision between the agent and the obstacles. An important distinction between our work and existing CMDPs [5], [9] is that they approximate the execution risk as a sum of step-wise immediate risks:

$$er^{\pi}_{approx}(s_t) = \sum_{i=t}^{h} r_b(s_i), \qquad (5)$$

which is an upper bound of $er^{\pi}(s_t)$, and solving for a policy following this definition leads to conservative results.

In this paper, we assume that the agent model and the reward function are hidden. Hence, we aim at solving the following planning problem:

**Risk-Bounded Motion Planning Problem:** Given a CC-MDP, with unknown transition and reward functions, i.e., $T, R$, we look for a policy $\pi$ to maximize the expected cumulative reward function in (1) with respect to the risk constraints in (2).

## IV. APPROACH

In this section, we propose a deep neural network as a function approximator to solve the motion planning problem defined in Section III. We first review the soft actor critic method, which is a popular reinforcement learning model that improves the stability of training compared to standard actor critic method. We then introduce our extended model, risk conditioned soft actor critic that includes a risk critic to estimate the probability of violating risk constraint, and show how we leverage the learned risk in policy optimization to learn risk-bounded policies. Finally, we present an algorithm to train our model through gradient descent.

### A. Soft Actor Critic

Soft Actor Critic (SAC) [24] is an off-policy actor critic deep reinforcement learning algorithm based on max entropy reinforcement learning. The objective is to find a policy that maximizes the maximum entropy:

$$\pi^* = \arg\max_{\pi} \ \sum_{t=0}^{h} \mathbb{E}_{(s_t, a_t) \sim \rho_{\pi}} \left[ R(s_t, a_t) + \alpha\mathcal{H}(\pi(\cdot|s_t)) \right] \ (6)$$

where temperature parameter $\alpha$ determines the relative importance of the entropy term versus the reward, and thus controls the stochasticity of the optimal policy [24]; $\mathcal{H}(\pi(\cdot|s_t))$ is the entropy of the policy $\pi$ at state $s_t$; $\rho_{\pi}(s_t, a_t)$ denotes the state-action marginals of the trajectory distribution induced by a policy $\pi(a_t|s_t)$.

As an actor critic algorithm, SAC alternates between policy evaluation, by estimating the value function for a policy, and policy improvement, by using the estimated value function to obtain a better policy. More specifically, SAC leverages a state value function $V_{\psi}(s_t)$ parameterized by $\psi$, a soft Q-function $Q_{\theta}(s_t, a_t)$ parameterized by $\theta$, and a policy function $\pi_{\phi}(a_t|s_t)$ parameterized by $\phi$. These functions are modeled as deep neural networks (i.e. $V_{\psi}(s_t)$, $Q_{\theta}(s_t, a_t)$), or Gaussian parameters from a deep neural network (i.e. $\pi_{\phi}(a_t|s_t)$).

*1) Actor Model:* According to [24], the policy parameters can be learned by directly minimizing the following:

$$J_{\pi}(\phi) = \mathbb{E}_{s_t \sim \mathcal{D}} \left[ \log \pi_{\phi}(f_{\phi}(s_t)|s_t) - Q_{\theta}(s_t, f_{\phi}(s_t)) \right], \ (7)$$

where $\mathcal{D}$ is the replay buffer representing the distribution of previously sampled states and actions, $\pi_{\phi}$ is defined implicitly in terms of $f_{\phi}(s_t)$, which is a neural network transformation that is used to reparameterize the policy, such that

$$a_t = f_{\phi}(s_t). \qquad (8)$$

*2) Critic Model:* The soft Q-function parameters can be trained to minimize the soft Bellman residual:

$$J_Q(\theta) = \mathbb{E}_{(s_t, a_t) \sim \mathcal{D}} \left[ \frac{1}{2} \left( Q_{\theta}(s_t, a_t) - \hat{Q}(s_t, a_t) \right)^2 \right], \quad (9)$$

with

$$\hat{Q}(s_t, a_t) = R(s_t, a_t) + \gamma\mathbb{E}_{s_{t+1} \sim p} V_{\psi}(s_{t+1}), \qquad (10)$$

where the soft value function $V_\psi$ is another deep neural network that is trained to approximate the value function. We refer to [24] for more details of $V_\psi$.

### B. Risk Conditioned Soft Actor Critic

In risk conditioned actor critic, we aim to learn an extra risk critic model $Q_{er_\zeta}$, parameterized by $\zeta$, that estimates the expected execution risk starting at a given state and acting according to the policy, defined in Eq. (4).

*1) Risk Critic Model:* We use a risk critic model to estimate the execution risk, and it can be trained to minimize the L2 loss to the actual execution risk computed from data:

$$J_{Q_{er}}(\zeta) = \mathbb{E}_{(s_t, a_t) \sim \mathcal{D}} \left[ \frac{1}{2} \left( Q_{er_\zeta}(s_t, a_t) - \hat{Q}_{er}(s_t) \right)^2 \right],$$
(11)

with

$$\hat{Q}_{er}(s_t) = r_b(s_t) + (1 - r_b(s_t)) \mathbb{E}_{s_{t+1} \sim D} \left[ \hat{Q}_{er}(s_{t+1}) \right]. \quad (12)$$

The gradient is computed as follows:

$$\hat{\nabla}_\zeta J_{Q_{er}}(\zeta) = \nabla_\zeta Q_{er_\zeta}(s_t, a_t)(Q_{er_\zeta}(s_t, a_t) - \hat{Q}_{er}(s_t)).$$
(13)

*2) Risk-Bounded Actor Model:* Given the risk critic model, we can update the loss function for the actor model as:

$$J_\pi(\phi) = \mathbb{E}_{s_t \sim \mathcal{D}} \left[ \log \pi_\phi(f_\phi(s_t, \Delta) | s_t) - Q_\phi(s_t, f_\phi(s_t, \Delta)) \right.$$
$$\left. + \lambda_{er} ReLU \left( Q_{er_\zeta}(s_t, f_\phi(s_t, \Delta)) - \Delta \right) \right],$$
(14)

where $ReLU(x)$ returns $x$ if $x \geq 0$ and 0 otherwise. This extra term in the policy loss adds an penalty if the estimated execution risk is larger than the risk bound, with a coefficient $\lambda_{er}$ to balance performance and safety. Compared to standard SAC, the transformation function $f_\phi$ is updated to include $\Delta$ as its input so that the policy is conditioned on the upper risk bound:

$$a_t = f_\phi(s_t, \Delta). \quad (15)$$

### C. Algorithm

The algorithm to train a risk conditioned soft actor critic is illustrated in Alg. 1. The key difference to the standard soft actor critic algorithm [24] is that we need to train a separate *risk critic* that estimates the execution risk. This requires us to collect the immediate risk $r_b$ at each state, so that we can use it as a supervisory cue to train the risk critic network. In addition, we generate random upper risk bound samples $\Delta$ from a uniform distribution, so that our policy network can learn to produce risk-bounded policies conditioned on different bounds. This, in practice, allows the online adjustment of the upper risk bound $\Delta$, subsequently changing the aggressiveness of the agent's actions. In the gradient descent step, we compute the gradient of each network, and we refer to [24] for the detailed derivation of the gradients for the policy network and soft $Q$-function network.

---

**Algorithm 1:** Risk conditioned SAC algorithm.

Initialize parameter vectors $\phi$, $\theta$, $\zeta$, $\Delta \sim U[0, 1]$.
**for** each iteration **do**
  **for** each environment step **do**
    $\mathbf{a}_t \sim \pi_\phi(\mathbf{a}_t | \mathbf{s}_t, \Delta)$
    $\mathbf{s}_{t+1} \sim p_T(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$
    $\mathcal{D} \leftarrow \mathcal{D} \cup \{(\mathbf{s}_t, \mathbf{a}_t, R(\mathbf{s}_t, \mathbf{a}_t), r_b(\mathbf{s}_t), \Delta, \mathbf{s}_{t+1})\}$
  **end for**
  **for** each gradient step **do**
    $\theta \leftarrow \theta - \lambda_Q \hat{\nabla}_\theta J_Q(\theta)$
    $\zeta \leftarrow \zeta - \lambda_{Q_{er}} \hat{\nabla}_\zeta J_{Q_{er}}(\zeta)$
    $\phi \leftarrow \phi - \lambda_\pi \hat{\nabla}_\phi J_\pi(\phi)$
  **end for**
**end for**

---

## V. RESULTS

In this section, we introduce two sets of experiments to validate our proposed model. The first set of experiments provide a detailed study comparing our method and a state-of-the-art MILP-based baseline in a maze environment that simulates robotics rescue tasks; the second set of experiments showcase the performance of our method in more complicated scenarios, including nonlinear dynamics and larger state space, which usually pose great challenges for MILP-based methods. In each experiment, we define the problem setup and introduce model details, followed by results and discussions.

### A. Maze Environment

*1) Problem Setup:* In order to show that our model provides risk conditioned policies efficiently, we introduce a *OneObstacle* maze environment that includes an obstacle between the start location and the goal location, as illustrated in Fig. 3. The agent state space $S \in \mathbb{R}^2$ represents the continuous position $(x, y)$ bounded by the maze boundary $[0, 10] \times [0, 10]$, and the action space $A \in \mathbb{R}^2$ represents the linear velocity $(v_x, v_y)$ with the magnitude bounded by 1. The dynamics of the agent follows linear change of states:

$$\dot{x} = v_x, \quad \dot{y} = v_y. \quad (16)$$

We consider an additive Gaussian noise to the agent state with 0 mean and standard deviation of 1 along each axis, to model the uncertainty. The unbounded uncertainty prevents us from finding a policy that is risk-free; instead, we can only find a policy that bounds the probability of violating the risk constraint, in which the risk is defined as the probability of the agent colliding with any obstacles. While there exist a number of efficient risk assessment methods [25], [26], we use a naive Monte Carlo sampling method to compute the collision risk for its simplicity.

*2) Model Details:* To find risk conditioned policies, we learn a risk conditioned SAC model[1] described in Sec. IV-B. The model includes two $Q$ function estimators as 3-layer

---

[1]The source code will be open sourced in the near future at https://github.com/cyrushx/risk_sac.

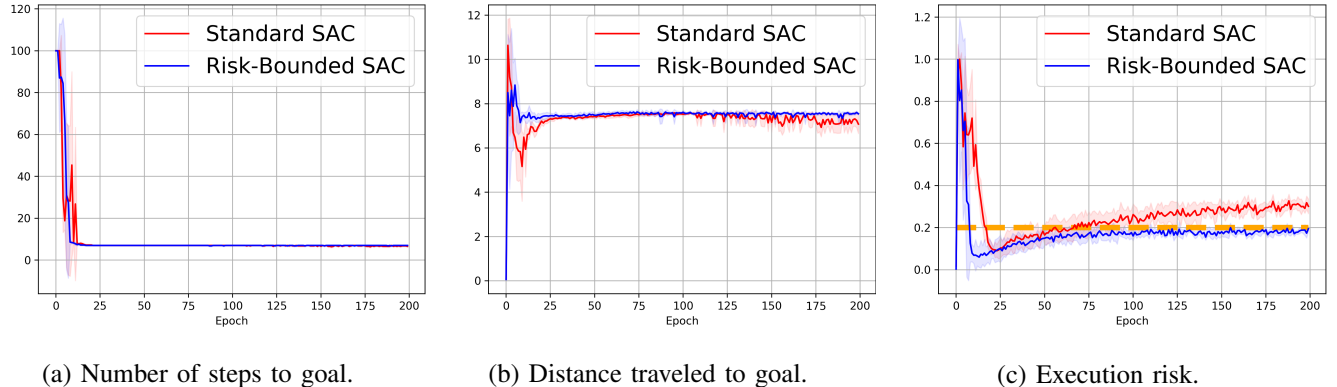| (a) Number of steps to goal. | (b) Distance traveled to goal. | (c) Execution risk. |

Fig. 2: Evolution of evaluation metrics over runs from 5 randoms seeds, between standard SAC and our model, given a fixed upper risk bound $\Delta = 0.2$ in *OneObstacle* maze environment. Shaded area represents one standard deviation. (a) Our model converges to the same number of time steps as SAC. (b) Our model converges to a distance slightly larger than SAC, due to additional risk constraints. (c) Our model converges to risk level below the risk bound, as visualized by the orange line.

multilayer perceptrons (MLPs) with hidden dimensions of 256, and two target $Q$ function estimators as 3-layer MLPs with the same structure. The policy network is composed of 3-layer MLPs with hidden dimensions of 256 and outputs a mean vector and a log standard deviation vector, which are further applied with an invertible squashing function (*tanh*) to generate the bounded distribution for the resulting action. In addition, our model includes a risk estimator with the same structure as the $Q$ function estimators.

Our model is implemented in PyTorch, based on the open-source rlkit library[2]. At training time, we train for 200 epochs with a learning rate of 3e-4, a batch size of 256, and a reward discount factor of 0.99. The coefficient of risk penalty term $\lambda_{er}$ is 10. At test time, we evaluate the path generated from the learned policy and obtain the evaluation metrics, including distance traveled to goal and execution risk. The execution risk is computed using Monte Carlo sampling methods with 500 samples. We train and test our model in the same environment, as customary in many reinforcement learning tasks despite being prune to overfitting [27].

*3) Risk-Bounded Policy with a Fixed Upper Risk Bound:* We start by validating the performance of our model given a fixed risk bound $\Delta = 0.2$ as input during training. This bound is selected arbitrarily – in practice, it can be specified based on risk-averse level or requirement for plan efficiency. In addition, we train a standard soft actor critic (SAC) method that naively minimizes the cost function without reasoning about risk.

The comparisons between our model and standard SAC over 5 random seeds are illustrated in Fig. 2. We note that our model has been trained successfully to generate policies subject to the risk bound, as visualized by the orange line in Fig. 2(c). At the same time, it produces the policy with the same number of time steps and slightly worse traveled

distance, compared to the standard SAC model. The worse traveled distance is due to the trade-off between performance and safety.

*4) Risk Conditioned Policy:* Next, we train our model by providing random risk bounds in the data so that it can learn risk-bounded policies conditioned on different upper risk bound levels. At test time, we visualize the planned trajectories by varying values for $\Delta$, as shown in Fig. 3. In the same figure, we visualize the risk-bounded policies generated by a state-of-the-art MILP-based solver, IRA, given the same upper risk bound levels and environment configurations. When the bounds become smaller, both models generate paths that are further away from the obstacle to improve safety. Both our model and IRA are run on the same hardware without GPU to make a fair comparison.

In Table I, we compare the distance traveled to goal and computational time between our method and IRA. The results show that on average, our method is able to reduce the computational time by 93.87%, demonstrating great time efficiency by using deep neural networks. Furthermore, it improves the plan quality, in terms of distance traveled to goal, by approximately 3.56%, compared to IRA that produces conservative plans by assigning the overall risk constraint into each time step through a union bound. The table also shows the actual execution risk for each plan, which verifies that our plans are bounded by the desired $\Delta$ values. Overall, our method successfully produces risk-bounded policies efficiently without sacrificing the plan quality.

### B. Dubins Car Model

*1) Problem Setup:* In this experiment, we want to demonstrate the performance of our model with nonlinear dynamics that are usually difficult to be modeled by a MILP, which is limited to linear systems. More specifically, we consider the
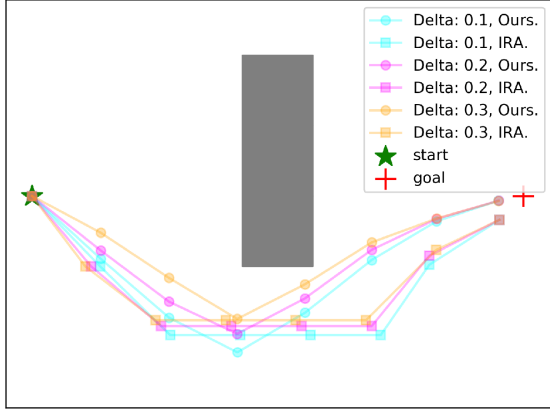
Fig. 3: Visualization of paths generated by our model (circles) and IRA (squares), under 3 risk bounds represented in different colors, in *OneObstacle* maze.

| $\Delta$ | IRA | | | Ours | | |
|---|---|---|---|---|---|---|
| | Distance[m] | Time[s] | Risk | Distance[m] | Time[s] | Risk |
| 0.1 | 8.20 | 0.16 | 0.024 | 8.04 | 0.0128 | 0.055 |
| 0.2 | 8.04 | 0.15 | 0.051 | 7.75 | 0.0085 | 0.125 |
| 0.3 | 7.93 | 0.16 | 0.071 | 7.52 | 0.0075 | 0.242 |

TABLE I: Comparison between our model and IRA in *OneObstacle* maze. Our model achieves better plan quality and computational time, subject to the risk bound.

Dubins car model [28] that is commonly used for real world ground robotics platforms [29]:

$$\dot{x} = v\cos(\theta), \quad \dot{y} = v\sin(\theta), \quad (17)$$

$$\dot{\theta} = u_\theta, \quad \dot{v} = u_v. \quad (18)$$

The test environment is a standard *TwoRooms* maze [30] that is composed of two rooms connected by two paths, as visualized in Fig. 4.

*2) Model Details:* We train a model with the same structure and parameters as in Sec. V-A.2, except the output from the policy network changes from velocities to angular velocities and accelerations for Dubins car model. The model is trained for 500 epochs with a risk penalty term of 20, due to the nonlinear dynamics. Despite the longer convergence time, our system achieves similar run-time complexity, by leveraging a model of the same size.

*3) Risk Conditioned Policy:* While it is usually infeasible to handle nonlinear dynamics in linear programs, we show that our model is capable of finding risk-bounded policies when the dynamics are nonlinear, as visualized in Fig. 4. The metrics are summarized in Table II. We notice that under different risk bounds (i.e. 0.2 and 0.3), the risk-bounded plans have similar distances, which can be explained by the fact that the distance does not strictly decrease when the risk bound increases, due to the nonlinearity of the dynamics model.

### C. FlyTrapBig Maze

*1) Problem Setup:* In this experiment, we validate our model in a larger maze, named *FlyTrapBig* [30], which has
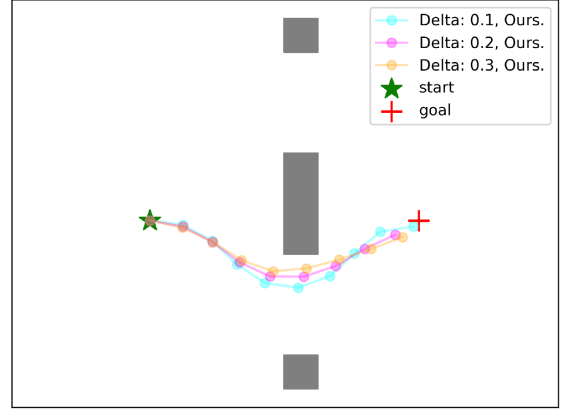


Fig. 4: Visualization of paths generated by our model with Dubins car model, under different risk bounds, in *TwoRooms*.
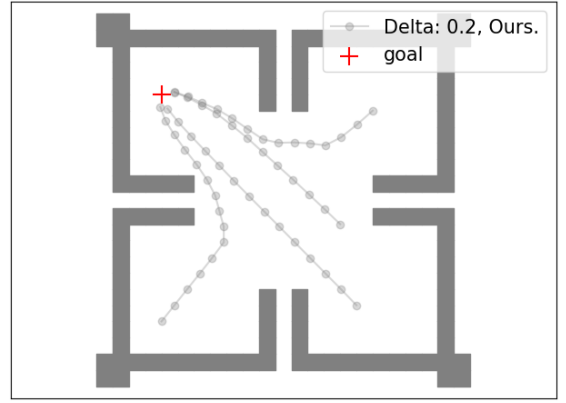


Fig. 5: Visualization of paths generated by our model in *FlyTrapBig* maze starting at different locations, under the same risk bound of 0.2. Our model successfully generalizes to random start locations in a large maze.

| $\Delta$ | Distance[m] | Time[s] | Risk |
|---|---|---|---|
| 0.1 | 9.00 | 0.0094 | 0.076 |
| 0.2 | 8.00 | 0.0071 | 0.142 |
| 0.3 | 8.00 | 0.0072 | 0.214 |

TABLE II: Dubins car model results in *TwoRooms*. Our model produces risk-bounded plans under different risk bounds.

four times larger state size compared to the *OneObstacle* maze. This requires more exploration in the space, and more importantly, requires more time steps to get to the goal, which poses great difficulty to the existing constrained MDP methods [7], [9]. This is because these methods need to allocate the global risk bound into many individual time steps, which could end up returning infeasible solutions. On the other hand, our model avoids conservatism by using an accurate definition of risk and generates feasible solutions even when the number of steps to the goal is large.

*2) Model Details:* We train a model with the same structure and parameters as in Sec. V-A.2, except with a larger epoch size of 1200 and a larger risk penalty coefficient $\lambda_{er}$ of 20, due to the larger state space and increased magnitude

of plan rewards, respectively.

*3) Risk Conditioned Policy:* We show that our model can generate risk-bounded policies starting at random positions in a large maze, as visualized in Fig. 5. When there exist obstacles between the start and the goal (i.e. when starting at the lower left or upper right room), our model outputs a risk-bounded path that balances between plan quality and safety.

## VI. CONCLUSIONS

In conclusion, we propose a risk conditioned soft actor critic model that generates risk-bounded policies in the context of chance-constrained MDPs. Our model leverages a risk critic to estimate the execution risk at a given state and acting according to the policy, and adds a penalty term to the policy loss if the estimated risk is greater than the risk bound to produce risk-bounded policies. By providing upper risk bound as part of the input, our method is able to generate plans with different risk-averse levels for the agent on the fly. We demonstrate the advantage of our model in terms of both time complexity and path quality compared to a MILP-based baseline in a simple maze environment, and further validate its performance with more complex agent dynamics and larger state space, which are usually hard to handle by MILP-based methods. Future work includes validating our model in more complicated settings (i.e. a larger state space, a variety of geometries, real-world experiments) and extending our algorithm to handle dynamic obstacles.

## VII. ACKNOWLEDGMENT

## REFERENCES

[1] A. Rajeswaran, V. Kumar, A. Gupta, G. Vezzani, J. Schulman, E. Todorov, and S. Levine, "Learning complex dexterous manipulation with deep reinforcement learning and demonstrations," in *Robotics: Science and Systems*, 2017.

[2] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot *et al.*, "Mastering the game of go with deep neural networks and tree search," *nature*, vol. 529, no. 7587, pp. 484–489, 2016.

[3] J. R. Birge and F. Louveaux, *Introduction to stochastic programming*. Springer Science & Business Media, 2011.

[4] E. Altman, *Constrained Markov decision processes*. CRC Press, 1999, vol. 7.

[5] M. Ono and B. C. Williams, "An efficient motion planning algorithm for stochastic dynamic systems with constraints on probability of failure." in *Association for the Advancement of Artificial Intelligence (AAAI)*, 2008, pp. 1376–1382.

[6] P. Santana, S. Thiébaux, and B. Williams, "Rao*: An algorithm for chance-constrained pomdp's," in *Association for the Advancement of Artificial Intelligence (AAAI)*, vol. 30, no. 1, 2016.

[7] J. Achiam, D. Held, A. Tamar, and P. Abbeel, "Constrained policy optimization," in *International Conference on Machine Learning*, 2017, pp. 22–31.

[8] M. Pfeiffer, S. Shukla, M. Turchetta, C. Cadena, A. Krause, R. Siegwart, and J. Nieto, "Reinforced imitation: Sample efficient deep reinforcement learning for mapless navigation by leveraging prior demonstrations," *IEEE Robotics and Automation Letters*, vol. 3, no. 4, pp. 4423–4430, 2018.

[9] M. Yu, Z. Yang, M. Kolar, and Z. Wang, "Convergent policy optimization for safe reinforcement learning," in *Conference on Neural Information Processing Systems (Neurips)*, 2019.

[10] E. A. Feinberg and A. Shwartz, "Constrained discounted dynamic programming," *Mathematics of Operations Research*, vol. 21, no. 4, pp. 922–945, 1996.

[11] D. A. Dolgov and E. H. Durfee, "Stationary deterministic policies for constrained mdps with multiple rewards, costs, and discount factors," in *IJCAI*, 2005, pp. 1326–1331.

[12] B. Luders, M. Kothari, and J. How, "Chance constrained rrt for probabilistic robustness to environmental uncertainty," in *AIAA guidance, navigation, and control conference*, 2010, p. 8160.

[13] X. Huang, A. Jasour, M. Deyo, A. Hofmann, and B. C. Williams, "Hybrid risk-aware conditional planning with applications in autonomous vehicles," in *2018 IEEE Conference on Decision and Control (CDC)*. IEEE, 2018, pp. 3608–3614.

[14] X. Huang, S. Hong, A. Hofmann, and B. C. Williams, "Online risk-bounded motion planning for autonomous vehicles in dynamic environments," in *Proceedings of the International Conference on Automated Planning and Scheduling*, vol. 29, 2019, pp. 214–222.

[15] S. Hong, S. U. Lee, X. Huang, M. Khonji, R. Alyassi, and B. C. Williams, "An anytime algorithm for chance constrained stochastic shortest path problems and its application to aircraft routing," in *International Conference on Robotics and Automation (ICRA)*, 2021.

[16] E. A. Feinberg, "Constrained discounted markov decision processes and hamiltonian cycles," *Mathematics of Operations Research*, vol. 25, no. 1, pp. 130–140, 2000.

[17] Y. Chow, M. Ghavamzadeh, L. Janson, and M. Pavone, "Risk-constrained reinforcement learning with percentile risk criteria," *The Journal of Machine Learning Research*, vol. 18, no. 1, pp. 6070–6120, 2017.

[18] O. Andersson, M. Wzorek, and P. Doherty, "Deep learning quadcopter control via risk-aware active learning," in *Association for the Advancement of Artificial Intelligence (AAAI)*, vol. 31, no. 1, 2017.

[19] H. Krasowski, X. Wang, and M. Althoff, "Safe reinforcement learning for autonomous lane changing using set-based prediction," in *2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2020, pp. 1–7.

[20] Y. S. Shao, C. Chao, S. Kousik, and R. Vasudevan, "Reachability-based trajectory safeguard (rts): A safe and fast reinforcement learning safety layer for continuous control," *arXiv preprint arXiv:2011.08421*, 2020.

[21] A. Levy, G. Konidaris, R. Platt, and K. Saenko, "Learning multi-level hierarchies with hindsight," in *International Conference on Learning Representations*, 2018.

[22] D. Ghosh, A. Gupta, and S. Levine, "Learning actionable representations with goal conditioned policies," in *International Conference on Learning Representations*, 2018.

[23] S. Nasiriany, V. H. Pong, S. Lin, and S. Levine, "Planning with goal-conditioned policies," in *Conference on Neural Information Processing Systems (Neurips)*, 2019.

[24] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," in *International Conference on Machine Learning*. PMLR, 2018, pp. 1861–1870.

[25] E. Schmerling and M. Pavone, "Evaluating trajectory collision probability through adaptive importance sampling for safe motion planning," in *Robotics: Science and Systems*, 2017.

[26] A. Wang, X. Huang, A. Jasour, and B. Williams, "Fast risk assessment for autonomous vehicles using learned models of agent futures," in *Robotics: Science and Systems*, 2020.

[27] K. Cobbe, O. Klimov, C. Hesse, T. Kim, and J. Schulman, "Quantifying generalization in reinforcement learning," in *International Conference on Machine Learning*. PMLR, 2019, pp. 1282–1289.

[28] L. E. Dubins, "On curves of minimal length with a constraint on average curvature, and with prescribed initial and terminal positions and tangents," *American Journal of mathematics*, vol. 79, no. 3, pp. 497–516, 1957.

[29] P. R. Giordano and M. Vendittelli, "Shortest paths to obstacles for a polygonal dubins car," *IEEE Transactions on Robotics*, vol. 25, no. 5, pp. 1184–1191, 2009.

[30] B. Eysenbach, R. Salakhutdinov, and S. Levine, "Search on the replay buffer: Bridging planning and reinforcement learning," in *Conference on Neural Information Processing Systems (Neurips)*, 2019.