

# MADR: MPC-guided Adversarial DeepReach

Ryan Teoh<sup>1,\*</sup>, Sander Tonkens<sup>2,\*</sup>, William Sharpless<sup>2</sup>, Aijia Yang<sup>2</sup>,  
Zeyuan Feng<sup>3</sup>, Somil Bansal<sup>3</sup>, and Sylvia Herbert<sup>2</sup>

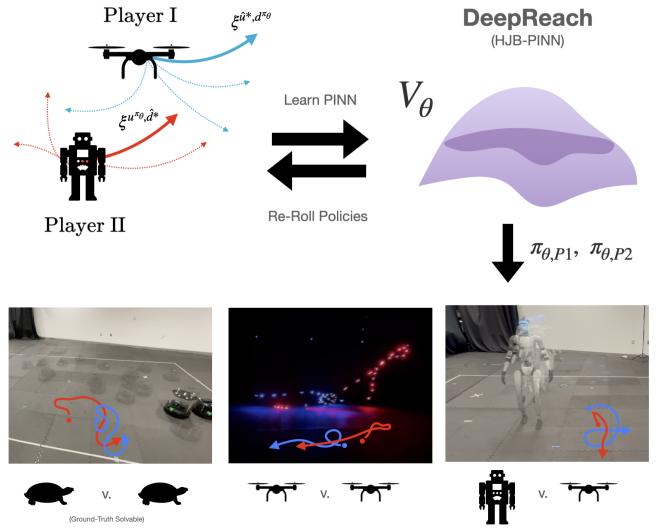
**Abstract**— Hamilton-Jacobi (HJ) Reachability offers a framework for generating safe value functions and policies in the face of adversarial disturbance, but is limited by the curse of dimensionality. Physics-informed deep learning is able to overcome this infeasibility, but itself suffers from slow and inaccurate convergence, primarily due to weak PDE gradients and the complexity of self-supervised learning. A few works, recently, have demonstrated that enriching the self-supervision process with regular supervision (based on the nature of the optimal control problem), greatly accelerates convergence and solution quality, however, these have been limited to single player problems and simple games. In this work, we introduce MADR: MPC-guided Adversarial DeepReach, a general framework to robustly approximate the two-player, zero-sum differential game value function. In doing so, MADR yields the corresponding optimal strategies for both players in zero-sum games as well as safe policies for worst-case robustness. We test MADR on a multitude of high-dimensional simulated and real robotic agents with varying dynamics and games, finding that our approach significantly out-performs state-of-the-art baselines in simulation and produces impressive results in hardware.

## I. INTRODUCTION

Zero-sum differential games provide a unifying framework for modeling interactions between an ego robot and an adversarial entity—whether that entity is a passive disturbance (e.g., wind, sensor noise) or an active agent (e.g., an opponent robot or attacker). To ensure robustness, zero-sum formulations assume that the other player or disturbance is adversarial, i.e. the ego robot aims to maximize a reward function (e.g., time to target, safety margin), while the disturbance or adversary attempts to minimize it.

Hamilton-Jacobi (HJ) Reachability is a popular approach for solving zero-sum differential games with respect to target achievement and obstacle avoidance [1], [2]. HJ reachability requires solving a partial differential equation (HJ-PDE) to compute a value function; this value function implicitly provides a) the optimal control policy for the robot and adversary, and b) the reachable set: the set of states from which a system can reach (or avoid) a target (or obstacle) over time using these optimal policies. HJ reachability is particularly well-suited to safety-critical applications due to its ability to model worst-case scenarios in a game-theoretic formulation [?].

HJ reachability has been applied to a broad range of robotic systems, including collision avoidance for aerial vehicles, multi-agent coordination, and motion planning under uncertainty [1], [3]. However, classical approaches rely on



**Fig. 1:** A graphical abstract of MADR and the robotic experiments. In this work, we propose enriching self-supervised learning of HJ-PDE's, i.e. DeepReach, with supervision given by the best sampled game rollout (top left), where the opponents policy is defined by the current value approximation (top right). We demonstrate this approach performs across games with general dynamics (bottom), namely with TurtleBots, Drones and Humanoid experiments.

grid-based dynamic programming and scale poorly with state dimension due to the curse of dimensionality [4]. This limits their use to systems with fewer than six states in practice.

To address this challenge, learning-based approaches have been explored for scalability. Methods such as DeepReach [5] replace the grid with a function approximator trained to satisfy the HJ-PDE in a physics-informed learning (PINN) framework. This approach significantly improves scalability and has been demonstrated on systems of 50 dimensions [6].

DeepReach has primarily focused on control-only cases at high dimensions, with only limited work addressing robustness or more complex dynamics [7], [8]. A few works have demonstrated that introducing supervision greatly improves the learned approximation. In [6], the Hopf formula is used to solve a linearization of the game for a supervision loss, however, this struggles with nonlinear or angular dimensions, which frequently appear in robotics. More recently, [9] employed sample-based Model Predictive Control (MPC) as supervisory signal, finding that this significantly improves performance in high-dimensions but is limited to the control-only case.

<sup>1</sup>University of California, Los Angeles, <sup>2</sup>University of California, San Diego, <sup>3</sup>Stanford University. \*Authors contributed equally.

### A. Contributions

Our method, MADR, proposes solving zero-sum differential games using a value-only DeepReach approach combined with adversarial MPC roll-outs. Our key insight is to define the opponent's policy through the current value gradient approximation, which allows robust sampling of improved ego policies and the associated optimal cost of their trajectories. Unlike actor-critic methods, where policy performance depends on simultaneous learning of both actor and critic, our value-only approach ensures that high-quality policies emerge directly from accurately-learned value functions. To achieve this, we mitigate co-learning issues by collecting separate datasets in which the sampling (and associated adversarial gameplay under the current value gradient) occurs for each agent individually. This enables robust learning of a single value function for both players, leading to effective policies in adversarial scenarios.

We empirically demonstrate that this value-informed adversarial supervision improves both the fidelity and safety margins of learned reachable sets, leading to better robustness against disturbances and adversarial agents. The effectiveness of MADR is validated through extensive simulation experiments and a variety of hardware demonstrations, highlighting improved safety performance in complex, high-dimensional systems exposed to significant disturbances or adversarial players. In addition, we compare MADR against other current approaches in head-to-head policy matchup comparisons. Ultimately, our approach bridges the gap between the rigorous, theory-driven framework of HJ reachability and the practical use of optimal control to handle adversarial or worst-case scenarios in real-world environments.

### B. Related Work

**Linear Quadratic Games (LQG)** are a class of differential games where multiple players aim to minimize a quadratic cost function over time, subject to linear dynamics, either cooperatively or non-cooperatively. In the non-cooperative scenario, players often seek Nash equilibria, which can be computed by solving coupled Riccati equations [10] in limited scenarios, but practically in iteratively linearizations [11]. The cooperative case typically involves joint optimization to minimize a shared cost function [12]. Recent advancements have extended LQGs to incorporate uncertainties and stochastic disturbances, enhancing their applicability to real-world scenarios [13]. However, they remain relatively challenging in zero-sum settings and face limitations when applied to complex dynamical systems, where linearization of the dynamics can lead to significant inaccuracies.

**Model Predictive Control (MPC)** is a model-based optimal control strategy. At each timestep, it solves a constrained optimization problem over a finite prediction horizon, a method known as receding horizon control. MPC is well-suited for safety-critical applications because it produces a control sequence that is guaranteed by the model to satisfy state and input constraints over the entire prediction horizon [14]. However, MPC has two significant drawbacks: its

high computational burden and its sensitivity to model mismatch. Sampling-based variants, such as MPPI [15], improve robustness to model inaccuracies by evaluating thousands of stochastic rollouts. Additionally, robust, stochastic, and distributed formulations further enhance performance and efficiency in uncertain environments [16].

**Adversarial Reinforcement Learning (RL)** formulates control in uncertain environments as a two-player game, where a controller learns to counteract worst-case disturbances or adversarial agents [17], [18]. This framework improves robustness to unmodeled dynamics, environmental perturbations, and intelligent opponents. Modern algorithms, often based on actor-critic methods, enable scalable training in high-dimensional systems. However, these approaches are often brittle, as adversarial training can induce instability, high variance, or convergence issues [19]. Additionally, the reliance on additive payoffs in these methods makes them ill-suited for safety-critical tasks, thus motivating the integration of reachability-based value functions for learning in such domains [20], [21].

**ISAACS** (Iterative Soft Adversarial Actor-Critic for Safety) is a Safety Bellman based representative framework that formulates safety-critical control as a zero-sum game between a control policy and a disturbance policy [22], [23], [24]. The disturbance acts as an adversary attempting to drive the system into unsafe regions, while the controller learns to maintain safety. ISAACS trains both players jointly via soft actor-critic methods and leverages the learned value function to implement a runtime safety filter. If the current task policy selects an unsafe action, the safety filter can override it with a robust fallback generated from the learned safety policy. This framework enables real-time safety assurance in high-dimensional systems. However, ISAACS relies on actor-critic style sampling to learn the value function, which can be unstable or sample-inefficient in complex domains. Additionally, it prioritizes robustness against adversarial disturbances by optimizing the disturbance policy given a control policy, and has primarily been applied to worst-case disturbance scenarios, rather than general two-player games.

**Organization:** We begin in Section II with a brief overview of Hamilton-Jacobi (HJ) reachability analysis and the computation of backward reachable tubes (BRTs). In Section III, we introduce the proposed MADR: MPC-guided Adversarial DeepReach framework and detail the algorithmic approach for learning reachability solutions under adversarial MPC supervision. Section IV presents simulation results demonstrating the effectiveness of MADR, followed by hardware deployment results in Section V. Finally, we summarize our findings and contributions in Section VI.

## II. PRELIMINARIES

### A. Problem Formulation

We consider a dynamical system with state  $x \in \mathbb{R}^n$ , control input  $u \in \mathcal{U}$ , and disturbance (or adversary)  $d \in \mathcal{D}$ , governed by control-and-disturbance affine dynamics

$$\dot{\xi}(t) = f(\xi(t)) + g(\xi(t))\mathbf{u}(t) + w(\xi(t))\mathbf{d}(t), \quad (1)$$

where  $\mathbf{u} : \mathbb{T} \rightarrow \mathcal{U}$  and  $\mathbf{d} : \mathbb{T} \rightarrow \mathcal{D}$  are time-varying control and disturbance signals. Here,  $\mathcal{U}$  and  $\mathcal{D}$  denote compact sets of admissible control and disturbance inputs, defined by physical or operational constraints. Let Player I choose  $u$  while Player II chooses  $d$ .

In this setting, Player I seeks to avoid entering a *failure set*  $\mathcal{F} \subset \mathbb{R}^n$ , despite the opposing actions of Player II. This failure set is the set of states that Player I must avoid to remain safe. This can include situations such as being captured by Player II, or colliding with obstacles. Conversely, Player II acts adversarially to drive the system toward  $\mathcal{F}$ . This interaction defines an *avoid game*, where the objective of Player I is to maintain safety over a finite time horizon  $\mathbb{T} = [t, T]$ .

Throughout this work, we solve such games by characterizing the *safe set*  $\mathcal{S} \subset \mathbb{R}^n$ : the set of states from which Player I can guarantee safety against the worst-case disturbance within the specified horizon.

### B. Hamilton–Jacobi Formulation

To characterize the safe set and corresponding robust control policy, we begin by defining a *boundary function* that defines the failure set  $\mathcal{F} \subset \mathbb{R}^n = \{x : \ell(x) \leq 0\}$ . This is typically constructed as a signed distance function to the failure set. The objective function is defined as

$$J(x, t, \mathbf{u}, \mathbf{d}) = \min_{s \in [t, T]} \ell(\xi_{x,t}^{\mathbf{u}, \mathbf{d}}(s)), \quad (2)$$

where  $\xi_{x,t}^{\mathbf{u}, \mathbf{d}}(\cdot)$  denotes the system trajectory starting from  $x$  at time  $t$  and state  $x$ .

The *value function*  $V(x, t)$  encodes the safety of each state under worst-case disturbances and optimal control by optimizing over (2). It is defined as

$$V(x, t) = \min_{\mathbf{d}[\mathbf{u}] \in \Xi} \max_{\mathbf{u} \in U} J(x, t, \mathbf{u}, \mathbf{d}), \quad (3)$$

where  $\mathbf{d}[\mathbf{u}] : U \rightarrow D$  denotes a nonanticipative strategy such that Player II may choose to play  $d$  based on the current and past values of the control  $\mathbf{u}$ , but not on its future values.

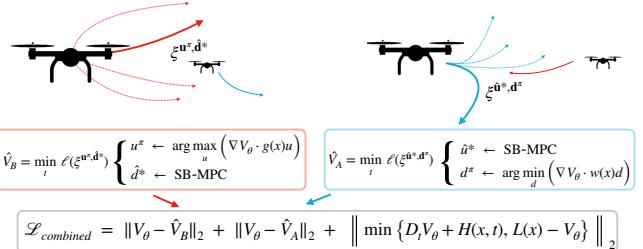
A positive value  $V(x, t) > 0$  indicates that the robot starting from state  $x$  at time  $t$  can avoid entering the failure set under all admissible disturbances. Therefore, the *safe set* at time  $t$  is given by

$$\mathcal{S}(t) = \{x \in \mathbb{R}^n \mid V(x, t) > 0\}. \quad (4)$$

The value function is computed by solving a Hamilton–Jacobi–Isaacs variational inequality (HJI–VI):

$$\begin{aligned} \min \left\{ \frac{\partial V}{\partial t}(x, t) + \max_{u \in U} \min_{d \in D} \nabla V(x, t) \cdot f(x), \right. \\ \left. \ell(x) - V(x, t) \right\} = 0, \quad V(x, T) = \ell(x). \end{aligned} \quad (5)$$

The variational inequality ensures that once a state enters the failure set, it remains marked as unsafe for all earlier times, thus correctly capturing the avoid for all time game objective. Grid-based methods can compute the exact value function but quickly become intractable in high-dimensional systems. To address this, we use a semi-supervised learning



**Fig. 2:** Graphical depiction of the MPC formulation combining both player’s rollouts in the loss function for the proposed MPG-guided adversarial PINN training.

approach (see Section III) to approximate the value function, preserving the key properties of the BRT while remaining computationally feasible.

### C. Online Control Using the BRT

Once the value function is computed offline, its gradients can be used online by each player to determine optimal control actions. Specifically, the evader (safety-maximizing player) chooses  $u^*(x, t) = \arg \max_{u \in \mathcal{U}} \nabla V(x, t) \cdot g(x)u$ .

The disturbance / pursuer (adversary) chooses  $d^*(x, t) = \arg \min_{d \in \mathcal{D}} \nabla V(x, t) \cdot w(x)d$ , with  $g(x)$  and  $w(x)$  from (1). The sign of the gradient  $\nabla_x V$  along these directions determines whether increasing or decreasing the input increases safety. These rules allow real-time control by evaluating the gradient of the precomputed value function, avoiding the need to solve HJ-PDE describing the optimal control problem online. Optionally, this value function can act directly as a *safety filter* [25], constraining control inputs to prevent entering unsafe regions under disturbances.

## III. ADVERSARIAL MPC-GUIDED REACHABILITY LEARNING FRAMEWORK

Building on the Hamilton–Jacobi reachability concepts introduced above, we present a framework for learning the value function in high-dimensional systems subject to adversarial inputs. Rather than solving (5) over a grid, our approach leverages self-supervised learning of the HJI–VI loss, enhanced with a supervised adversarial sampling-based MPC procedure that is generated in parallel to efficiently approximate the value function. We will first introduce the sampling-based MPC method that optimizes the objective (2). Then we describe how this dataset will be used to augment the training loss of DeepReach. These losses are combined, as illustrated in Fig 2, to guide the learned value function toward robust, safety-informed behavior in adversarial scenarios.

### A. Sampling-Based MPC Dataset

To generate the MPC dataset, we solve the following discrete-time version of the zero-sum game:

$$\begin{aligned} \hat{V}(x, t) = \min_{\mathbf{d}} \max_{\mathbf{u}} \min_{h \in \{0, 1, \dots, H\}} \ell(\xi_h) \\ \text{s.t. } \xi_{h+1} = \bar{f}(\xi_h) + \bar{g}(\xi_h, u_h) + \bar{w}(\xi_h, d_h), \\ \xi_0 = x, \quad u_h \in U, \quad d_h \in D. \end{aligned} \quad (6)$$

**Algorithm 1** Sampling-based MPC dataset [Control Perspective] [Disturbance Perspective]

---

**Input:** MPC dataset size  $|D_{\text{MPC}}|$ , horizon  $H$ , step size  $\Delta t$ , dynamics model  $f$ , learned value function  $V_\theta$ , constraint function  $\ell$ , sample size  $N$ , refinement iterations  $K$

**Output:**  $D_{\text{MPC}} = \{(x_j, t_j, \hat{V}(x_j, t_j))\}$

**Initialization:** Set best cost  $J^* \leftarrow -\infty$   $J^{(i)} \leftarrow \infty$   $D_{\text{MPC}} = \emptyset$

**for**  $j \leftarrow 1$  **to**  $|D_{\text{MPC}}|$  **do**

$x_j \sim \text{Uniform}(\mathcal{X})$ ,  $t_j \sim \text{Uniform}(0, T)$

**for**  $k \leftarrow 1$  **to**  $K$  **do**

**for**  $i \leftarrow 1$  **to**  $N$  **do**

$x_0^{(i)} \leftarrow x_j$

**for**  $h \leftarrow 0$  **to**  $H - 1$  **do**

$u^{(i)} \sim \mathcal{N}(\mu_u, \sigma_u^2)$

$d^{(i)} \leftarrow \arg \min_d \nabla V_\theta(x_h^{(i)}) \cdot w(x_h^{(i)}) d$

$u^{(i)} \leftarrow \arg \max_u \nabla V_\theta(x_h^{(i)}) \cdot g(x_h^{(i)}) u$

$d^{(i)} \sim \mathcal{N}(\mu_d, \sigma_d^2)$

$x_{h+1}^{(i)} = x_h^{(i)} +$

$\bar{f}(x_h^{(i)}) \Delta + \bar{g}(x_h^{(i)}, u^{(i)}) \Delta + \bar{w}(x_h^{(i)}, d^{(i)}) \Delta$

$J^{(i)} = \min_h \ell(x_h^{(i)})$  **if**  $H < \mathbb{T}$  **then**

$V_{\text{togo}}^{(i)} \leftarrow V_\theta(x_H^{(i)})$   $J^{(i)} \leftarrow \min\{J^{(i)}, V_{\text{togo}}^{(i)}\}$

**Update best cost:** **if**  $J^{(i)} > J^*$   $J^{(i)} < J^*$  **then**

$\mu_u \leftarrow u^{(i)}$   $\mu_d \leftarrow d^{(i)}$   $J^* \leftarrow J^{(i)}$

$\hat{V}(t_j, x_j) \leftarrow J^*$

**return**  $D_{\text{MPC}} \leftarrow D_{\text{MPC}} \cup (x_j, t_j, \hat{V}(x_j, t_j))$

---

Here  $h \in \{0, 1, \dots, H\}$  denotes the time steps between  $t$  and  $T$ ,  $\xi$  is the system trajectory,  $\mathbf{u} := [u_0, \dots, u_H]$  is the control sequence, and  $\mathbf{d} := [d_0, \dots, d_H]$  is the disturbance sequence. The functions  $\bar{f}, \bar{g}, \bar{w}$  are the discretized dynamics that can be obtained from the continuous dynamics using first-order Euler approximation. This game can be solved in many ways; we will use sampling-based MPC (SB-MPC). Directly solving the resulting two-player game using MPC is challenging due to the need for co-optimization over both control and disturbance sequences. Instead, we generate two distinct datasets—one with SB-MPC controlling the agent while the disturbance follows a policy, and one with the roles reversed—allowing efficient learning of the value function. In practice, we generate two complementary datasets: 1) SB-MPC for control combined with policy-driven disturbance rollouts, and 2) policy-driven control combined with SB-MPC disturbance rollouts. Our MPC algorithm is summarized in Alg. 1 and described below.

1) *Control Dataset:* We first review how to generate the MPC dataset for the control input, shown in black and blue text. The algorithm takes as inputs the MPC dataset size,  $|D_{\text{MPC}}|$ , time horizon  $H$  and step size  $\Delta t$ , dynamics model  $f$ , current learned value function  $V_\theta$ , constraint function  $\ell$ , number of trajectories  $N$  and number of refinement iterations  $K$ . We first uniformly sample states from the state space

$\mathcal{X}$  and times from the time horizon  $\mathbb{T}$ . Over the course of  $K$  refinement iterations, we then generate  $N$  trajectories starting at the sampled initial state  $x_0$ . For each timestep  $\Delta t$  over a horizon  $H$ , we first generate the control input by sampling a Gaussian distribution around a nominal mean control  $\mu_u$  and variance  $\sigma_u^2$ , both defined by the user. In practice, this is typically set to  $\mu_u = 0$  and the variance is set to the maximum control bound. The robot then samples the disturbance by taking the action that minimizes the gradient of the current learned value function. The state, control, and disturbance are then propagated over  $\Delta t$  via the dynamics.

Next, the objective function is set as the minimum of the constraint function over the trajectory. If the horizon  $H$  is less than the total time horizon, the robot takes the min between the objective and the value function evaluated at the end of the horizon  $H$  (i.e. cost-to-go). Finally, we update the best solution by making the new nominal control the best control action that we sampled so far. After all refinement steps, the value function estimate at the initial state is set to the highest cost of all rollouts.

2) *Disturbance Dataset:* To generate the MPC dataset for the disturbance, we follow the same algorithm but using the red text rather than the blue. The disturbance is now generated by sampling from a Gaussian distribution around a nominal mean disturbance  $\mu_d$  and variance  $\sigma_d^2$ , also defined by the user. The robot then samples the control by taking the action that maximizes the gradient of the current learned value function. At the end of all refinement steps, the initial state's value function's estimate is set to the lowest cost of all rollouts.

These two datasets provide supervised learning signals for DeepReach with the following loss for both:

$$\mathcal{L}_{\text{MPC}} = \sum_{j=1}^{|D_{\text{MPC}}|} l_{\text{MPC}}(x_j, t_j, \hat{V}(x_j, t_j); \theta) \quad (7)$$

$$l_{\text{MPC}}(x, t, \hat{V}(x, t); \theta) = \|\hat{V}(x, t) - V_\theta(x, t)\|_2,$$

with  $\mathcal{L}_{\text{MPC}} = \mathcal{L}_{\text{MPC}, \mathbf{u}} + \mathcal{L}_{\text{MPC}, \mathbf{d}}$  combining the sampling-based rollouts from the control and disturbance perspective.

*Remark 1:* In addition to sampling-based MPC, standard Model Predictive Path Integral control (MPPI) could also be used to generate these datasets by sampling stochastic rollouts. Within our codebase, the user can select either SB-MPC or regular MPPI as the guiding policy; however, in practice, we use SB-MPC, which assigns full weight to the best rollout to approximate optimal control under uncertainty.

## B. Learning Value Function using MPC-Guided DeepReach

To guide the learning of the value function, we directly use lines 5-16 of Algorithm 2 in [9]. This approach leverages a dataset of MPC rollouts together with PDE-based supervision to train a value function network that is amenable to disturbances. We briefly review key components of [9], Algorithm 2, highlighting key differences.

Firstly, as the MPC dataset collection, Alg. 1, relies on the gradient of the estimated value function we do not consider

a pretraining phase and collect the first MPC dataset after a set period of the curriculum training.

Next, Algorithm 2 of [9] relies on collecting data iteratively over a horizon that terminates at the current time in the curriculum  $t_{\text{curr}}$  (hence starts from a time further in the curriculum). We employ the same technique but fix the evaluation of the learned value function’s gradient to be evaluated at time  $t_{\text{curr}}$  to remain in distribution with respect to the curriculum training.

### C. A special pursuit-evasion filter for suboptimality in long-horizon games

For two-player games with equally equipped agents, the value function is positive (safe for the evader) for the majority of the state space. Specifically, the pursuer’s objective is to minimize the minimum cost over a pre-specified horizon, assuming an optimal evader. As such, as the time-horizon increases, and the error of a learned system compounds, the approximated solution quality tends to deteriorate and erode pursuer performance.

Therefore, we propose endowing the pursuer with a second policy as a filter, which prioritizes staying close to the evader when unable to “catch” the evader and switch to the classic pursuit-evasion policy when a “catch” is achievable. The “following” game has the following cost function:

$$J_{\text{follow}}(x, t, \mathbf{u}, \mathbf{d}) = \max_{s \in [t, T]} \ell(\xi_{x,t}^{\mathbf{u}, \mathbf{d}}(s)), \quad (8)$$

and the associated value function is:

$$V_{\text{follow}}(x, t) = \min_{\mathbf{d}[\mathbf{u}] \in \Xi} \max_{\mathbf{u} \in U} J_{\text{follow}}(x, t, \mathbf{u}, \mathbf{d}). \quad (9)$$

This formulation is akin to [3], which considers a worst-case tracking bound between two systems. While the resulting policy for the evader is not necessarily performant (e.g., it might be optimal under this policy for the evader to reach a state  $x$  with  $\ell(x) < 0$  during the trajectory if this trajectory achieves a higher cost  $J$  over the full trajectory), the policy for the pursuer tries to minimize the maximum boundary function, i.e. distance, over a trajectory. Such a policy induces the pursuer to stay close to the evader, but does not prioritize “catching” the evader. As such, we introduce a filtered strategy for the pursuer:

$$d(x) = \begin{cases} \arg \max_d H_d & \text{if } |H_d| \geq \epsilon \\ \arg \max_d H_{d, \text{follow}}, & \end{cases} \quad (10)$$

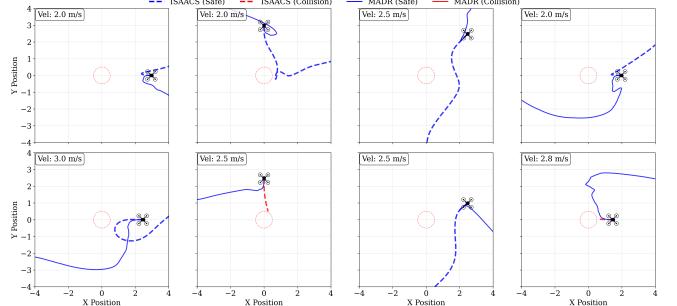
where  $H_d = \nabla V_{\theta}(x, T) \cdot w(x)d$  and  $\epsilon \ll 1$ . We denote this policy as MADR-FOLLOW for the pursuer.

## IV. SIMULATION RESULTS

### A. Simulation Performance Evaluation Metrics

We benchmark our proposed approach against a set of baselines and evaluate performance across all case studies. The primary baselines and evaluation aspects include:

**Baselines:** Vanilla DeepReach [5], DP Grid-Based Methods [?], ISAACS [22], sampling-based MPC [26], and our Adversarial MPC-Guided Policies (MADR).



**Fig. 3:** Representative trajectory rollouts comparing MADR and ISAACS in the 13D quadrotor environment under wind disturbances. The solid lines correspond to MADR, while the dashed lines correspond to ISAACS.

**Metrics:** Recovered volume and IOU with ground truth for safe/unsafe regions and policy matchup tables: time-to-capture and capture rate under adversarial disturbances for evader/pursuer policy pairs.

These metrics enable a comprehensive assessment of the strength of MADR, demonstrated by highly comparable performance to the ground truth solution (for low-dimensional settings) and outcompeting existing baseline’s learned policies across all experiments.

### B. Implementation Details

For all baselines, we employ a three-layer sinusoidal neural network, a standard architecture in machine learning for capturing complex, high-frequency functions, with 512 neurons per layer. The networks are trained using the Adam optimizer with a learning rate of  $\alpha = 2 \times 10^{-5}$  on an NVIDIA GeForce RTX 4090 GPU. The sampling-based MPC method uses a time step of  $\Delta = 0.02s$  and 10 iterative sampling steps, with  $N = 100$  perturbed control sequences generated per step. The fine-tuning loss weight  $\lambda_{\text{FT}} = 100$  and the dataset refinement horizon  $H_R = 0.2s$  for all experiments.

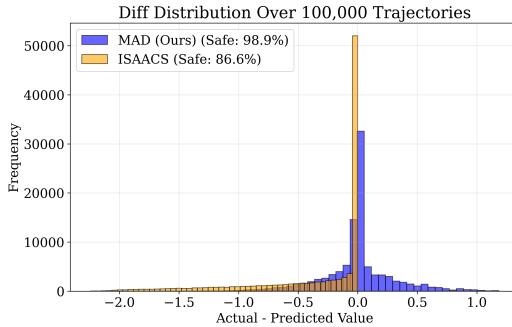
In each case study, we use the same number of training iterations across all baselines to ensure a fair comparison. However, Vanilla DeepReach does not involve a fine-tuning phase, see [9], as we observe that fine-tuning encourages overfitting in these baselines and degrades their performance. Detailed training parameters are provided in Table I.

**TABLE I:** Training time, number of epochs, and MPC dataset size for MADR across different systems.

Baseline / System	Training Time (hrs)	Epochs	MPC Dataset Size
Quadrotor (13D)	4.2	110,000	10,000
Dubins (6D)	3.6	150,000	20,000
Drones (20D)	5.4	250,000	30,000

### C. 13D Quadrotor with Disturbances

We first consider a robustness scenario, simulating a full 13-dimensional, quaternion-based quadrotor under wind disturbances, including forces along the  $x$ ,  $y$ , and  $z$  axes and torques about the rotational axes. The drone must avoid a central pillar of radius 0.50 in the state space despite a high initial velocity.

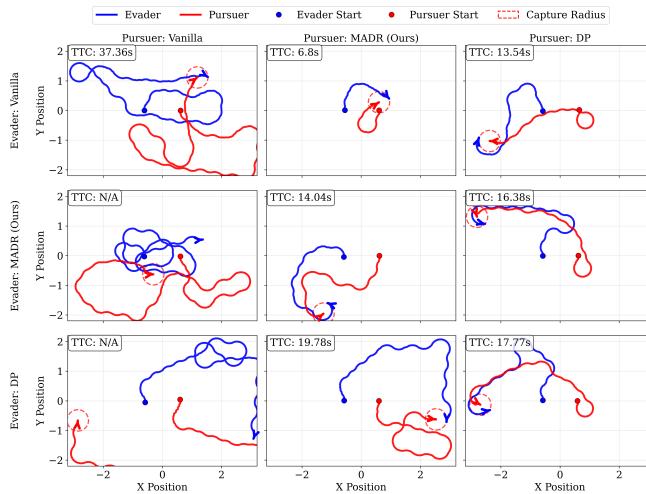


**Fig. 4:** Distribution of actual–predicted cost differences and safe rates over 100,000 trajectories, comparing MADR and ISAACS based on their value functions. Negative values indicate overestimation of safety.

Figure 3 shows representative trajectory rollouts as the drone approaches the cylinder at high initial speeds, illustrating the effectiveness of MADR. Even under worst-case disturbances, our method consistently avoids the central pillar, whereas ISAACS fails in some cases. Figure 4 presents the distribution of actual versus predicted trajectory costs for ISAACS and MADR, highlighting each method’s safety estimation. MADR achieves a 98.9% safe rate compared with 86.6% for ISAACS and provides significantly more accurate value function estimates, rarely overestimating safety.

#### D. 6D Dubins Pursuit Evasion Game

We then consider a zero-sum game for which we can obtain a ground truth comparison using dynamic programming (DP). Namely, we consider a Dubins pursuit–evasion game in which each player is described by three states:  $x$ ,  $y$ , and  $\theta$ . The evader aims to escape the pursuer, which has a capture radius of 0.36, based on TurtleBot size. The players are constrained to a bounded state space of  $[-3, 3]$  in  $x$  and  $[-2, 2]$  in  $y$ . Both players have a maximum turn rate of 1.9 rad/s and move at a constant velocity of 0.5 m/s. All methods are trained for a time horizon of 3 seconds.



**Fig. 5:** TurtleBot hardware trajectories from one initial condition, showing qualitative differences across policies. Non-capture trajectories are truncated.

Capture rates for each evader–pursuer policy pairing are summarized in Table III. As shown in the table, MADR performs near DP-optimal, significantly outperforming Vanilla DeepReach and ISAACS. Additionally, Table II compares each method’s unsafe set volume and Intersection-over-Union (IOU) with the ground-truth value function, demonstrating that MADR comes within 0.03% of the ground-truth solution. It maintains a slightly conservative unsafe set while outperforming ISAACS and Vanilla DeepReach in both accuracy and safe estimation.

**TABLE II:** Comparison of methods: unsafe set volume within  $x \in [-1.5, 1.5]$ ,  $y \in [-1.5, 1.5]$  and Intersection-over-Union (IOU) with DP ground truth. Higher IOU indicates better recovery of the true backward reachable tube (BRT).

Method	Unsafe Set Volume (%)	IOU with Ground Truth (DP)
DP (True BRT)	7.51	1.000
MADR (Ours)	7.82	0.997
Vanilla	4.60	0.969
ISAACS	0.37	0.928

**TABLE III:** Capture rates (%) over 100 safe initial states ( $0 < V < 0.1$ ) based on the DP BRT ground truth. Rows correspond to Player I’s policy and columns to Player II’s policy. Cells are color-coded to highlight performance: blue indicates a stronger Evader, while red indicates a stronger Pursuer.

Pursuer ( $\uparrow$ )		DP	MADR-FOLLOW	MADR	Vanilla	Isaacs
Evader ( $\downarrow$ )	Pursuer ( $\uparrow$ )	DP	MADR-FOLLOW	MADR	Vanilla	Isaacs
DP	DP	10	6	10	0	6
MADR	MADR	17	13	17	1	7
Vanilla	Vanilla	71	70	58	45	63
Isaacs	Isaacs	60	56	55	20	32

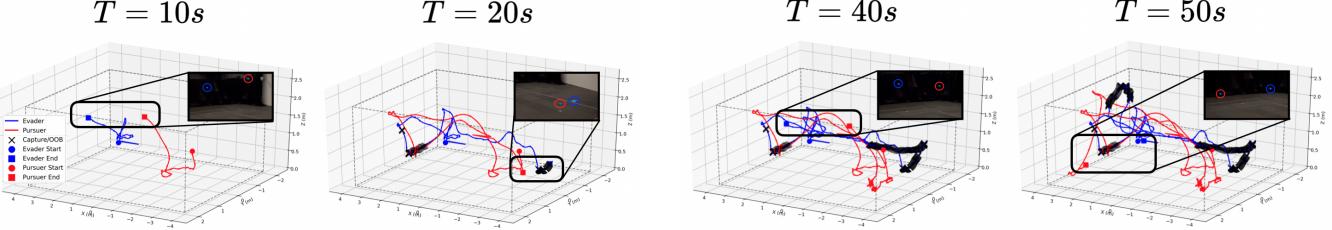
#### E. 20D Drone Pursuit–Evasion Game

In this scenario, we consider a drone pursuit–evasion game, where each drone is described by 10 states as detailed in [8]. The evader’s capture boundary is defined by a halfellipse centered under the pursuer to represent the evader trying to avoid the downwash from the pursuer, similar to [27]. All policies are trained over a 1-second time horizon. The players are constrained to a bounded state space of  $x \in [-4, 4]$ ,  $y \in [-2, 2]$ , and  $z \in [0.2, 2.0]$ . Each drone has identical dynamics, with a maximum desired pitch/roll angle of 0.15 rad, maximum thrust of 14 N, and maximum linear and angular velocities of 2.0 m/s.

Table IV reports capture rates for each policy combination over 100 rollouts from random initial states. As illustrated, our evader successfully avoids capture from all policies except our own pursuer, while capturing all other policies. Moreover, the evader policy consistently outperforms all competing policies, and our pursuer performs comparably to MPC.

**TABLE IV:** Capture rates (%) for each evader–pursuer policy pairing in the 20D pursuit–evasion drone game, averaged over 100 rollouts from random initial conditions with a 3-second horizon.

Pursuer ( $\uparrow$ )		MADR-FOLLOW	MPC	Vanilla	MADR
Evader ( $\downarrow$ )	Pursuer ( $\uparrow$ )	MADR-FOLLOW	MPC	Vanilla	MADR
MADR	MADR	7	1	0	4
MPC	MPC	25	29	21	21
Vanilla	Vanilla	9	11	8	10



**Fig. 6:** A hardware demonstration of the drone vs. drone pursuit-evasion game with MADR. From left to right, four trajectory snapshots show the position of the evader (blue) and the pursuer (red) up to that time point with captures or out-of-bounds (OOB) marked by an ‘x’ (black). Both agents employ our proposed policy, with the pursuer using the follow-filtered augmentation.

## V. HARDWARE RESULTS

Having validated our approach in simulation (Section IV), we now evaluate its performance on the physical platform. This section presents experiments using hardware to assess real-world scalability.

### A. 6D Dubins Pursuit Evasion Game

To evaluate our approach in hardware, we implemented the 6D Dubins pursuit–evasion game on a pair of TurtleBots. Unlike the simulation studies (see Section IV-D), which focused on short horizons of 3 seconds, the hardware experiments probe longer-term behavior and robustness. In particular, we consider rollouts lasting up to 500 seconds, allowing us to assess whether value functions learned over short horizons remain effective for longer-horizon interactions, capturing stability effects and strategic patterns that are not visible in brief simulations.

Quantitative results on time to capture are reported in Table V, with representative trajectories highlighting qualitative differences in strategy shown in Figure 5. As shown, our policy significantly outperforms Vanilla DeepReach and achieves performance comparable to DP optimal, with similar times to capture. The trajectories further demonstrate long-term planning on the pursuer side, successfully cornering evaders. Despite the short-horizon learned value function, our method produces meaningful results even over 500-second experiments.

**TABLE V:** Average time-to-capture (seconds) for each evader-pursuer policy pairing, evaluated over six selected initial states with large relative separation.

Pursuer ( $\downarrow$ ) \ Evader ( $\uparrow$ )	DP	MADR	MADR-FOLLOW	Vanilla
DP	31.1	84.0	37.1	355.0
MADR	14.8	53.8	25.5	243.5
Vanilla	11.0	11.3	7.6	71.7

### B. 20D Drone Hardware Platform

To ground the 20D pursuit–evasion formulation in a realistic setting, we implement our framework on a pair of Crazyflie 2.1 quadrotors, each modeled with the 10-dimensional state space (See Section IV-E).

All neural network queries are performed offboard on a workstation. The Crazyflies receive only low-level combined thrust and desired attitude commands via a wireless link

at 50 Hz. For precise localization during experiments, the drones operate in a motion-capture arena providing real-time position and velocity measurements, which is fused with the internal IMU to provide full-state observability.

Each Crazyflie uses its onboard PID-based attitude controller to track the commanded thrust and torques. The actuation limits are consistent with the Crazyflie hardware: maximum thrust of approximately 15 N per rotor, roll and pitch angles constrained to  $\pm 0.15$  rad, and angular velocity limits of 2.0 rad/s.

We plot a representative roll-out, showcasing long horizon behavior in Figure 6. We compared the standard pursuit–evasion MADR policy to MADR-FOLLOW, finding that across near and distant initial positions MADR-FOLLOW achieved capture for 26.4% of the 3-minute trajectory on average versus 11.6% for MADR. Qualitatively, both policies induce diving attacks and dodge evasions (pictured in Figure 1) but the follow-filtered pursuit manages to recover better from distant points. See the supplemental for the trajectory videos.

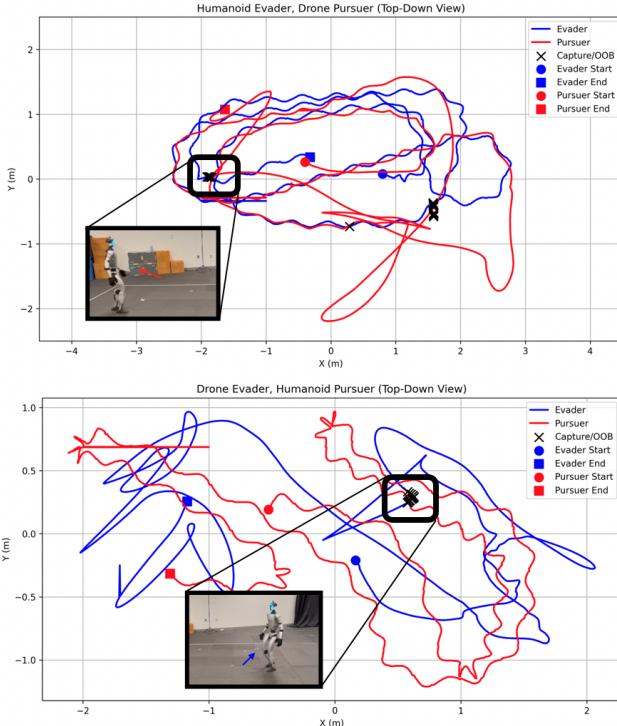
### C. Humanoid–Drone Pursuit–Evasion Experiments

To interrogate the robustness of our method, we conducted an experiment facing MADR against a Unitree G1 humanoid driven by a human-operator. In this game, we model the humanoid with quadrotor-like dynamics, although the opponent is weaker in reality due to its dynamic limitations, but one which we consider conservatively. A Crazyflie mounted on the humanoid provided state estimation via motion capture, while the aerial drone agent executed the MADR-FOLLOW strategy for the pursuer role and MADR for the evader role.

Two of six example trajectories are plotted in Figure 7, and the full videos may be found in the supplemental. In the pursuit case, the drone autonomously follows the rear of the humanoid and at another point rapidly lunges at its abdomen (pictured in the top of Figure 7). In the evader case, the drone frequents the lower heights, tending to “side-step” the approaching humanoid, but is occasionally knee’d (pictured in the bottom of Figure 7).

## VI. CONCLUSIONS

In this work, we proposed MADR, a robust learning framework that augments the DeepReach process with supervision from adversarial model predictive control (MPC). By leveraging worst-case control–disturbance interactions as



**Fig. 7:** Two hardware demonstrations of the teleoperated humanoid vs. MADR drone pursuit-evasion game. On top, the teleoperated humanoid evades, while on bottom, it pursues. In either case, the position of the evader (blue) and the pursuer (red) are plotted with captures marked by an ‘x’ (black).

expert guidance, MADR enables accurate approximation of Hamilton–Jacobi (HJ) reachability solutions and backward reachable tubes (BRTs) in the presence of both disturbances and adversarial agents. Our results highlight two key contributions. First, we show that MADR scales effectively across a range of dynamical systems and dimensions, and validate its practicality through both simulation and hardware experiments. Second, we demonstrate that MADR robustly handles diverse disturbance types, from structured environmental effects to two-agent zero-sum games. Together, these results underscore MADR’s potential for real-world deployment in safety-critical autonomous systems, where reliable decision-making under uncertainty is essential.

## REFERENCES

- [1] I. M. Mitchell, A. M. Bayen, and C. J. Tomlin, “A time-dependent hamilton–jacobi formulation of reachable sets for continuous dynamic games,” *IEEE Transactions on Automatic Control*, vol. 50, no. 7, pp. 947–957, 2005.
- [2] M. Chen, S. Herbert, and C. J. Tomlin, “Fast reachable set approximations via state decoupling disturbances,” in *Proc. IEEE Conf. on Decision and Control*, 2016.
- [3] S. Herbert, M. Chen, S. Han, W. Ma, and C. J. Tomlin, “Fastrack: A modular framework for real-time motion planning with guaranteed tracking,” in *Proc. IEEE Conf. on Decision and Control*, 2017.
- [4] S. Bansal, M. Chen, S. Herbert, and C. J. Tomlin, “Hamilton–Jacobi reachability: A brief overview and recent advances,” in *Proc. IEEE Conf. on Decision and Control*, 2017.
- [5] S. Bansal and C. J. Tomlin, “Deepreach: A deep learning approach to high-dimensional reachability,” in *Proc. IEEE Conf. on Robotics and Automation*, 2021.
- [6] W. Sharpless, Z. Feng, S. Bansal, and S. Herbert, “Linear supervision for nonlinear, high-dimensional neural control and differential games,” in *Learning for Dynamics & Control*, 2025.
- [7] J. Borquez, K. Nakamura, and S. Bansal, “Parameter-conditioned reachable sets for updating safety assurances online,” in *Proc. IEEE Conf. on Robotics and Automation*, 2023.
- [8] S. Tonkens, N. U. Shinde, A. Begzadie, M. C. Yip, J. Cortes, and S. L. Herbert, “From space to time: Enabling adaptive safety with learned value functions via disturbance recasting,” in *Conf. on Robot Learning*, 2025.
- [9] Z. Feng, L. Qiu, and S. Bansal, “Bridging model predictive control and deep learning for scalable reachability analysis,” in *Robotics: Science and Systems*, 2025.
- [10] G. P. Papavassiliopoulos, J. Medanic, and J. Cruz, “On the existence of nash strategies and solutions to coupled riccati equations in linear-quadratic games,” *Journal of Optimization Theory and Applications*, vol. 28, no. 1, pp. 49–76, 1979.
- [11] D. Fridovich-Keil, E. Ratner, L. Peters, A. D. Dragan, and C. J. Tomlin, “Efficient iterative linear-quadratic approximations for nonlinear multi-player general-sum differential games,” in *Proc. IEEE Conf. on Robotics and Automation*, 2020.
- [12] Y. Zhao and Q. Zhu, “Stackelberg meta-learning based control for guided cooperative lqg systems,” *IFAC-Papers Online*, vol. 56, no. 2, p. 10120–10125, 2023.
- [13] L. Falconi, A. Ferrante, and M. Zorzi, “Distributionally robust lqg control under distributed uncertainty,” *Automatica*, vol. 174, p. 112128, 2025.
- [14] F. Borrelli, A. Bemporad, and M. Morari, *Predictive Control for Linear and Hybrid Systems*. Cambridge Univ. Press, 2017.
- [15] G. Williams, A. Aldrich, and E. A. Theodorou, “Model predictive path integral control: From theory to parallelized implementation,” in *Proc. IEEE Conf. on Robotics and Automation*, 2016.
- [16] P. J. Campo and M. Morari, “Robust model predictive control,” in *American Control Conference*, 1987.
- [17] L. Pinto, J. Davidson, R. Sukthankar, and A. Gupta, “Robust adversarial reinforcement learning,” in *Int. Conf. on Machine Learning*, 2017.
- [18] Z. Li, C. Hu, S. E. Li, J. Cheng, and Y. Wang, “Robust safe reinforcement learning under adversarial disturbances,” in *Proc. IEEE Conf. on Decision and Control*, 2023.
- [19] R. Lowe, Y. Wu, A. Tamar, J. Harb, P. Abbeel, and I. Mordatch, “Multi-agent actor-critic for mixed cooperative-competitive environments,” in *Conf. on Neural Information Processing Systems*, 2017.
- [20] J. F. Fisac, N. F. Lugovoy, V. Rubies-Royo, S. Ghosh, and C. J. Tomlin, “Bridging hamilton–jacobi safety analysis and reinforcement learning,” in *Proc. IEEE Conf. on Robotics and Automation*, 2019.
- [21] O. So, C. Ge, and C. Fan, “Solving minimum-cost reach avoid using reinforcement learning,” in *Conf. on Neural Information Processing Systems*, 2024.
- [22] J. Hsu, J. Nguyen, and J. F. Fisac, “Isaacs: Iterative soft adversarial actor-critic for safety,” in *Learning for Dynamics & Control*, 2023.
- [23] D. P. Nguyen, K.-C. Hsu, W. Yu, J. Tan, and J. F. Fisac, “Gameplay filters: Robust zero-shot safety through adversarial imagination,” in *Conf. on Robot Learning*, 2024.
- [24] J. Wang, H. Hu, D. P. Nguyen, and J. F. Fisac, “Magics: Adversarial rl with minimax actors guided by implicit critic stackelberg for convergent neural synthesis of robot safety,” in *Workshop on Algorithmic Foundations of Robotics*, 2024.
- [25] L. Brunke, M. Greeff, A. W. Hall, Z. Yuan, S. Zhou, J. Panerati, and A. P. Schoellig, “Safe learning in robotics: From learning-based control to safe reinforcement learning,” *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 5, no. 1, pp. 411–444, 2022.
- [26] J. Sacks and B. Boots, “Learning to optimize in model predictive control,” in *Proc. IEEE Conf. on Robotics and Automation*, 2022.
- [27] J. Li, D. Lee, J. Lee, K. S. Dong, S. Sojoudi, and C. Tomlin, “Certifiable reachability learning using a new lipschitz continuous value function,” *IEEE Robotics and Automation Letters*, vol. 10, no. 4, pp. 3582–3589, 2025.