

# Safe Beyond the Horizon: Efficient Sampling-based MPC with Neural Control Barrier Functions

Ji Yin<sup>1\*</sup>, Oswin So<sup>2\*</sup>, Eric Yang Yu<sup>2</sup>, Chuchu Fan<sup>2</sup>, and Panagiotis Tsiotras<sup>1</sup>

**Abstract**—A common problem when using model predictive control (MPC) in practice is the satisfaction of safety specifications beyond the prediction horizon. While theoretical works have shown that safety can be guaranteed by enforcing a suitable terminal set constraint or a sufficiently long prediction horizon, these techniques are difficult to apply and thus are rarely used by practitioners, especially in the case of general nonlinear dynamics. To solve this problem, we impose a tradeoff between exact recursive feasibility, computational tractability, and applicability to “black-box” dynamics by learning an approximate discrete-time control barrier function and incorporating it into a variational inference MPC (VIMPC), a sampling-based MPC paradigm. To handle the resulting state constraints, we further propose a new sampling strategy that greatly reduces the variance of the estimated optimal control, improving the sample efficiency, and enabling real-time planning on a CPU. The resulting Neural Shield-VIMPC (NS-VIMPC) controller yields substantial safety improvements compared to existing sampling-based MPC controllers, even under badly designed cost functions. We validate our approach in both simulation and real-world hardware experiments.

## I. INTRODUCTION

Model Predictive Control (MPC) is a versatile control approach widely used in robotics applications such as autonomous driving [71], bio-inspired locomotion [9, 30], or manipulation of deformable objects [54], to name just a few. These methods address safety by incorporating state and control constraints into the finite-horizon optimization problem, ensuring that the system remains safe over the prediction horizon [9, 44, 54, 69]. However, safety of the system beyond the prediction horizon is often overlooked by practitioners, potentially leading to the violations of safety constraints at future timesteps.

This is a well-known problem in the field of MPC. The question of whether a sequence of safe control actions can always be found under an MPC controller has been studied extensively in the literature under the name of recursive feasibility [37, 47, 11, 27, 45]. A simple method of achieving recursive feasibility is by enforcing a control-invariant terminal set constraint at the end of the prediction horizon [37, 11]. However, it is difficult to find such a control-invariant set for general nonlinear systems. Methods such as bounding the system dynamics [11, 10, 21] often result in over-conservative sets, while methods that numerically solve the Hamilton-Jacobi

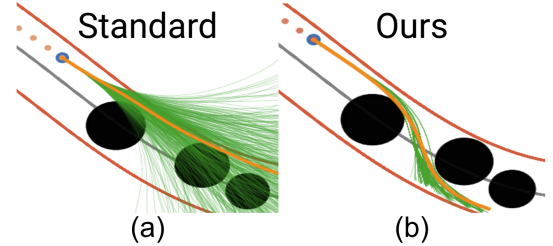


Fig. 1: Conventional sampling vs. the proposed Resampling-Based Rollout (RBR) approach. In this AutoRally example, the blue dot must sample future trajectories and compute an optimal control while avoiding the black obstacles. Our method rewires all sampled trajectories to be safe, resulting in a more accurate sampling distribution, whereas the standard approach samples from a Gaussian distribution and wastes computations on unsafe trajectories.

PDE [46] scale exponentially with the number of state variables, and hence this approach is computationally infeasible for systems with more than, say, 5 state variables [48].

Even without considering recursive feasibility, nonlinear *constrained* optimization is a challenging problem. Traditionally, constraints are handled using techniques such as interior-point, sequential quadratic programming and augmented Lagrangian type methods [51], which mostly rely on accurate linear or quadratic approximations of the cost and constraints. However, these gradient-based methods can get stuck in local minima or fail to converge when the problem is highly nonlinear. Moreover, first and second-order gradient information is needed to solve these optimization problems, which is often not available for “black-box” systems. Consequently, MPC controllers that rely on these underlying nonlinear constrained optimizers [51, 33, 25] inherit the same limitations.

One way to address these previous limitations is to use sampling-based optimization methods, which have become popular within both the robotics and model-based reinforcement learning communities. In particular, variational inference MPC (VIMPC) [52, 40, 69, 6, 55, 56] has emerged as a popular family of methods that pose the control problem as an inference problem instead of an optimization problem through the control as inference framework [4, 67, 60], and perform variational inference to approximate the resulting intractable optimal control distribution. Popular MPC controllers within this family include Model Predictive Path Integral (MPPI) [71] and the Cross-Entropy Method (CEM) [61, 69], which have been applied to robotics tasks such as autonomous driving [71],

<sup>1</sup>Ji Yin and Panagiotis Tsiotras are with D. Guggenheim School of Aerospace Engineering, Georgia Institute of Technology, Atlanta, GA. Email: {jyin81,tsiotras}@gatech.edu

<sup>2</sup>Oswin So, Eric Yang Yu and Chuchu Fan are with the Department of Aeronautics and Astronautics, Massachusetts Institute of Technology, Cambridge, MA. Email: {oswinso,eyyu,chuchu}@mit.edu

\* Equal contribution

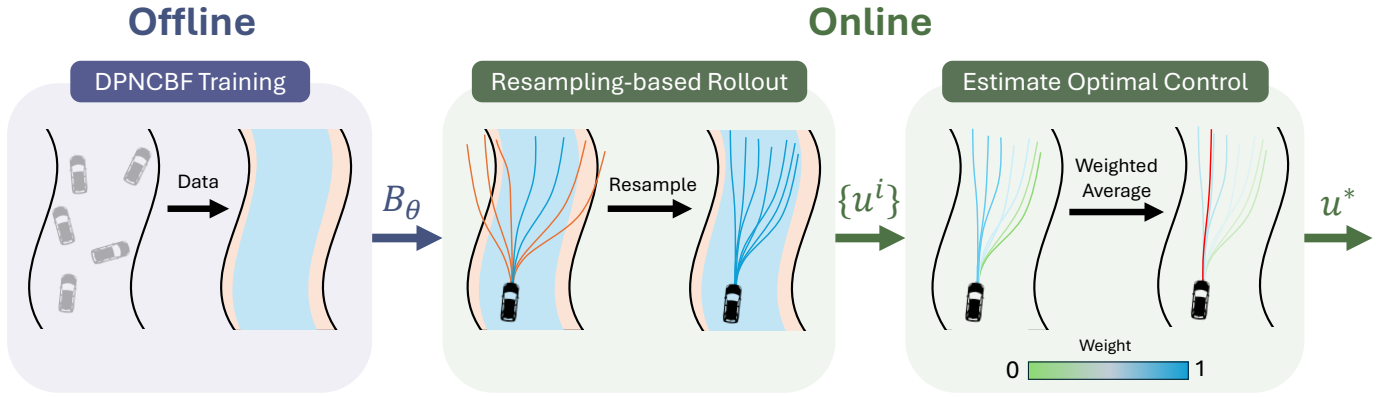


Fig. 2: **Overview of Neural Shield VIMPC (NS-VIMPC).** Offline, using DPNCBFs, we collect a dataset and train a NN approximation  $B_\theta$  of a DCBF.  $B_\theta$  is used to impose the DCBF descent condition (16) as a state constraint. Online, we modify the VIMPC architecture to use resampling-based rollouts to improve the sampling distribution of the Monte Carlo estimator in the presence of the DCBF state constraints. In addition, we integrate the DCBF safety condition (14b) into the optimization objective function as in (18).

bipedal locomotion [9], manipulation of deformable objects [54], and in-hand manipulation [49], among many others. Using this framework, constraints can be easily handled by appropriately manipulating the posterior control distribution without facing the same optimization challenges as traditional gradient-based methods. However, the problem of safety beyond the prediction horizon still remains a challenge and has not been addressed by existing works that employ VIMPC.

In this work, we present a novel sampling-based MPC approach that provides safety beyond the prediction horizon by using control barrier functions to enforce the control-invariant set constraint for VIMPC controllers. To tackle the challenge of finding control-invariant sets for general nonlinear systems, we extend our previous work on learning neural network approximations of control barrier functions using policy neural control barrier functions (PNCBF) [64] to the discrete-time case. Inspired by particle filtering and sequential Monte Carlo methods, we further propose a novel sampling strategy for handling state constraints in VIMPC that significantly improves the sampling efficiency and enables real-time planning on a CPU. Results from both simulation and hardware experiments suggest that the resulting Neural Shield-VIMPC (NS-VIMPC) controller outperforms existing MPC baselines in terms of safety, even under adversarially tuned cost functions.

**Contributions.** We summarize our contributions below.

- We extend policy neural CBF (PNCBF) to the discrete-time case and propose a novel approach to train a discrete-time PNCBF (DPNCBF) using policy evaluation.
- We propose Resampling-Based Rollout (RBR), a novel sampling strategy for handling state constraints in VIMPC inspired by particle filtering, which significantly improves the sampling efficiency by lowering the variance of the estimated optimal control.
- Simulation results on two benchmark tasks show the efficacy of NS-VIMPC compared to existing sampling-based MPC controllers in terms of safety and sample efficiency.

- Hardware experiments on AutoRally [24], a 1/5 scale autonomous driving platform, demonstrate the robustness of NS-VIMPC to unmodeled dynamical disturbances under adversarially tuned cost functions.

## II. RELATED WORK

**Sampling-based MPC.** Sampling-based MPC has become a popular alternative to traditional MPC methods that use gradient-based solvers, in part due to the advent of parallel computing and the recent advances in GPU hardware. By not using gradients, sampling-based MPC can be applied to any problem without having specific requirements on the problem structure. MPPI [71] is a popular sampling-based MPC approach that formulates a variational inference problem, then solves it in the case of the Gaussian distribution, having strong connections to stochastic optimal control [72] and maximum entropy control [63]. Separately, the Cross-Entropy Method (CEM) [61] has become popular in the reinforcement learning community [69], in part due to its simplicity. Both approaches were recently shown to be part of the VIMPC family of MPC algorithms [52]. Consequently, some works have looked at expanding the VIMPC family to include different choices of divergences [69] and sampling distributions beyond Gaussians [40, 52].

**Safety in Sampling-based MPC.** Research efforts such as [80] and [3] assess the risk associated with uncertain areas in the state space during the exploration phase, and enhance MPPI's safety by incorporating a risk penalty into its cost function. While these methods empirically enhance safety, they lack formal assurances. Another class of MPPI alternatives, for instance, [72, 57], leverage an auxiliary tracking controller to follow the MPPI output trajectories, improving robustness against unforeseeable disruptions, but these improvements are limited when the simulation-to-reality gap is significant. More recently, the use of Control Barrier Functions (CBF) [66, 79] has offered formal safety guarantees for MPPI controllers.

Nonetheless, [66] is restricted to linear systems with saturating CBFs, and also neglects control input limits. This can lead to Gaussian control distributions constrained by a CBF-based chance constraint, increasing the likelihood of local optima due to limited exploration. In contrast, [79] and [81] incorporate a CBF into the cost function and resolve a local optimization problem to reinforce CBF safety conditions for more complex nonlinear systems. However, the approach in [79] and [81] uses a distance-based CBF while ignoring the control limits. Our proposed Neural Shield MPPI (NS-MPPI) controller, as a specific form of the proposed NS-VIMPC framework, effectively resolves the critical issues identified in [66] and [79].

**Different Proposal Distributions for VIMPC.** There are many works that investigate the sampling distribution of VIMPC. The MPPI variant in [78] uses covariance steering to assign a terminal covariance to the sampling distribution, but this relies on expert knowledge on how the covariance should be designed. Reference [55] learns a normalizing flow to approximate the optimal distribution, but it does not address the problem of recursive feasibility on its own. Another direction looks at changing the effective sampling distribution by modifying the underlying dynamics system to be more amenable to calculations. The works [72, 57] leverage an auxiliary tracking controller, thus changing the sampling distribution to be the output of the stable tracking controller. However, this requires the construction of such an auxiliary tracking controller, which can be difficult to perform for arbitrary nonlinear discrete-time systems. In contrast, our proposed resampling-based rollouts (RBR) can be viewed as a way of easily improving the proposal distribution *without* the need for any specialized problem structure.

**Duality between Control and Inference.** The proposed resampling strategy in our work is similar in spirit to [75, 58] in that factor graphs, a method used originally for estimation, is adapted for control purposes. In [84] a linear optimal control was used to improve the performance of particle filters. However, to the best of our knowledge, our approach is the first to adopt the resampling mechanism from particle filters to improve the performance of sampling-based MPC controllers.

**Control Barrier Functions.** Designing CBFs with control-invariant safe sets is not a trivial task. As a result, previous works only seek saturating CBFs that do not consider input constraints [73, 43, 76]. Some works use hand-tuned CBFs to prevent saturation [70, 12], which can lead to overly conservative safe sets. Other works develop CBFs with input constraints for specific types of systems [13, 14]. The emerging neural CBFs [85, 83, 59, 86] allow for more general dynamics utilizing machine learning techniques; however, they can still saturate the control limits. **Our novel Discrete-time Policy Neural Control Barrier Function (DPNCBF) approach trains a neural CBF conditioned on the full system state while respecting control limits, enabling the discovery of a large control-invariant set. By formulating the problem as a discrete-time avoid task, DPNCBF customizes safety based on any given**

**control policy, and can integrate seamlessly with discrete-time controllers, enhancing its practical utility. Table I compares the proposed DPNCBF with other reachability-based methods.**

TABLE I: Relationship between DPNCBF and other reachability methods.

Method	Avoid	Discrete	(Arbitrary) Policy-Conditioned
Bansal and Tomlin [5]	✓	✗	✗
Fisac et al. [22]	✓	✗	✗
Hsu et al. [32]	✗	✓	✗
So et al. [64]	✓	✗	✓
He et al. [31]	✗	✓	✓
<b>DPNCBF (ours)</b>	✓	✓	✓

The proposed novel NS-VIMPC control framework utilizes DPNCBF to improve the efficiency of sampling by redirecting control sequences to previously identified safe states, thereby focusing sampling efforts in secure regions. This efficient sampling method will be discussed in detail later in Section IV. As a result, NS-VIMPC tackles the challenge of exploration versus exploitation by adaptively distributing computational resources, governed by an innovative neural CBF.

### III. PROBLEM FORMULATION

We consider the discrete-time, nonlinear dynamics

$$x_{k+1} = f(x_k, u_k), \quad (1)$$

with state  $x \in \mathcal{X} \subseteq \mathbb{R}^{n_x}$  and control  $u \in \mathcal{U} \subseteq \mathbb{R}^{n_u}$ . Let  $\mathcal{U}^K$  denote the set of control trajectories of length  $K$ . Following the MPC setup, we assume that a cost function  $J : \mathcal{U}^K \rightarrow \mathbb{R}$  encoding the desired behavior of the system is given. Moreover, we consider state constraints defined by an *avoid set*  $\mathcal{A} \subset \mathcal{X}$  described as the superlevel set of some specification function  $h$ , i.e.,

$$\mathcal{A} := \{x \in \mathcal{X} \mid h(x) > 0\}. \quad (2)$$

The goal is then to find a sequence of controls  $\mathbf{u} = \{u_0, u_1, \dots, u_{K-1}\} \in \mathcal{U}^K$  that minimizes the cost function  $J$  while satisfying the system dynamics (1) and safety constraints  $x_k \notin \mathcal{A}$  for all  $k \geq 0$ .

#### A. Variational Inference MPC For Sampling-based Optimization

To solve the above optimization problem, we deviate from traditional MPC solvers and use a variational inference MPC formulation to solve the problem via sampling-based optimization. To this end, we make use of the control-as-inference framework [41] to model the problem. Specifically, let  $o$  be a binary variable that indicates “optimality” such that, for all controls  $\mathbf{u} \in \mathcal{U}^K$ ,

$$p(o = 1 \mid \mathbf{u}) \propto \exp(-J(\mathbf{u})), \quad (3)$$

We assume a prior  $p_0(\mathbf{u})$  on the control trajectory. The “optimal” distribution can then be obtained via the posterior distribution

$$p(\mathbf{u} \mid o = 1) = \frac{p(o = 1 \mid \mathbf{u})p_0(\mathbf{u})}{p(o = 1)} = Z^{-1} \exp(-J(\mathbf{u}))p_0(\mathbf{u}), \quad (4)$$

where  $Z := \int \exp(-J(\mathbf{u}))p_0(\mathbf{u})d\mathbf{u}$  denotes the unknown normalization constant. Since sampling from  $p(\mathbf{u} | o = 1)$  is intractable, we use variational inference to approximate the posterior  $p(\mathbf{u} | o = 1)$  with a tractable distribution  $q_{\mathbf{v}}(\mathbf{u})$  parametrized by some vector  $\mathbf{v}$  by minimizing the forward KL divergence, i.e.,

$$\min_{\mathbf{v}} \mathbb{D}_{\text{KL}}(p(\mathbf{u} | o = 1) \parallel q_{\mathbf{v}}(\mathbf{u})). \quad (5)$$

In the special case of  $q_{\mathbf{v}}$  being a Gaussian distribution with mean  $\mathbf{v}$  and a fixed control input covariance  $\Sigma$ , intrinsic to the robotic system of interest [71], we can solve (5) in closed-form to obtain the optimal  $\mathbf{v}^*$  as (see Section B1 for details)

$$\mathbf{v}^* = \mathbb{E}_{p(\mathbf{u}|o=1)}[\mathbf{u}]. \quad (6)$$

While this expectation cannot be readily computed because  $p(\mathbf{u} | o = 1)$  is intractable to sample from, we can use importance sampling to change the sampling distribution to some other distribution  $r(\mathbf{u})$  that is easier to sample from, leading to

$$\mathbf{v}^* = \mathbb{E}_{r(\mathbf{u})} \left[ \frac{p(\mathbf{u} | o = 1)}{r(\mathbf{u})} \mathbf{u} \right] \quad (7)$$

$$= \mathbb{E}_{r(\mathbf{u})} \left[ \underbrace{\frac{Z^{-1} \exp(-J(\mathbf{u}))p_0(\mathbf{u})}{r(\mathbf{u})}}_{:=\omega(\mathbf{u})} \mathbf{u} \right]. \quad (8)$$

Using samples  $\mathbf{u}^1, \dots, \mathbf{u}^N$  drawn from  $r$ , we compute a Monte Carlo estimate  $\hat{\mathbf{v}}$  of the optimal control sequence  $\mathbf{v}^*$  (see Section B2 for details) as follows,

$$\hat{\mathbf{v}} = \sum_{i=1}^N \tilde{\omega}^i \mathbf{u}^i, \quad (9)$$

$$\tilde{\omega}^i := \frac{\omega(\mathbf{u}^i)}{\sum_{j=1}^N \omega(\mathbf{u}^j)}. \quad (10)$$

Note that the  $Z$  in  $\omega(\mathbf{u})$  is canceled out in the computation of  $\tilde{\omega}^i$  in (10) and hence can be ignored.

**Remark 1** (Self-normalized importance sampling). Note that the weights  $\tilde{\omega}^i$  in (9) are **not** from the regular importance sampling estimates in (8) due to the normalization by the sum of the weights in (10). Instead, (9) is a self-normalized importance sampling estimator (SNIS), which uses an *estimate* of the weights  $\omega$  but results in a biased, though asymptotically unbiased, estimator. This fact is not present in many existing works on both VIMPC (e.g., [52, 69]) and MPPI (e.g., [71]). See Section B2 for more details.

**Remark 2** (Connections to MPPI and CEM). In the case where we choose  $p_0(\mathbf{u}) = q_0(\mathbf{u})$  and  $r(\mathbf{u}) = q_{\bar{\mathbf{v}}}(\mathbf{u})$  for some previous estimate of the optimal control sequence  $\bar{\mathbf{v}}$ , the above variational inference MPC framework reduces to MPPI (see Section B3 for details). Moreover, Cross-Entropy Method (CEM) also falls in the VIMPC framework [52, 69].

## B. Constraint Handling In Variational Inference MPC

One advantage of sampling-based MPC is that it is simple to incorporate hard constraints. One can include an indicator function in the cost function that heavily penalizes constraint violations (e.g., see [71, 72, 6, 8]). Specifically, for some large constant  $C > 0$ , we can modify the cost function as

$$J_{\text{new}}(\mathbf{u}) = J(\mathbf{u}) + C \sum_{k=0}^K \mathbb{1}_{x_k \in \mathcal{A}}. \quad (11)$$

From the inference perspective (3), we can interpret  $J_{\text{new}}$  (11) as saying that  $o = o_{\text{cost}} \wedge o_{\text{constraint}}$ , where,

$$p(o_{\text{cost}} = 1 | \mathbf{u}) \propto \exp(-J(\mathbf{u})), \quad (12)$$

$$p(o_{\text{constraint}} = 1 | \mathbf{u}) \propto \exp\left(-C \sum_{k=0}^K \mathbb{1}_{x_k \in \mathcal{A}}\right). \quad (13)$$

However, a problem with (13) is that it looks at safety only within the prediction horizon and does not consider the probability of staying safe beyond the prediction horizon. To tackle this problem, we will use a discrete-time control barrier function (DCBF), which we introduce in the next section, to enforce that the states remain within a control-invariant set.

## C. Discrete-time Control Barrier Functions (DCBF)

A discrete-time control barrier function (DCBF) [79] associated to the avoid set  $\mathcal{A}$  is a function  $B : \mathcal{X} \rightarrow \mathbb{R}$  such that<sup>1</sup>

$$B(x) > 0, \quad \forall x \in \mathcal{A}, \quad (14a)$$

$$B(x) \leq 0 \implies \inf_{u \in \mathcal{U}} B(f(x, u)) - B(x) \leq -\alpha(B(x)), \quad (14b)$$

where  $\alpha$  is an extended class- $\kappa$  function [77]. As in [79], we restrict our attention to the class of linear extended class- $\kappa$  functions, i.e.,

$$\alpha(B(x)) = a \cdot B(x), \quad a \in (0, 1). \quad (15)$$

The following theorem from [79] proves that a controller satisfying the condition (14b) renders the sublevel set  $\mathcal{S} = \{x | B(x) \leq 0\}$  forward-invariant.

**Theorem 1** ([79, Property 3.1]). *Any control policy  $\pi : \mathcal{X} \rightarrow \mathcal{U}$  satisfying the condition*

$$B(f(x, \pi(x))) - B(x) \leq -\alpha(B(x)), \quad (16)$$

*renders the sublevel set  $\mathcal{S} = \{x | B(x) \leq 0\}$  forward-invariant.*

Hence, one way to guarantee recursive feasibility, and thus safety beyond the prediction horizon, is to enforce the condition (16) at every time step of the optimization problem. Another point to note is that  $\mathcal{S}$  is a control-invariant set, and thus the constraint  $B(x_K) \leq 0$  can be imposed as a terminal state constraint to guarantee recursive feasibility for MPC as is done classically [37, 11]. Thus, one can try to enforce these constraints by incorporating them into the cost function as

<sup>1</sup>Note that we use the opposite sign convention as compared to [79].



in (11) [79], i.e., modify the cost according to one of the following options,

$$J_{\text{new}}(\mathbf{u}) = J(\mathbf{u}) + C \sum_{k=0}^K \mathbb{1}_{B(f(x,u)) - B(x) > -\alpha(B(x))}, \quad (17)$$

$$J_{\text{new}}(\mathbf{u}) = J(\mathbf{u}) + C \sum_{k=0}^K \left[ B(f(x,u)) - B(x) + \alpha(B(x)) \right]_+. \quad (18)$$

However, this approach has two problems. First, for sufficiently large  $C$ , samples that violate the DCBF constraint will have a normalized weight of near zero, rendering these samples useless. Conservative DCBFs may cause the majority of the samples to violate the constraint (16) and hence have zero weight, resulting in poor estimates of the optimal control and wasted computation. Second, while constructing a function  $B$  that satisfies (14a) is relatively simple, it is much harder to construct a function  $B$  that also satisfies (14b), contrary to the case with (continuous-time) control barrier functions [64, 83] (see also Remark 3 below). Consequently, many works that integrate control barrier functions into MPC often only propose functions for which (14a) holds and not (14b) [79], rendering the safety guarantees of Theorem 1 invalid.

In the next section, we address these two problems via a novel resampling method that reuses computations from zero-weight samples and learns a DCBF that tries to respect both (14a) and (14b) using policy value functions.

**Remark 3** (Differences between discrete-time and continuous-time control barrier functions under unbounded controls). Note that in the continuous-time case where having an unbounded control space  $\mathcal{U}$ , control-affine dynamics, and a non-zero  $\frac{\partial \dot{B}}{\partial u}$  are sufficient to guarantee that  $B$  is a valid control barrier function. This is because (14b) generally holds for a function  $B$  that satisfies (14a) in the continuous-time case under these assumptions. However, the same does not apply to discrete-time under general nonlinear dynamics  $f$  since (16) is nonlinear in  $u$ , let alone the fact that robotic systems in real life are unable to exert infinite forces and hence generally do not have unbounded controls.

**Remark 4** (Constraint satisfaction in the variational inference framework). One potential issue with incorporating DCBF constraints into the cost function  $J$  is that, while  $p(\mathbf{u} \mid o = 1)$  may have zero density on the set of controls that violate the DCBF constraint, this does not necessarily hold for the approximating distribution  $q$  since we are minimizing the *forward* KL divergence [34]. In particular, the mean  $\mathbf{v}$  of  $q_{\mathbf{v}}$ , which is the control to be used, may not satisfy the DCBF constraint. One way to guarantee that  $\mathbf{v}$  does satisfy the DCBF constraint is to assume that the set of controls that satisfy the DCBF constraint is itself convex (see Appendix F). However, this is an unrealistic assumption that is often violated by state constraints such as obstacle avoidance.

Despite these shortcomings, we observed in our experimental results that this method of enforcing DCBF constraints on  $p(\mathbf{u} \mid o = 1)$  indeed drastically improved safety. We leave further theoretical exploration of this issue as future work.

#### IV. NEURAL SHIELD VIMPC

In this section, we propose Neural Shield VIMPC (NS-VIMPC), a sampling-based MPC paradigm that efficiently samples trajectories using a DCBF modeled using a neural network. We illustrate the proposed NS-VIMPC algorithm in Fig. 2.

##### A. Approximating DCBF Using Neural Policy Value Functions

Let  $x_k^\pi$  denote the state at time  $k$  following the control policy  $\pi : \mathcal{X} \rightarrow \mathcal{U}$ . Define the policy value function  $V^{h,\pi}$  as,

$$V^{h,\pi}(x_0) := \max_{k \geq 0} h(x_k^\pi). \quad (19)$$

We have the following theorem.

**Theorem 2.**  $V^{h,\pi}$  satisfies (14a) and (14b) and is a DCBF.

*Proof:* From the definition of  $V^{h,\pi}$  (19), we have that

$$V^{h,\pi}(x_k) \geq h(x_k), \quad (20)$$

$$V^{h,\pi}(x_k) \geq V^{h,\pi}(f(x_k, \pi(x_k))). \quad (21)$$

Using the definition of the avoid set  $\mathcal{A}$  (2) and (20), it follows that  $V^{h,\pi}(x) > 0$  for all  $x \in \mathcal{A}$ , satisfying the first condition (14a) of a DCBF. When  $V^{h,\pi}(x_k) \leq 0$ , we have  $-\alpha(V^{h,\pi}(x_k)) \geq 0$ , and (21) implies that,

$$V^{h,\pi}(f(x_k, \pi(x_k))) - V^{h,\pi}(x_k) \leq 0 \leq -\alpha(V^{h,\pi}(x_k)). \quad (22)$$

Since  $\pi(x_k) \in \mathcal{U}$ , this implies the second condition (14b). Thus, the policy value function  $V^{h,\pi}$  is a DCBF. ■

Although we have constructed a DCBF from (19), the challenge is that the policy value function  $V^{h,\pi}$  cannot be easily evaluated at arbitrary states since the maximization in (19) is taken over an infinite horizon. To fix this, we train a neural network approximation  $V_\theta^{h,\pi}$  of  $V^{h,\pi}$ , extending the approach of [64] to the discrete-time case. To begin, we first rewrite (19) in a dynamic programming form,

$$V^{h,\pi}(x_0) = \max \left\{ \max_{0 \leq k \leq T} h(x_k^\pi), V^{h,\pi}(x_T^\pi) \right\}. \quad (23)$$

We can then train a neural network  $V_\theta^{h,\pi}$  to approximate the value function  $V^{h,\pi}$  by minimizing the loss

$$L(\theta) = \left\| V_\theta^{h,\pi}(x_0) - \max \left\{ \max_{0 \leq k \leq T} h(x_k^\pi), V_\theta^{h,\pi}(x_T^\pi) \right\} \right\|^2$$

over all states  $x_0$ . One problem, however, is that the minimizer of (IV-A) is not unique. For example, if  $h(x) \leq \bar{h}$  for all  $x$ , then  $V_\theta^{h,\pi} = \bar{h}$  is a minimizer of (IV-A) but does not necessarily satisfy (19). To fix this, we follow the approach of [22, 62] and, inspired by reinforcement learning [65], introduce a discount factor  $\gamma \in (0, 1)$  to define the *discounted* value function

$$V^{h,\pi,\gamma}(x_k) = \max \left\{ h(x_k^\pi), (1-\gamma)h(x_k) + \gamma V^{h,\pi,\gamma}(x_{k+1}^\pi) \right\}, \quad (24)$$

and the corresponding loss  $L$ ,

$$L(\theta) = \left\| V_{\theta}^{h,\pi,\gamma}(x_k) - \hat{V}_{\theta}^{h,\pi,\gamma}(x_k) \right\|^2, \quad (25)$$

$$\hat{V}_{\theta}^{h,\pi,\gamma}(x_k) = \max \left\{ h(x_k^{\pi}), (1 - \gamma)h(x_k^{\pi}) + \gamma V_{\theta}^{h,\pi,\gamma}(x_{k+1}^{\pi}) \right\}. \quad (26)$$

Finally, instead of using the learned  $V_{\theta}^{h,\pi,\gamma}$  directly in (16), we first take the maximum with  $h$  and use  $\tilde{V}_{\theta}^{h,\pi,\gamma}$  defined as

$$\tilde{V}_{\theta}^{h,\pi,\gamma}(x) := \max \{ h(x), V_{\theta}^{h,\pi,\gamma}(x) \}. \quad (27)$$

This guarantees that  $\tilde{V}_{\theta}^{h,\pi,\gamma}(x) \geq h(x)$  and hence the zero sublevel set of  $\tilde{V}_{\theta}^{h,\pi,\gamma}$  will be a subset of  $h$ . Hence, imposing the state constraint  $\tilde{V}_{\theta}^{h,\pi,\gamma}(x) \leq 0$  will, at the very least, prevent violations of the original state constraints during the prediction horizon, and potentially also induce a state constraint that is closer to the true control-invariant set than the original sublevel set of  $h$ .

**Remark 5** (Neural Network Verification of DCBFs). We emphasize that our goal here is to obtain a good approximation of a DCBF  $B_{\theta}$  using a neural policy value function  $V_{\theta}^{h,\pi,\gamma}$ , and not necessarily to obtain a true DCBF. Verifying whether the learned  $V_{\theta}^{h,\pi,\gamma}$  is a true DCBF requires neural network verification which can be intractable or inconclusive (see Appendix 1 in [36] on the NP-completeness of the NN-verification problem). This is especially true in the discrete-time case where the condition (16) may not be affine in the control  $u$ .

Nevertheless, as we show later, empirical results show that using an approximation of a DCBF is sufficient for enabling the use of much shorter prediction horizons without sacrificing safety.

### B. Efficient Sampling Using Resampling-based Rollouts

We tackle the problem of wasted samples with zero weights by drawing inspiration from the sequential Monte Carlo [18] and particle filter [28] literature, and by performing a per-timestep resampling during the rollout, which we call Resampling-Based Rollouts (RBR). Specifically, the control-as-inference problem formulation (3) gives us a temporal decomposition of  $p(o_{\text{constraint}} = 1 \mid \mathbf{u})$  (13) in the case of state constraints when  $C \rightarrow \infty$ , as follows

$$p(o_{\text{constraint}} = 1 \mid \mathbf{u}) \propto \prod_{k=0}^K \mathbb{1}_{x_k \notin \mathcal{A}}. \quad (28)$$

Hence, we can write the posterior  $p(\mathbf{u} \mid o = 1)$  as

$$p(\mathbf{u} \mid o = 1) \propto p(o_{\text{cost}} = 1 \mid \mathbf{u}) p(\mathbf{u}) \left( \prod_{k=0}^K \mathbb{1}_{x_k \notin \mathcal{A}} \right). \quad (29)$$

By treating the term on the right as a “measurement model,” it allows us to use particle filtering [28] to solve the control-as-inference problem. The update of the particle filter weights  $\hat{w}_k^i$  at time step  $k$  for particle  $i$  is written as

$$\hat{w}_k^i = \hat{w}_{k-1}^i \mathbb{1}_{x_k \notin \mathcal{A}}. \quad (30)$$

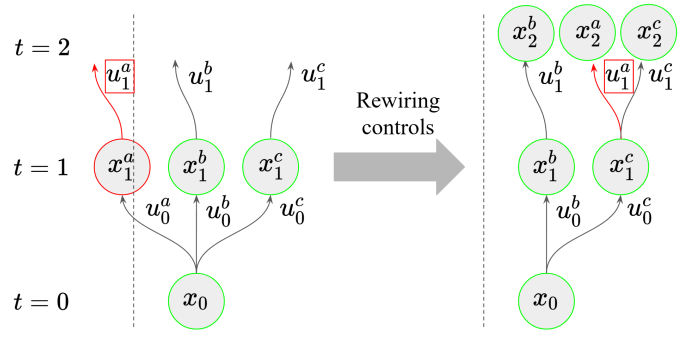


Fig. 3: **Resampling-based Rollouts (RBR)**. Inspired by particle filtering, at each step of the trajectory rollout, we uniformly resample any samples that violate the state constraints among the set of safe samples. In this example, since the state  $x_1^a$  violates constraints, it has a weight  $\hat{w}_1^a = 0$ . The control  $u_1^a$  is rewired to state  $x_1^c$  (since  $\hat{w}_1^b = \hat{w}_1^c = 1$ ,  $x_1^b$  and  $x_1^c$  has equal chance of getting  $u_1^a$ ), resulting in a new sample state  $x_2^a$ .

Then, we *resample* the particles with probability proportional to their weights  $\hat{w}_k^i$  to obtain a new set of particles  $\mathbf{u}^i$ . Due to the indicator function in the weight update (30), the weights are either 0 or 1. Assuming there exists a particle that satisfies the state constraint, applying systematic resampling [38] results in “rewiring” particles that violate the constraint to particles that still maintain safety, reusing the computation from zero-weight samples (see Fig. 3). **Note that we check state constraint  $N(K - 1)$  times for potential resampling during a rollout (once for each particle, all time steps except the last). When all particles violate the constraint, we do not resample.** In this case, since we want to minimize constraint violations, we still use the cost term (18) such that trajectories with higher constraint violations have higher costs.

While we can prove that this resampling is unbiased from the particle filter perspective, we also provide a more direct proof of this fact without the analogy to particle filters.

**Theorem 3.** For a probability density function  $f$  and set  $S$ , let  $a$  be sampled from the conditional density  $f(x \mid x \in S)$ , and let  $b$  be sampled from the unconditional density  $f$ , such that  $a$  and  $b$  are independent. Define the “rewired” random variable  $\tilde{b}$  to be equal to  $b$  if  $b \in S$  and  $a$  otherwise, i.e.,

$$\tilde{b} = \mathbb{1}_{b \in S} b + \mathbb{1}_{b \notin S} a \quad (31)$$

Then,  $\tilde{b}$  is also sampled from the conditional density  $f(x \mid x \in S)$ , such that  $\tilde{b}$  and  $a$  have the same distribution,  $\tilde{b} \stackrel{d}{=} a$ .

The proof is given in Section D2. Consequently, we can use  $b$  in a Monte Carlo estimator and still obtain unbiased estimates, as shown in the following Corollary (proof is given in Section D3).

**Corollary 1.** For any function  $w$ , the Monte Carlo estimate of  $\mathbb{E}[w(x)]$  under the conditional density  $f(x \mid x \in S)$  using

random variables  $\mathbf{a}$  and  $\tilde{\mathbf{b}}$  is unbiased, i.e.,

$$\mathbb{E}\left[\frac{1}{2}w(\mathbf{a}) + \frac{1}{2}w(\tilde{\mathbf{b}})\right] = \mathbb{E}[w(\mathbf{x})]. \quad (32)$$

Intuitively, resampling improves the efficiency of the Monte Carlo estimator, since most of the samples do not violate the state constraints and hence contribute to the weighted sum with non-zero weight. We can also theoretically prove that resampling improves the variance of the resulting Monte Carlo estimator. In the following theorem, we show how resampling reduces the *exponential* growth of the variance on the prediction horizon to a constant factor in the limit as the number of samples  $N$  goes to infinity.

**Theorem 4.** Let the horizon  $K > 0$ , consider  $\mathcal{U} = [-1, 1]$ , and define the avoid set and the dynamics such that  $\mathbf{u} \in [0, 1]^K$  is safe, and unsafe otherwise. Let the prior distribution  $p(\mathbf{u})$  be uniform on  $[-1, 1]^K$  and let  $p(o = 1 | \mathbf{u})$  be the indicator function for  $\mathbf{u} \in [0, 1]$  such that the posterior distribution  $p(\mathbf{u} | o = 1)$  is uniform on  $[0, 1]^K$ , and the proposal distribution  $r(\mathbf{u}) = 1/2^K$  is uniform on  $\mathcal{U}$ . Then, the variance of the Monte Carlo estimator of the optimal control law

$$\hat{\mathbf{v}} = \frac{1}{N} \sum_{i=1}^N \frac{p(\mathbf{u}^i | o = 1)}{r(\mathbf{u}^i)} \mathbf{u}^i, \quad (33)$$

grows exponentially in  $K$ , i.e.,

$$\text{Var}[\hat{v}_k] = \frac{1}{N} \left( \frac{1}{3} 2^K - \frac{1}{4} \right), \quad k = 1, \dots, K. \quad (34)$$

Using resampling, the variance is upper-bounded by,

$$\text{Var}[\hat{v}_{k, \text{resample}}] = O \left( \left( \frac{1}{1 - 2^{-N}} \right)^K + (1 - 2^{-N})^K \right). \quad (35)$$

The proof of Theorem 4 is given in Section D4. As shown in Theorem 4, although the variance is still exponential in  $K$  using resampling, the base of the exponential decreases to 1 exponentially in the number of samples  $N$ . In other words, in the limit as  $N \rightarrow \infty$ , the variance of the estimator using the proposed RBR is bounded by a constant factor. This novel approach reduces the variance of the Monte Carlo estimator of the optimal control law in a way that mirrors the relationship between Sequential Importance Sampling (i.e., particle filters without resampling), where the variance increases exponentially with the horizon length, and Sequential Monte Carlo (i.e., particle filters with resampling), where the (asymptotic) variance only increases linearly [19].

Another method to theoretically quantify the improvement in the variance of the estimator is via the *effective sample size* (ESS) [19, 20]. ESS is defined as the ratio between the variance of the estimator with  $N$  samples from the target and the variance of the SNIS estimator [20]. It can be interpreted as the number of samples simulated from the target pdf that would provide an estimator with variance equal to the performance

of the  $N$ -sample SNIS estimator. However, since the ESS is computationally intractable, the *approximation* (made formal in [20])

$$\widehat{ESS} := \frac{1}{\sum_{n=1}^N w_n^2}, \quad (36)$$

is more often used in practice. The following theorem shows that performing RBR results in either the same or higher  $\widehat{ESS}$ , given certain assumptions (proof is given in Section D5).

**Theorem 5.** Let  $\mathbf{w} = [w_1, \dots, w_m, 0, \dots, 0]$  denote the unnormalized weight vector without resampling, where the last  $N - m$  entries are zero due to violating the safety constraints. Let  $\mathbf{w}' = \mathbf{w} + \mathbf{c} = [w_1, \dots, w_m, c_{m+1}, \dots, c_N]$  denote the unnormalized weight vector resulting from safe resampling, where  $\mathbf{c} = [0, \dots, 0, c_{m+1}, \dots, c_N]$ . Let  $\tilde{\mathbf{w}} = \mathbf{w} / \|\mathbf{w}\|_1$  and  $\tilde{\mathbf{w}}' = \mathbf{w}' / \|\mathbf{w}'\|_1$  denote the normalized weights. Suppose that the weights of the resampled trajectories  $\mathbf{c}$  are not “drastically larger” than the weights of the original trajectories  $\mathbf{w}$ , i.e.,

$$\|\mathbf{c}\|_1 \leq 2 \frac{N}{N-1} \frac{\|\mathbf{w}\|_2^2}{\|\mathbf{w}\|_1}. \quad (37)$$

Then, the  $\widehat{ESS}$  using RBR is no smaller than the  $\widehat{ESS}$  without resampling, and is strictly greater if the inequality in (37) is strict. In other words,

$$\frac{1}{\|\tilde{\mathbf{w}}\|_2^2} \leq \frac{1}{\|\tilde{\mathbf{w}}'\|_2^2}. \quad (38)$$

### C. Summary of NS-VIMPC

We now summarize the NS-VIMPC algorithm, shown in Fig. 2.

Offline, we first approximate a DCBF by using the DPNCBF algorithm to learn the policy value function for a user-specified policy (Section IV-A). In our experiments, we chose this to be Shield MPPI (S-MPPI) [79]. Online, we use the learned DPNCBF to enforce the DPNCBF constraint (16) to try to enforce safety beyond the prediction horizon. During sampling, we use RBR to efficiently sample control sequences  $\{\tilde{\mathbf{u}}^i\}$  from the raw samples  $\{\mathbf{u}^i\}$  to satisfy the DPNCBF constraint, thus improving the sample efficiency of the Monte Carlo estimate of the optimal control. The sampled control sequences  $\{\tilde{\mathbf{u}}^i\}$  are then used to compute the estimate of the optimal control  $\hat{\mathbf{v}}$  using (9). As in MPC fashion, we only execute the first control  $\hat{\mathbf{v}}_0$ , and  $\hat{\mathbf{v}}$  is then used as the parameter vector of the sampling distribution  $q_{\mathbf{v}}$  for the next iteration.

## V. SIMULATIONS

We first performed simulation experiments to better understand the performance of the proposed Neural Shield VIMPC (NS-VIMPC) controller. Although many sampling-based MPC controllers fall under the VIMPC family with different choices of the prior  $p_0$  and  $r$ , we choose to instantiate the MPPI algorithm (see Section B1 for details), and call the resulting controller Neural Shield-MPPI (NS-MPPI).

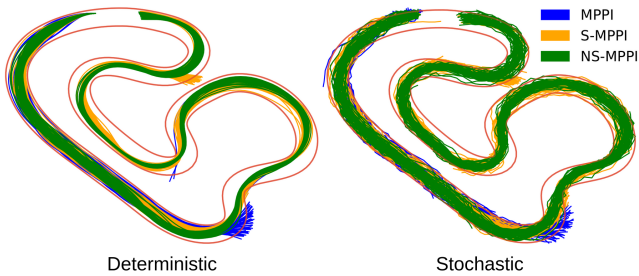


Fig. 4: **AutoRally Trajectories.** We visualize the trajectories of the three MPPI baselines under a challenging target velocity of  $15 \text{ ms}^{-1}$ . Both MPPI and S-MPPI veer off course and crash while NS-MPPI stays within the track even under Gaussian disturbances.

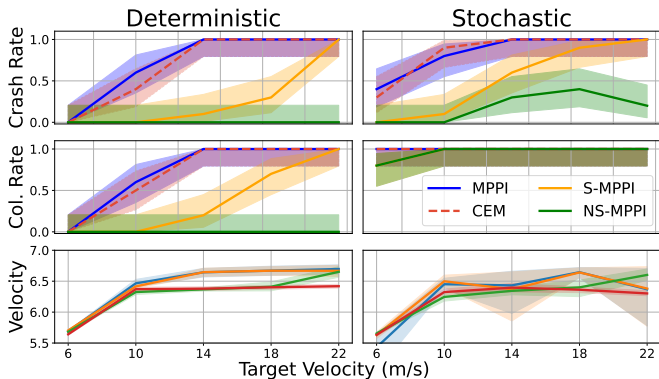


Fig. 5: **Varying target velocities on AutoRally.** Our method NS-MPPI achieves the lowest crash and collision (truncated to Col.) rates under both deterministic and stochastic dynamics. Although the collision rate is close to 1 for every method in the stochastic environment, NS-MPPI achieves a crash rate of near 0.

**Baseline methods.** We compared NS-MPPI against the following sampling-based MPC methods.

- Baseline MPPI (MPPI) [71], which forward simulates a set of randomly sampled trajectories for optimal control.
- Shield MPPI (S-MPPI) [79], which extends MPPI by taking  $h$  in (2) to be a DCBF and by adding the DCBF constraint violation into the cost as in (18).
- CEM [7], which samples trajectories similar to the baseline MPPI but with the weight  $\omega = 1$  for only the  $k$ -lowest cost trajectories and 0 otherwise. This corresponds to an average of the  $k$ -lowest cost trajectories.

In all simulations, we use S-MPPI as the control policy  $\pi$  to learn the DPNCBF. We provide further details on the simulation experiments in Appendix A.

#### A. Simulations on AutoRally

We first compare all methods on the AutoRally [26] testbed, a 1/5 scale autonomous racing car. The goal for this task is to track a given fixed velocity without exiting the track.

The vehicle *collides* when it contacts the track boundary and *crashes* when it fully exceeds the track boundary.

We tested our algorithm under both deterministic dynamics and stochastic dynamics with a Gaussian state disturbance added at each timestep. We visualize the resulting trajectories in Fig. 4 with a target velocity of  $15 \text{ ms}^{-1}$ . This is a very high target velocity, as previous works considered at most velocities of  $8 \text{ ms}^{-1}$  [79] or  $9 \text{ ms}^{-1}$  [23]. Under this challenging speed, both the MPPI and the S-MPPI frequently veer off course. In contrast, the proposed NS-MPPI successfully retains the vehicle within the confines of the track, thereby ensuring safety.

Next, we vary the target velocities, showing the resulting crash rate, collision rate, and velocity in Fig. 5 over 20 trials. With higher velocities, the vehicle has less time to turn, increasing the likelihood of leaving the track and colliding or crashing. We see that NS-MPPI consistently outperforms all other methods in both settings without having a significantly lower average velocity.

#### B. Simulations on Drone

We next tested our algorithm on Drone, a simulated planar quadrotor that incorporates ground effects arising from the intricate interaction between the blade airflow and the ground surface. The goal is for the drone to navigate as fast as possible through a narrow corridor close to the ground. We vary the control horizon for 200 sampled trajectories and plot the results in Fig. 6. Only NS-MPPI has a crash rate of zero over all control horizons. S-MPPI is only able to consistently avoid crashes with longer control horizons. To understand why this is the case, we visualize the trajectories for NS-MPPI and S-MPPI with a control horizon of 10 in Fig. 6b. Due to the DCBF constraint, NS-MPPI starts descending to avoid the obstacles even though none of the original collision constraints have been violated. On the other hand, S-MPPI is unaware of the obstacles beyond the prediction horizon and keeps moving to the right. This suggests that enforcing the DCBF constraint even with only an approximate  $V_\theta^{h,\pi,\gamma}$  is beneficial in enforcing safety with very short control horizons.

To assess the impact of the proposed RBR method, we also implemented Shield MPPI using RBR supervised by a distance-based CBF (S-MPPI+RBR). Interestingly, the success rate of this controller initially increases with the control horizon (as safety increases) but later decreases (as the drone fails to go through the narrow corridor but remains safe). Similarly, the success rate of the NS-MPPI also decreases with the control horizon. This is because the variance of the Monte Carlo estimator scales exponentially with the control horizon, as shown in Section IV-B.

To compare the proposed DPNCBF with existing reachability method for the discrete-time avoid problem, we use [48], a grid-based reachability analysis approach in continuous-time, with the finest discretization allowed by memory and apply the resulting maximal CBF to NS-MPPI, as demonstrated by Fig. 7(b). To the best of our knowledge, there is no existing open-source reachability solver in discrete-time. Simulations show that safety improves slightly at low control horizons



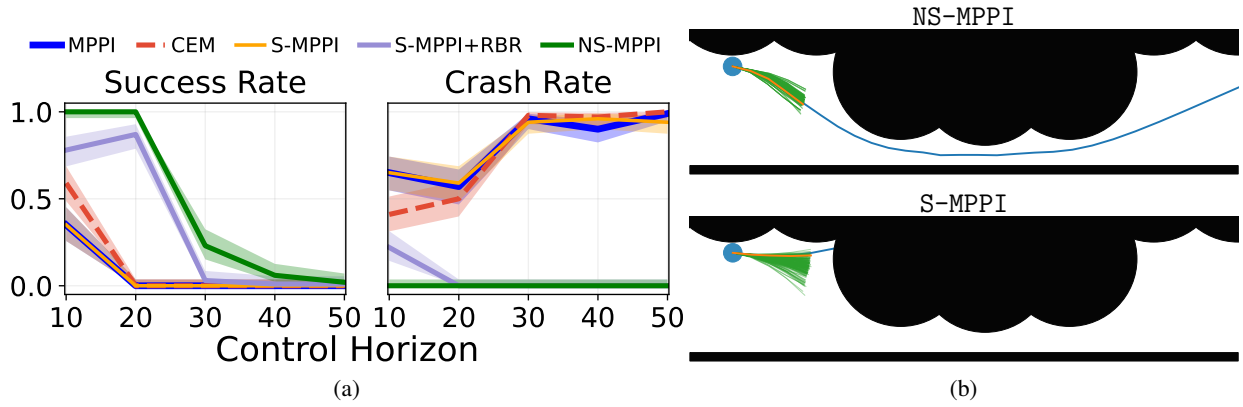


Fig. 6: **Varying control horizon on Drone.** (a) Only NS-MPPI has a crash rate of zero even with a control horizon of 10, demonstrating the benefit of enforcing the DCBF constraint for maintaining safety beyond the prediction horizon. (b) We compare the sampling trajectory  $u^i$  (in green) and estimated optimal trajectory  $\hat{v}$  (in orange) for NS-MPPI and S-MPPI for a control horizon of 10. Due to the DCBF constraint, NS-MPPI starts descending to avoid the obstacles even though none of the original collision constraints have been violated.

but remains imperfect due to errors from mismatches between discrete-time (zero-order hold) and continuous-time control. At longer horizons, success rates drop as state discretization errors mark the narrow corridor as unsafe, blocking traversal, as demonstrated by Fig. 7(a).

### C. Deeper Investigation Into RBR

We performed a case study to better understand the benefits of the proposed RBR method on sample efficiency.

**RBR improves sample efficiency by 5X.** To isolate the effects of RBR without the other improvements, we considered a new method that extends S-MPPI with RBR (S-MPPI w/ RBR), and compared it against the original S-MPPI across different numbers of sampled trajectories on AutoRally over 100 trials in Fig. 8. The results show that S-MPPI with RBR can achieve the same collision rates as S-MPPI without RBR while using 5 times *fewer* trajectories. This matches our expectations, both from the intuition that RBR results in a more closely aligned proposal distribution (e.g., see Fig. 1), and from the theoretical results in Section IV-B that show that RBR improves the quality of the estimator.

**RBR achieves larger  $\widehat{ESS}$ .** RBR theoretically improves the  $\widehat{ESS}$ , as mentioned in Section IV-B. We verified this claim empirically on the AutoRally and plot the results in Fig. 9. While MPPI and S-MPPI exhibit values of  $\widehat{ESS}$  concentrated near 1, the  $\widehat{ESS}$  values for NS-MPPI are more uniformly distributed, indicating a more extensive utilization of the sampled trajectories and thus a better use of the available computational resources.

**RBR is especially beneficial in harder environments.** To further stress-test RBR, we ran 50 trials with 100 sampled trajectories on a harder variant of the AutoRally that contains obstacles of various sizes and compared its performance with other baseline methods. The cluttered AutoRally environment can become especially challenging when the sampled trajectory states fall within a non-convex set caused by close-to-center obstacles. First, consider the trajectories produced by S-MPPI

[79] and our NS-MPPI in Fig. 10. Since S-MPPI relies on a heuristic DCBF, it is unable to avoid crashing into the walls and obstacles, whereas NS-MPPI avoids crashes in almost all runs as shown in Fig. 10 (b). Furthermore, if we add efficient sampling to the S-MPPI, we observe a massive performance boost as seen in Fig. 10, but still not as good as NS-MPPI.

Finally, to demonstrate the enhancement in sampling efficiency achieved by the proposed sampling method, we provide visual representations of trajectory sampling distributions obtained from simulations in environments with obstacles in Fig. 1. As demonstrated by Fig. 1 (b), the RBR concentrates trajectory samples within a feasible narrow passage that satisfies the DCBF safety criteria. Conversely, the standard sampling approach, exemplified by Fig. 1 (a), produces a sparse distribution, misallocating samples to unsafe zones and thereby failing to sufficiently explore safe areas. This inadequacy leads to a collision with an obstacle.

**Enhanced sampling efficiency enables real-time planning on a CPU.** To demonstrate the computational benefits provided by RBR, we evaluated several VI-MPC controllers on a standard laptop CPU. This assessment sets a target tracking speed of 12 m/s, a horizon of  $K = 15$  and a small sample size of  $N = 30$  trajectories at each optimization iteration for fast computation. The resulting statistics are presented in Table II. From the data presented in Table II, it is evident that NS-MPPI attains the minimal crash rate, although it operates at a marginally reduced control update frequency, with S-MPPI next in line. This challenging task results in both MPPI and CEM controllers exhibiting crash and collision rates that approach 1. For most robotic applications, a control rate of 70 Hz is sufficient.

## VI. HARDWARE EXPERIMENTS

Finally, we deployed our method and other existing sampling-based MPC baselines on hardware using the AutoRally experimental platform [26]. Testing was conducted on a dirt track outlined by drainage pipes, as depicted in Fig. 11. All

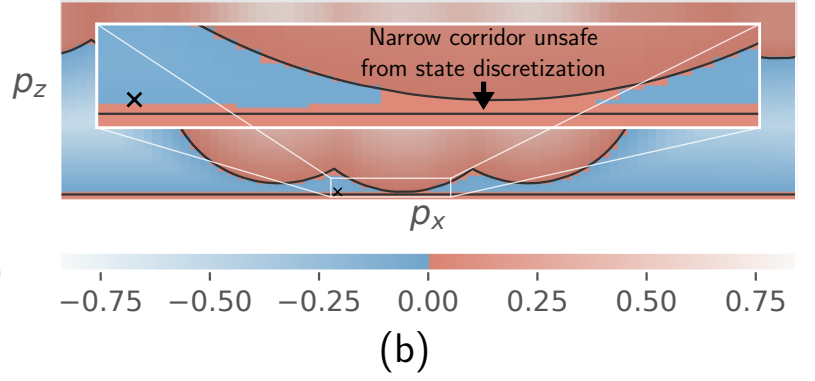
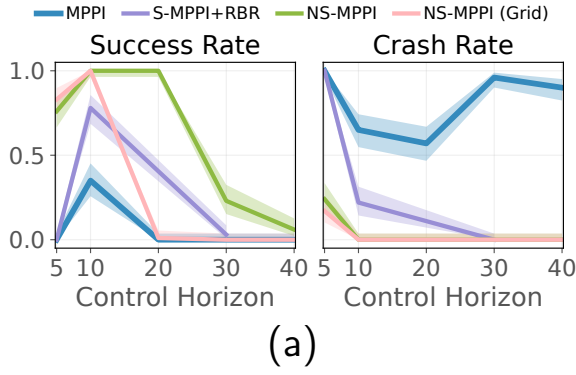


Fig. 7: **Comparison between DPNCBF and CBF obtained from a grid-based reachability method** (a) Using the  $V$  from reachability on Drone. (b) Due to state discretization, the narrow corridor is marked unsafe, causing the drone to be overly conservative and not cross the corridor to remain safe. [Note: Data in Fig.7(a) is not consistent with data in Fig.6(a). Fig.7(b) needs to have a different plot scale, current plot looks very different from Fig.6(b)]

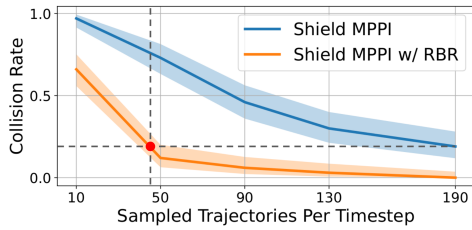


Fig. 8: **RBR reduces the samples needed by 5X.** On AutoRally, applying RBR to S-MPPI significantly reduces the required number of sampled trajectories (190 to 40) to achieve the same level of safety. From another perspective, RBR reduces the collision rate by 74% at 50 sampled trajectories.

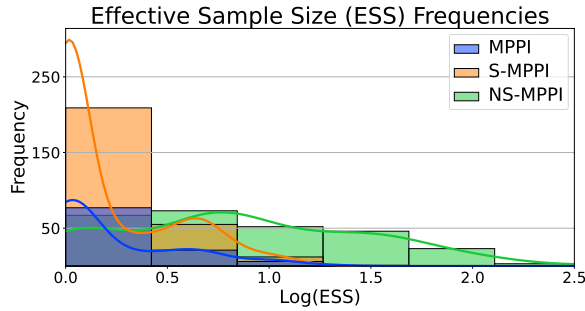


Fig. 9: **NS-MPPI has larger  $\widehat{ESS}$ .** The proposed NS-MPPI achieves larger effective sample sizes, verifying that RBR enables more efficient use of samples and computations.

TABLE II: Performance and Timing Comparison

Controller	Crash Rate	Collision Rate	Control Rate (Hz)
NS-MPPI	0.04	0.06	73.64
S-MPPI	0.46	0.70	107.31
MPPI	0.98	1.00	130.05
CEM	1.00	1.00	88.17

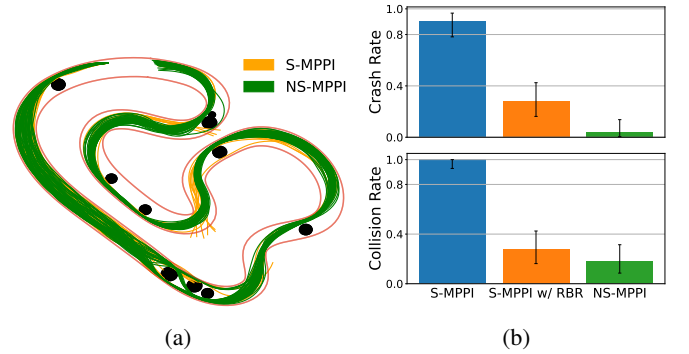


Fig. 10: **AutoRally with obstacles.** (a) While NS-MPPI consistently clears the entire track, S-MPPI often crashes into obstacles and walls. (b) **AutoRally Statistics.** The crash and collision rate statistics supporting the observations made on the left.

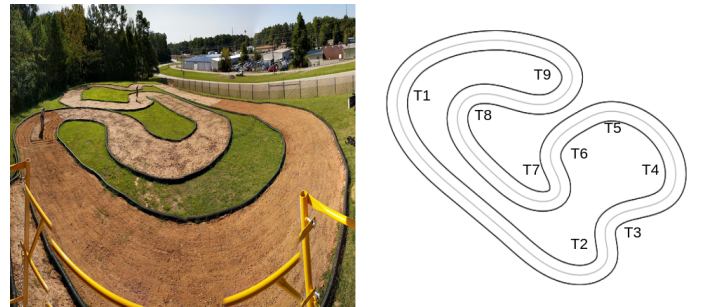


Fig. 11: **AutoRally Track Setup.** The site is equipped with a spacious carport, a 3-meters-tall observation tower, and a storage shed [24]. Each turn of the track is numbered as  $T_i$ , with  $T_1$  is the first turn and  $T_9$  is the last.

computations use the onboard Intel i7-6700 CPU and Nvidia GTX-750ti GPU [26]. See Appendix E for more details on the hardware experiments.

### A. Robustness Against Adversarial Costs

Robots depend on sensors such as cameras to gather information about their environment, and use the collected data to construct suitable cost functions. Deep Neural Networks (DNNs) are frequently utilized for object detection, classification, and semantic segmentation [87]. However, DNNs are susceptible to adversarial attacks [29], which can mislead the algorithms into generating incorrect outputs. These attacks often involve altering the appearance of the object of interest directly [42], eventually resulting in erroneous cost functions. In this hardware experiment, we tested the robustness of the proposed NS-MPPI to these erroneous or misspecified cost functions by using an *adversarial* cost function on the actual AutoRally hardware platform. Namely, instead of penalizing infeasible system states, we *reward* unsafe states by assigning negative costs, effectively incentivizing collisions. We compared against the S-MPPI and MPPI baselines.

The results are depicted in Fig. 12. As MPPI only uses cost information, it rapidly steers the AutoRally vehicle toward the closest boundary, resulting in a crash. S-MPPI crashes the vehicle at the first turn (T1) as well, but it covers a greater distance than MPPI. Only NS-MPPI achieves a collision-free completion of 10 laps. Nonetheless, there are noticeable oscillatory movements of the car, indicating susceptibility to attractive costs driving it toward boundary collisions.

### B. Safety Under Unsafe User Input

Exceeding speed limits is a significant contributor to road accidents [50] due to noncompliance by many drivers. The proposed Shield VIMPC control framework can smartly manage speeds supervised by the NCBF to achieve safety, even given unsafe user inputs. To this end, we tested MPPI, S-MPPI, and the proposed NS-MPPI on the AutoRally hardware using large, unsafe target velocities to examine their safety performance. The results are shown in Fig. 13. Designating high target velocities encourages maximizing speed, thus disrupting the equilibrium in the controllers' cost structure by diminishing the emphasis on cost penalties associated with obstacles. As a result, MPPI causes the vehicle to deviate from the path, resulting in a crash at the second turn (T2), while S-MPPI impacts at the seventh turn (T7). In contrast, NS-MPPI ensures safety, preventing any accidents.

## VII. LIMITATIONS

One limitation of our approach is that the DPNCBF is only an approximation of a DCBF. As such, the typical safety guarantees of DCBFs may not hold if the function is not a true DCBF. Moreover, as noted in Remark 4, the use of variational inference means that  $\mathbf{v}^*$  may not satisfy the state constraints despite the optimal distribution  $p(\mathbf{u} \mid o = 1)$  having zero density at states that violate the state constraints. While theoretically we were only able to show that  $\mathbf{v}^*$  stays safe under the assumption that  $\mathcal{U}_{\text{safe}}^K$  is convex, empirical results in Fig. 10 show that NS-MPPI can perform well despite the fact that  $\mathcal{U}_{\text{safe}}^K$  is typically not convex.

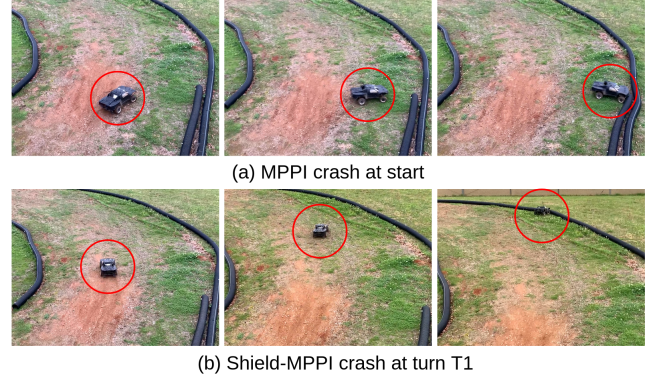


Fig. 12: **Adversarial cost.** With a cost function that rewards crashes, both MPPI and S-MPPI crash. Only NS-MPPI completes 10 laps collision-free.

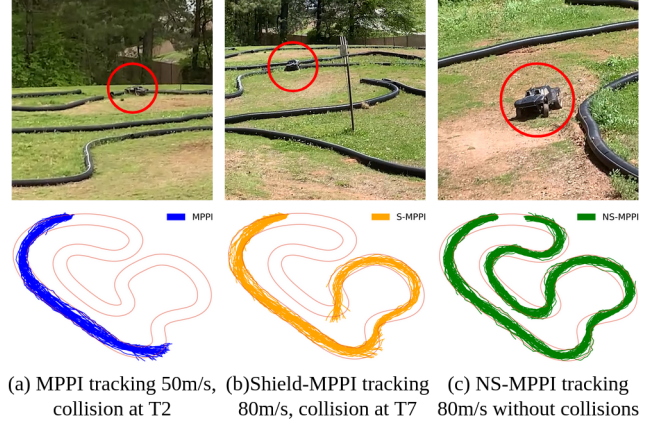


Fig. 13: **Safety under unsafe user inputs.** Using an erroneously large velocity target, both MPPI and S-MPPI crash. The visualization on the bottom row visualizes simulation trajectories that lead to crashes in the same turn.

## VIII. CONCLUSION

In this study, we have adapted the policy neural control barrier function (PNCBF) to a discrete-time setting and utilized the resulting discrete-time policy neural control barrier function (DPNCBF) to supervise the proposed resampling-based rollout (RBR) method. This novel method enhances sampling efficiency and safety for all variational inference model predictive controllers (VIMPCs). Leveraging RBR, we developed the neural shield VIMPC (NS-VIMPC) control framework and demonstrated its benefits for safe planning through the novel neural shield model predictive path integral (NS-MPPI) controller. We conducted tests of the NS-MPPI, benchmarking its performance against state-of-the-art sampling-based MPC controllers using both an autonomous vehicle and a drone. Our simulations and experimental data indicate that the NS-MPPI outperforms existing VIMPC methods in terms of safety and sampling efficiency. The improved sampling efficiency allows NS-MPPI to reach similar performance levels as other VIMPC methods, while using significantly fewer

sampled trajectories, facilitating real-time control on a CPU rather than necessitating costly GPU resources.

We have used the novel RBR method to concentrate sampled trajectories on safe regions. While this approach can significantly improve sampling efficiency in terms of safety, it does not improve other aspects of performance. To further enhance performance, we may, for instance, integrate a Control Lyapunov Function [16] or a Control Lyapunov Barrier Function [74] with the novel RBR approach to sample safe and higher-performance trajectories.

#### REFERENCES

- [1] Manuel Acosta and Stratis Kanarachos. Tire lateral force estimation and grip potential identification using neural networks, extended Kalman filter, and recursive least squares. *Neural Computing and Applications*, 30(11): 3445–3465, 2018.
- [2] Mohamadreza Ahmadi, Andrew Singletary, Joel W. Burdick, and Aaron D. Ames. Safe policy synthesis in multi-agent POMDPs via discrete-time barrier functions. In *IEEE 58th Conference on Decision and Control (CDC)*, pages 4797–4803, Nice, France, Dec. 11 – 13, 2019.
- [3] Ermano Arruda, Michael J. Mathew, Marek Kopicki, Michael Mistry, Morteza Azad, and Jeremy L. Wyatt. Uncertainty averse pushing with model predictive path integral control. In *IEEE-RAS 17th International Conference on Humanoid Robotics (Humanoids)*, pages 497–502, Birmingham, UK, Nov. 15 – 17, 2017.
- [4] Hagai Attias. Planning by probabilistic inference. In *International Workshop on Artificial Intelligence and Statistics*, pages 9–16, Key West, FL, Jan. 3 – 6, 2003.
- [5] Somil Bansal and Claire J Tomlin. Deepreach: A deep learning approach to high-dimensional reachability. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1817–1824. IEEE, 2021.
- [6] Lucas Barcelos, Alexander Lambert, Rafael Oliveira, Paulo Borges, Byron Boots, and Fabio Ramos. Dual Online Stein Variational Inference for Control and Dynamics. In *Proceedings of Robotics: Science and Systems*, Virtual, July 12 – 16, 2021.
- [7] Homanga Bharadhwaj, Kevin Xie, and Florian Shkurti. Model-predictive control via cross-entropy and gradient-based optimization. In *Proceedings of the 2nd Conference on Learning for Dynamics and Control*, volume 120 of *Proceedings of Machine Learning Research*, pages 277–286, Virtual, June 11 – 12, 2020.
- [8] Mohak Bhardwaj, Balakumar Sundaralingam, Arsalan Mousavian, Nathan D. Ratliff, Dieter Fox, Fabio Ramos, and Byron Boots. STORM: An integrated framework for fast joint-space model-predictive control for reactive manipulation. In *Conference on Robot Learning*, pages 750–759, Auckland, New Zealand, Dec. 14 – 18, 2022.
- [9] Camille Brasseur, Alexander Sherikov, Cyrille Collette, Dimitar Dimitrov, and Pierre-Brice Wieber. A robust linear MPC approach to online generation of 3d biped walking motion. In *IEEE-RAS 15th International Conference on Humanoid Robots (Humanoids)*, pages 595–601, Seoul, South Korea, Nov. 3 – 5, 2015.
- [10] José Manuel Bravo, Daniel Limón, Teodoro Alamo, and Eduardo F. Camacho. On the computation of invariant sets for constrained nonlinear systems: An interval arithmetic approach. *Automatica*, 41(9):1583–1589, 2005.
- [11] Wen-Hua Chen, John O’Reilly, and Donald J. Ballance. On the terminal region of model predictive control for nonlinear systems with input/state constraints. *International Journal of Adaptive Control and Signal Processing*, 17(3):195–207, 2003.
- [12] Andrew Clark. Verification and synthesis of control barrier functions. In *IEEE Conference on Decision and Control (CDC)*, pages 6105–6112, Austin, TX, Dec. 13 – 17, 2021.
- [13] Wenceslao Shaw Cortez and Dimos V. Dimarogonas. Safe-by-design control for Euler–Lagrange systems. *Automatica*, 146:110620, 2022.
- [14] Wenceslao Shaw Cortez, Xiao Tan, and Dimos V. Dimarogonas. A robust, multiple control barrier function framework for input constrained systems. *IEEE Control Systems Letters*, 6:1742–1747, 2021.
- [15] Li Danjun, Zhou Yan, Shi Zongying, and Lu Geng. Autonomous landing of quadrotor based on ground effect modelling. In *Chinese Control Conference (CCC)*, pages 5647–5652, Hangzhou, China, July 28 – 30, 2015.
- [16] Charles Dawson, Sicun Gao, and Chuchu Fan. Safe control with learned certificates: A survey of neural Lyapunov, barrier, and contraction methods for robotics and control. *IEEE Transactions on Robotics*, 39(3):1749–1767, 2023.
- [17] Christopher De Sa, Megan Leszczynski, Jian Zhang, Alana Marzoev, Christopher R Aberger, Kunle Olukotun, and Christopher Ré. High-accuracy low-precision training. *arXiv preprint arXiv:1803.03383*, 2018.
- [18] Arnaud Doucet, Nando De Freitas, and Neil Gordon. An introduction to sequential Monte Carlo methods. *Sequential Monte Carlo Methods in Practice*, pages 3–14, 2001.
- [19] Arnaud Doucet, Adam M. Johansen, et al. A tutorial on particle filtering and smoothing: Fifteen years later. *Handbook of Nonlinear Filtering*, 12(656-704):3, 2009.
- [20] Víctor Elvira, Luca Martino, and Christian P. Robert. Rethinking the effective sample size. *International Statistical Review*, 90(3):525–550, 2022.
- [21] Mirko Fiacchini, Teodoro Alamo, and Eduardo F. Camacho. On the computation of local invariant sets for nonlinear systems. In *IEEE Conference on Decision and Control*, pages 3989–3994, New Orleans, LA, Dec. 12 – 14, 2007.
- [22] Jaime F. Fisac, Neil F. Lugovoy, Vicenç Rubies-Royo, Shromona Ghosh, and Claire J. Tomlin. Bridging Hamilton-Jacobi safety analysis and reinforcement learning. In *International Conference on Robotics and Automation (ICRA)*, pages 8550–8556, Montreal, Canada, May 20 – 24, 2019.



- [23] Manan S. Gandhi, Bogdan Vlahov, Jason Gibson, Grady Williams, and Evangelos A. Theodorou. Robust model predictive path integral control: Analysis and performance guarantees. *IEEE Robotics and Automation Letters*, 6(2): 1423–1430, 2021.
- [24] GeorgiaTech. AutoRally. Online. Available: <https://autorally.github.io/>, 2022. Accessed: Nov. 9, 2024.
- [25] Philip E Gill, Walter Murray, and Michael A. Saunders. SNOPT: An SQP algorithm for large-scale constrained optimization. *SIAM Review*, 47(1):99–131, 2005.
- [26] Brian Goldfain, Paul Drews, Changxi You, Matthew Barulic, Orlin Velev, Panagiotis Tsiotras, and James M Rehg. AutoRally: An open platform for aggressive autonomous driving. *IEEE Control Systems Magazine*, 39(1):26–55, 2019.
- [27] Ravi Gondhalekar, Jun-ichi Imura, and Kenji Kashima. Controlled invariant feasibility—A general approach to enforcing strong feasibility in MPC applied to move-blocking. *Automatica*, 45(12):2869–2875, 2009.
- [28] Neil J. Gordon, David J. Salmond, and Adrian F.M. Smith. Novel approach to nonlinear/non-Gaussian Bayesian state estimation. In *IEE Proceedings F (Radar and Signal Processing)*, volume 140, pages 107–113, 1993.
- [29] Amira Guesmi, Muhammad Abdullah Hanif, Bassem Ouni, and Muhammad Shafique. Physical adversarial attacks for camera-based smart systems: Current trends, categorization, applications, research challenges, and future outlook. *IEEE Access*, 11:109617–109668, 2023.
- [30] Emily Hannigan, Bing Song, Gagan Khandate, Maximilian Haas-Heger, Ji Yin, and Matei Ciocarlie. Automatic snake gait generation using model predictive control. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 5101–5107, Paris, France, May 31 – Aug. 31, 2020.
- [31] Tairan He, Chong Zhang, Wenli Xiao, Guanqi He, Changliu Liu, and Guanya Shi. Agile But Safe: Learning Collision-Free High-Speed Legged Locomotion. In *Proceedings of Robotics: Science and Systems*, Delft, Netherlands, July 2024. doi: 10.15607/RSS.2024.XX.059.
- [32] Kai-Chieh Hsu, Vicenç Rubies-Royo, Claire Tomlin, and Jaime F Fisac. Safety and Liveness Guarantees through Reach-Avoid Reinforcement Learning. In *Proceedings of Robotics: Science and Systems*, Virtual, July 2021. doi: 10.15607/RSS.2021.XVII.077.
- [33] Wilson Jallet, Antoine Bambade, Nicolas Mansard, and Justin Carpentier. Constrained differential dynamic programming: A primal-dual augmented Lagrangian approach. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 13371–13378, Kyoto, Japan, Oct. 23 – 27, 2022.
- [34] Ghassen Jerfel, Serena Wang, Clara Wong-Fannjiang, Katherine A. Heller, Yian Ma, and Michael I. Jordan. Variational refinement for importance sampling using the forward kullback-leibler divergence. In *Uncertainty in Artificial Intelligence*, pages 1819–1829, 2021.
- [35] Zichao Jiang, Junyang Jiang, Qinghe Yao, and Gengchao Yang. A neural network-based pde solving algorithm with high precision. *Scientific Reports*, 13(1):4479, 2023.
- [36] Guy Katz, Clark Barrett, David L. Dill, Kyle Julian, and Mykel J. Kochenderfer. Reluplex: An efficient SMT solver for verifying deep neural networks. In *Computer Aided Verification: 29th International Conference (CAV)*, pages 97–117, Heidelberg, Germany, July 24–28, 2017.
- [37] Eric C. Kerrigan and Jan M. Maciejowski. Invariant sets for constrained nonlinear discrete-time systems with application to feasibility in model predictive control. In *Proceedings of the 39th IEEE conference on decision and control*, volume 5, pages 4951–4956, Sydney, Australia, Dec. 12 – 15, 2000.
- [38] Genshiro Kitagawa. Monte Carlo filter and smoother for non-Gaussian nonlinear state space models. *Journal of Computational and Graphical Statistics*, 5(1):1–25, 1996.
- [39] Jacob Knaup, Kazuhide Okamoto, and Panagiotis Tsiotras. Safe high-performance autonomous off-road driving using covariance steering stochastic model predictive control. *IEEE Transactions on Control Systems Technology*, 31(5):2066–2081, 2023.
- [40] Alexander Lambert, Fabio Ramos, Byron Boots, Dieter Fox, and Adam Fishman. Stein variational model predictive control. In *Proceedings of the Conference on Robot Learning*, volume 155 of *Proceedings of Machine Learning Research*, pages 1278–1297, 16–18 Nov, 2021.
- [41] Sergey Levine. Reinforcement learning and control as probabilistic inference: Tutorial and review. *arXiv preprint arXiv:1805.00909*, 2018.
- [42] Juncheng Li, Frank Schmidt, and Zico Kolter. Adversarial camera stickers: A physical camera-based attack on deep learning systems. In *International Conference on Machine Learning*, pages 3896–3904, Long Beach, CA, Jun. 10 – 15, 2019.
- [43] Lars Lindemann and Dimos V. Dimarogonas. Control barrier functions for signal temporal logic tasks. *IEEE Control Systems Letters*, 3(1):96–101, 2018.
- [44] Björn Lindqvist, Sina Sharif Mansouri, Ali-akbar Aghamohammadi, and George Nikolakopoulos. Nonlinear MPC for collision avoidance and control of UAVs with dynamic obstacles. *IEEE Robotics and Automation Letters*, 5(4):6001–6008, 2020.
- [45] Johan Löfberg. Oops! I cannot do it again: Testing for recursive feasibility in MPC. *Automatica*, 48(3):550–555, 2012.
- [46] Kostas Margellos and John Lygeros. Hamilton–Jacobi formulation for reach–avoid differential games. *IEEE Transactions on Automatic Control*, 56(8):1849–1861, 2011.
- [47] David Q. Mayne, James B. Rawlings, Christopher V. Rao, and Pierre O.M. Scokaert. Constrained model predictive control: Stability and optimality. *Automatica*, 36(6):789–814, 2000.
- [48] Ian M. Mitchell. The flexible, extensible and efficient toolbox of level set methods. *Journal of Scientific Computing*, 35:300–329, 2008.

- [49] Anusha Nagabandi, Kurt Konolige, Sergey Levine, and Vikash Kumar. Deep dynamics models for learning dexterous manipulation. In *Conference on Robot Learning*, pages 1101–1112, Virtual, Nov. 16 – 18, 2020.
- [50] National Highway Traffic Safety Administration (NHTSA). Speeding. <https://www.nhtsa.gov/risky-driving/speeding>. Accessed: Jan. 27, 2025.
- [51] Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer, 1999.
- [52] Masashi Okada and Tadahiro Taniguchi. Variational inference MPC for Bayesian model-based reinforcement learning. In *Conference on robot learning*, pages 258–272, Virtual, Nov. 16 – 18, 2020.
- [53] Art B. Owen. Monte Carlo theory, methods and examples, 2013.
- [54] Thomas Power and Dmitry Berenson. Keep it simple: Data-efficient learning for controlling complex systems with simple models. *IEEE Robotics and Automation Letters*, 6(2):1184–1191, 2021.
- [55] Thomas Power and Dmitry Berenson. Variational Inference MPC using Normalizing Flows and Out-of-Distribution Projection. In *Proceedings of Robotics: Science and Systems*, New York City, NY, June 27 – July 1, 2022.
- [56] Thomas Power and Dmitry Berenson. Learning a generalizable trajectory sampling distribution for model predictive control. *IEEE Transactions on Robotics*, 40: 2111–2127, 2024.
- [57] Jintasit Pravitra, Kasey A. Ackerman, Chengyu Cao, Naira Hovakimyan, and Evangelos A. Theodorou. L1-adaptive MPPI architecture for robust and agile control of multirotors. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 7661–7666, Las Vegas, NV, Oct. 25 – 29, 2020.
- [58] Mohamad Qadri, Paloma Sodhi, Joshua G. Mangelson, Frank Dellaert, and Michael Kaess. InCOpt: Incremental constrained optimization using the Bayes tree. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 6381–6388, Kyoto, Japan, Oct. 23 – 27, 2022.
- [59] Zengyi Qin, Kaiqing Zhang, Yuxiao Chen, Jingkai Chen, and Chuchu Fan. Learning safe multi-agent control with decentralized neural barrier certificates. *arXiv preprint arXiv:2101.05436*, 2021.
- [60] Konrad Rawlik, Marc Toussaint, and Sethu Vijayakumar. On stochastic optimal control and reinforcement learning by approximate inference. In *International Joint Conference on Artificial Intelligence*, Beijing, China, Aug. 3 – 9, 2013.
- [61] Reuven Rubinfeld. The cross-entropy method for combinatorial and continuous optimization. *Methodology and Computing in Applied Probability*, 1:127–190, 1999.
- [62] Oswin So and Chuchu Fan. Solving stabilize-avoid optimal control via epigraph form and deep reinforcement learning. In *Robotics: Science and Systems*, Daegu, South Korea, July 10 – 14, 2023.
- [63] Oswin So, Ziyi Wang, and Evangelos A. Theodorou. Maximum entropy differential dynamic programming. In *International Conference on Robotics and Automation (ICRA)*, pages 3422–3428, Philadelphia, PA, May 23 – 27, 2022.
- [64] Oswin So, Zachary Serlin, Makai Mann, Jake Gonzales, Kwesi Rutledge, Nicholas Roy, and Chuchu Fan. How to train your neural control barrier function: Learning safety filters for complex input-constrained systems. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 11532–11539, Yokohama, Japan, May 13 – 17, 2024.
- [65] Richard S. Sutton and Andrew G. Barto. Reinforcement learning: An introduction, 2018.
- [66] Chuyuan Tao, Hunmin Kim, Hyungjin Yoon, Naira Hovakimyan, and Petros Voulgaris. Control barrier function augmentation in sampling-based control algorithm for sample efficiency. In *American Control Conference (ACC)*, pages 3488–3493, Atlanta, GA, June 8 – 10, 2022.
- [67] Marc Toussaint and Amos Storkey. Probabilistic inference for solving discrete and continuous state Markov decision processes. In *Proceedings of the 23rd International Conference on Machine Learning*, pages 945–952, Pittsburgh, PA, June 25 – 29, 2006.
- [68] Efstathios Velenis, Emilio Frazzoli, and Panagiotis Tsioras. Steady-state cornering equilibria and stabilisation for a vehicle during extreme operating conditions. *International Journal of Vehicle Autonomous Systems*, 8(2-4): 217–241, 2010.
- [69] Ziyi Wang, Oswin So, Jason Gibson, Bogdan Vlahov, Manan Gandhi, Guan-Hong Liu, and Evangelos Theodorou. Variational Inference MPC using Tsallis Divergence. In *Proceedings of Robotics: Science and Systems*, Virtual, July 12 – 16, 2021.
- [70] Tianhao Wei and Changliu Liu. Safe control with neural network dynamic models. In *Learning for Dynamics and Control Conference*, pages 739–750, Palo Alto, CA, June 23 – 24, 2022.
- [71] Grady Williams, Paul Drews, Brian Goldfain, James M. Rehg, and Evangelos A. Theodorou. Information-theoretic model predictive control: Theory and applications to autonomous driving. *IEEE Transactions on Robotics*, 34(6):1603–1622, 2018.
- [72] Grady Williams, Brian Goldfain, Paul Drews, Kamil Saigol, James M. Rehg, and Evangelos A. Theodorou. Robust sampling based model predictive control with sparse objective information. In *Robotics: Science and Systems*, volume 14, Pittsburgh, PA, June 26 – 30, 2018.
- [73] Sean Wilson, Paul Glotfelter, Li Wang, Siddharth Mayya, Gennaro Notomista, Mark Mote, and Magnus Egerstedt. The Robotarium: Globally impactful opportunities, challenges, and lessons learned in remote-access, distributed control of multirobot systems. *IEEE Control Systems Magazine*, 40(1):26–44, 2020.
- [74] Zhe Wu, Fahad Albalawi, Zhihao Zhang, Junfeng Zhang, Helen Durand, and Panagiotis D. Christofides. Control

- Lyapunov-barrier function-based model predictive control of nonlinear systems. In *American Control Conference (ACC)*, pages 5920–5926, Milwaukee, WI, June 27 – 29, 2018.
- [75] Mandy Xie, Alejandro Escontrela, and Frank Dellaert. A factor-graph approach for optimization problems with dynamics constraints. *arXiv preprint arXiv:2011.06194*, 2020.
- [76] Bin Xu and Koushil Sreenath. Safe teleoperation of dynamic UAVs through control barrier functions. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 7848–7855, Brisbane, Australia, May 21 – 26, 2018.
- [77] Xiangru Xu, Paulo Tabuada, Jessy W. Grizzle, and Aaron D. Ames. Robustness of control barrier functions for safety critical control. *IFAC-PapersOnLine*, 48(27): 54–61, 2015.
- [78] Ji Yin, Zhiyuan Zhang, Evangelos Theodorou, and Panagiotis Tsiotras. Trajectory distribution control for model predictive path integral control using covariance steering. In *International Conference on Robotics and Automation (ICRA)*, pages 1478–1484, Philadelphia, PA, May 23 – 27, 2022.
- [79] Ji Yin, Charles Dawson, Chuchu Fan, and Panagiotis Tsiotras. Shield model predictive path integral: A computationally efficient robust MPC method using control barrier functions. *IEEE Robotics and Automation Letters*, 8(11):7106–7113, 2023.
- [80] Ji Yin, Zhiyuan Zhang, and Panagiotis Tsiotras. Risk-aware model predictive path integral control using conditional value-at-risk. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 7937–7943, London, UK, May 29 – June 2, 2023.
- [81] Ji Yin, Panagiotis Tsiotras, and Karl Berntorp. Chance-constrained information-theoretic stochastic model predictive control with safety shielding. *63rd IEEE Conference on Decision and Control*, pages 653–658, Dec. 16–19, 2024.
- [82] Changxi You and Panagiotis Tsiotras. Vehicle modeling and parameter estimation using adaptive limited memory joint-state UKF. In *American Control Conference*, pages 322–327, Seattle, WA, May 24–26, 2017.
- [83] Hongzhan Yu, Chiaki Hirayama, Chenning Yu, Sylvia Herbert, and Sicun Gao. Sequential neural barriers for scalable dynamic obstacle avoidance. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 11241–11248, Detroit, MI, Oct. 1 – 5, 2023.
- [84] Qinsheng Zhang, Amirhossein Taghvaei, and Yongxin Chen. An optimal control approach to particle filtering. *Automatica*, 151:110894, 2023.
- [85] Songyuan Zhang, Kunal Garg, and Chuchu Fan. Neural graph control barrier functions guided distributed collision-avoidance multi-agent control. In *Conference on Robot Learning*, pages 2373–2392, Atlanta, GA, Nov. 6 – 9, 2023.
- [86] Songyuan Zhang, Oswin So, Kunal Garg, and Chuchu Fan. GCBF+: A neural graph control barrier function framework for distributed safe multi-agent control. *IEEE Transactions on Robotics*, 2025.
- [87] Bo Zhao, Jiashi Feng, Xiao Wu, and Shuicheng Yan. A survey on deep learning-based fine-grained object classification and semantic segmentation. *International Journal of Automation and Computing*, 14(2):119–135, 2017.

## APPENDIX A

### DETAILS ON THE SIMULATION EXPERIMENTS

#### A1 AutoRally Environment

The numerical simulations use the AutoRally platform, a 1/5-scale autonomous vehicle system capable of executing high-speed, limit-handling maneuvers in complicated environments [26]. Figure 14 shows the AutoRally hardware. In this study, we define a crash as an event when the vehicle makes contact with the track boundary, constructed from soft drainage pipes, resulting in the inability to continue driving. We define a collision as an event when the vehicle contacts the track boundary but can continue driving.



Fig. 14: AutoRally chassis and compute box [24].

#### A2 AutoRally Dynamics Modelling

We model the AutoRally using a rear-wheel drive, single-track bicycle model [39]. The vehicle system is represented in curvilinear coordinates using the centerline of the track as the reference curve. Using a curvilinear coordinate system provides a more intuitive interpretation of the vehicle's position and heading relative to the track, as compared to Cartesian coordinates.

To this end, we consider a nonlinear, continuous-time system,

$$\dot{x} = F(x, u), \quad (\text{A.1})$$

where the system state of the AutoRally is given by,

$$x = [v_x, v_y, \dot{\psi}, \omega_F, \omega_R, e_\psi, e_y, s]^\top, \quad (\text{A.2})$$

and where  $v_x$  is the longitudinal velocity,  $v_y$  is the lateral velocity, and  $\dot{\psi}$  is the yaw rate. The front and rear wheel angular velocities are denoted by  $\omega_F$  and  $\omega_R$ . The yaw error and the lateral distance error from the centerline of the track are denoted by  $e_\psi$  and  $e_y$ , respectively (see Fig. 15). The variable  $s$  is the curvilinear position along the track centerline. The control input to the system (A.1) is,

$$u = [\delta, T]^\top, \quad (\text{A.3})$$

where  $\delta$  is the steering angle and  $T$  denotes the values for throttle input (if positive) or braking input (if negative). The dynamics of a single-track dynamic bicycle model [39] used to model (A.1) are given by,

$$\dot{v}_x = \frac{f_{Fx} \cos \delta - f_{Fy} \sin \delta + f_{Rx}}{m} + v_y \dot{\psi}, \quad (\text{A.4a})$$

$$\dot{v}_y = \frac{f_{Fx} \sin \delta + f_{Fy} \cos \delta + f_{Ry}}{m} - v_x \dot{\psi}, \quad (\text{A.4b})$$

$$\ddot{\psi} = \frac{(f_{Fy} \cos \delta + f_{Fx} \sin \delta) \ell_F - f_{Ry} \ell_R}{I_z}, \quad (\text{A.4c})$$

$$\dot{\omega}_F = -\frac{r_F}{I_{\omega F}} f_{Fx}, \quad (\text{A.4d})$$

$$\dot{\omega}_R = \Theta(\omega_R, T), \quad (\text{A.4e})$$

$$\dot{e}_\psi = \dot{\psi} - \frac{v_x \cos e_\psi - v_y \sin e_\psi}{1 - \rho(s) e_y} \rho(s), \quad (\text{A.4f})$$

$$\dot{e}_y = v_x \sin e_\psi + v_y \cos e_\psi, \quad (\text{A.4g})$$

$$\dot{s} = \frac{v_x \cos e_\psi - v_y \sin e_\psi}{1 - \rho(s) e_y}, \quad (\text{A.4h})$$

where  $m$  and  $I_z$  are the vehicle's mass and moment of inertia about the vertical axis, respectively. The radius of the front wheel is  $r_F$ , and the moment of inertia of the front wheel about the front axle is  $I_{\omega F}$ . The curvature of the track centerline



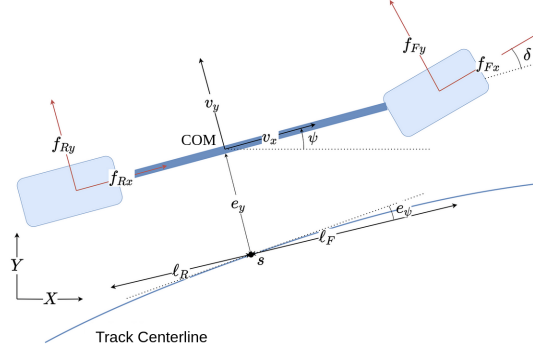


Fig. 15: Schematic of the dynamic bicycle model.

at position  $s$  is  $\rho(s)$ . The front and rear tire frictional forces are denoted by  $f_{Fx}, f_{Fy}, f_{Rx}, f_{Ry}$ , where the subscripts  $F, R$  indicate front and rear tires and  $x, y$  indicate longitudinal and lateral directions. These frictional forces are computed using the ellipse model in [68],

$$f_{ix} = f_{iz}\mu_{ix}, \quad (\text{A.5a})$$

$$f_{iy} = f_{iz}\mu_{iy} + \Phi_i(\delta, \alpha_i), \quad (\text{A.5b})$$

where  $i = F, R$  indicates front or rear wheels,  $f_{iz}$  is the normal force acting on the tire. The tire friction coefficients  $\mu_{ix}, \mu_{iy}$  are dependent on wheel slip angles and the intrinsic tire parameters.

Inspired by [1], we use a neural network to model the residual error  $\Phi_i(\delta, \alpha_i)$ , where  $\delta$  is steering angle,  $\alpha_i = \arctan(v_{iy}/v_{ix})$  is wheel slip angle.  $\Phi_i$  is trained using experimental data collected from the AutoRally. The rear wheel angular acceleration  $\dot{\omega}_R$  is given by (A.4e). Modelling  $\dot{\omega}_R$  is challenging, because it is determined by a variety of mechanical and physical conditions, including nonlinear tire behavior, dynamic load transfer, road surface irregularities and the complex interaction between these components. To this end, we use a data-driven approach and utilize a neural network  $\Theta(\cdot, \cdot)$  trained on experimental data to model the rear wheel angular acceleration  $\dot{\omega}_R$ .

The nonlinear, continuous-time system (A.1) is discretized and converted to a discrete-time system using Euler integration,

$$x_{k+1} = f(x_k, u_k) = x_k + F(x_k, u_k)\Delta t, \quad (\text{A.6})$$

where  $\Delta t = t_{k+1} - t_k$  is the discretization time interval. We employ an Unscented Kalman Filter (UKF) on the discrete-time system (A.6) to identify the vehicle and tire parameters using the approach in [82]. The resulting parameters of the vehicle model include the vehicle mass  $m = 22$  kg, moment of inertia  $I_z = 1.1$  kg  $\cdot$  m<sup>2</sup>,  $l_F = 0.34$  m,  $l_R = 0.23$  m,  $I_{\omega F} = 0.10$  kg  $\cdot$  m<sup>2</sup> and front wheel radius  $r_F = 0.095$  m. The parameters in Pacejka's Magic Formula used by the tire friction ellipse model [68] are computed as  $\text{tire}_B = 4.1$ ,  $\text{tire}_C = 0.95$ ,  $\text{tire}_D = 1.1$ .

### A3 Drone Dynamics Modelling

We model the two-dimensional drone dynamics using the following equations,

$$\dot{v}_x = -\frac{F}{m} \sin \theta, \quad (\text{A.7a})$$

$$\dot{v}_z = \frac{F}{m} \cos \theta - g, \quad (\text{A.7b})$$

$$\ddot{\theta} = \frac{\tau}{I}, \quad (\text{A.7c})$$

$$(\text{A.7d})$$

where the total thrust  $F$  and the torque  $\tau$  are given by,

$$F = F_1 + F_2, \quad \tau = \ell(F_2 - F_1), \quad (\text{A.8})$$

where  $\ell$  is the length from the center of the drone to the center of each rotor, and  $F_1, F_2$  are thrust forces generated by the left and right rotors, respectively. The thrust forces are computed considering the ground effect [15], given by,

$$F_1 = F_{in1} / \left(1 - \rho \left(\frac{r}{4z_{r1}}\right)\right), \quad (\text{A.9})$$

$$F_2 = F_{in2} / \left(1 - \rho \left(\frac{r}{4z_{r2}}\right)\right), \quad (\text{A.10})$$

where  $F_{in1}$ ,  $F_{in2}$  are the input thrust commands for the left and right rotors, and  $\rho$  is an intrinsic property of the drone. The radius of the rotors is  $r$  and the heights from the rotors to the ground are  $z_{r1}$  and  $z_{r2}$ , which are given by,

$$z_{r1} = z - \ell \sin \theta, \quad (\text{A.11})$$

$$z_{r2} = z + \ell \sin \theta, \quad (\text{A.12})$$

where  $\theta$  is the pitch angle.

## APPENDIX B PROOFS FOR VARIATIONAL INFERENCE MPC

### B1 Derivations of the Variational Inference Updates

Let  $q_{\mathbf{v}}$  be the density function of a Gaussian distribution with mean  $\mathbf{v}$  and covariance  $\Sigma$ , i.e.,

$$q_{\mathbf{v}}(\mathbf{u}) = Z_q^{-1} \exp \left( -\frac{1}{2} (\mathbf{u} - \mathbf{v})^\top \Sigma^{-1} (\mathbf{u} - \mathbf{v}) \right), \quad (\text{B.1})$$

where the normalization constant  $Z_q$  is independent of the mean  $\mathbf{v}$ . We wish to solve the variational inference problem

$$\min_{\mathbf{v}} \mathbb{D}_{\text{KL}}(p(\mathbf{u} \mid o = 1) \parallel q_{\mathbf{v}}(\mathbf{u})). \quad (\text{B.2})$$

Expanding and ignoring terms unrelated to  $\mathbf{v}$ , we get

$$\mathbb{D}_{\text{KL}}(p(\mathbf{u} \mid o = 1) \parallel q_{\mathbf{v}}(\mathbf{u})) = \int p(\mathbf{u} \mid o = 1) \log \left( \frac{p(\mathbf{u} \mid o = 1)}{q_{\mathbf{v}}(\mathbf{u})} \right) d\mathbf{u} \quad (\text{B.3})$$

$$= \int p(\mathbf{u} \mid o = 1) \log(p(\mathbf{u} \mid o = 1)) d\mathbf{u} - \int p(\mathbf{u} \mid o = 1) \log(q_{\mathbf{v}}(\mathbf{u})) d\mathbf{u} \quad (\text{B.4})$$

$$= - \int p(\mathbf{u} \mid o = 1) \log(q_{\mathbf{v}}(\mathbf{u})) d\mathbf{u} + O(1) \quad (\text{B.5})$$

$$= \int p(\mathbf{u} \mid o = 1) \left( \frac{1}{2} (\mathbf{u} - \mathbf{v})^\top \Sigma^{-1} (\mathbf{u} - \mathbf{v}) + \log Z_q \right) d\mathbf{u} + O(1) \quad (\text{B.6})$$

$$= \mathbb{E}_{p(\mathbf{u} \mid o=1)} \left[ \frac{1}{2} (\mathbf{u} - \mathbf{v})^\top \Sigma^{-1} (\mathbf{u} - \mathbf{v}) + \log Z_q \right] + O(1) \quad (\text{B.7})$$

$$= \mathbb{E}_{p(\mathbf{u} \mid o=1)} \left[ \frac{1}{2} (\mathbf{u} - \mathbf{v})^\top \Sigma^{-1} (\mathbf{u} - \mathbf{v}) \right] + O(1). \quad (\text{B.8})$$

Taking the derivative with respect to  $\mathbf{v}$  and setting to zero, we get

$$\mathbb{E}_{p(\mathbf{u} \mid o=1)} [\Sigma^{-1} (\mathbf{u} - \mathbf{v}^*)] = 0. \quad (\text{B.9})$$

Rearranging, yields

$$\mathbf{v}^* = \mathbb{E}_{p(\mathbf{u} \mid o=1)} [\mathbf{u}]. \quad (\text{B.10})$$

### B2 Derivations of the Self Normalized Importance Sampling Estimator

We next derive the self-normalized importance sampling (SNIS) estimator  $\hat{\mathbf{v}}$  in (9) of  $\mathbf{v}^*$  in (8). Note that, since  $Z = \int \exp(-J(\mathbf{u})) p_0(\mathbf{u}) d\mathbf{u}$  is unknown, we cannot compute the weights  $\omega(\mathbf{u})$  in (8) directly. Hence, the regular importance sampled Monte Carlo estimator

$$\hat{\mathbf{v}} = \frac{1}{N} \sum_{i=1}^N \omega(\mathbf{u}^i) \mathbf{u}^i, \quad (\text{B.11})$$

is not computable. Instead, we can use the SNIS estimator [53], which can be derived as follows. First, note that  $Z$  can be written as the expectation of  $\exp(-J(\mathbf{u}))$ , i.e.,

$$Z = \mathbb{E}_{p_0(\mathbf{u})} [\exp(-J(\mathbf{u}))]. \quad (\text{B.12})$$

Hence, one idea is to reuse the samples  $\mathbf{u}^i \sim r$  to compute an estimate  $\hat{Z}$  of  $Z$ , that is,

$$\hat{Z} = \frac{1}{N} \sum_{i=1}^N \exp(-J(\mathbf{u}^i)) p_0(\mathbf{u}^i) / r(\mathbf{u}^i), \quad (\text{B.13})$$

where  $\hat{Z}$  is a “normal” importance sampled Monte Carlo estimator of  $Z$ . If we use the estimate  $\hat{Z}$  to compute the weights  $\hat{\omega}(\mathbf{u})$  in the normal importance sampled Monte Carlo estimator  $\hat{\mathbf{v}}$  (B.11), we obtain the SNIS estimator as in (9) and (8)

$$\hat{\mathbf{v}} = \frac{1}{N} \sum_{i=1}^N \hat{\omega}(\mathbf{u}^i) \mathbf{u}^i, \quad (\text{B.14})$$

$$= \frac{1}{N} \sum_{i=1}^N \left( \frac{1}{\hat{Z}} \frac{\exp(-J(\mathbf{u}^i)) p_0(\mathbf{u}^i)}{r(\mathbf{u}^i)} \right) \mathbf{u}^i, \quad (\text{B.15})$$

$$= \frac{1}{N} \sum_{i=1}^N \left( \frac{1}{\frac{1}{N} \sum_{j=1}^N \exp(-J(\mathbf{u}^j)) p_0(\mathbf{u}^j) / r(\mathbf{u}^j)} \frac{\exp(-J(\mathbf{u}^i)) p_0(\mathbf{u}^i)}{r(\mathbf{u}^i)} \right) \mathbf{u}^i, \quad (\text{B.16})$$

$$= \frac{1}{N} \sum_{i=1}^N \left( \frac{1}{\frac{1}{N} \sum_{j=1}^N Z^{-1} \exp(-J(\mathbf{u}^j)) p_0(\mathbf{u}^j) / r(\mathbf{u}^j)} \frac{Z^{-1} \exp(-J(\mathbf{u}^i)) p_0(\mathbf{u}^i)}{r(\mathbf{u}^i)} \right) \mathbf{u}^i, \quad (\text{B.17})$$

$$= \frac{1}{N} \sum_{i=1}^N \left( \frac{1}{\frac{1}{N} \sum_{j=1}^N Z^{-1} \exp(-J(\mathbf{u}^j)) p_0(\mathbf{u}^j) / r(\mathbf{u}^j)} \omega(\mathbf{u}^i) \right) \mathbf{u}^i, \quad (\text{B.18})$$

$$= \frac{1}{N} \sum_{i=1}^N \left( \frac{\omega(\mathbf{u}^i)}{\frac{1}{N} \sum_{j=1}^N \omega(\mathbf{u}^j)} \right) \mathbf{u}^i, \quad (\text{B.19})$$

$$= \sum_{i=1}^N \left( \frac{\omega(\mathbf{u}^i)}{\sum_{j=1}^N \omega(\mathbf{u}^j)} \right) \mathbf{u}^i, \quad (\text{B.20})$$

$$= \sum_{i=1}^N \tilde{\omega}^i \mathbf{u}^i, \quad (\text{B.21})$$

where,

$$\tilde{\omega}^i = \frac{\omega(\mathbf{u}^i)}{\sum_{j=1}^N \omega(\mathbf{u}^j)}. \quad (\text{B.22})$$

Unfortunately, the use of the estimated  $\hat{Z}$  in the weights  $\omega(\mathbf{u})$  causes the SNIS estimator to be biased [53].

### B3 MPPI as a Special Case of the Variational Inference Update

**Theorem 6.** When  $p_0(\mathbf{u}) = q_0(\mathbf{u}) := q_{\mathbf{v}=0}(\mathbf{u})$  and  $r(\mathbf{u}) = q_{\bar{\mathbf{v}}}(\mathbf{u})$  for some previous estimate of the optimal control sequence  $\bar{\mathbf{v}}$ , the equation for the weights  $\omega(\mathbf{u})$  in the variational inference MPC update (8) reduces to that of MPPI.

*Proof:* Let  $q_{\mathbf{v}}$  be the density function for a Gaussian distribution with mean  $\mathbf{v}$  and covariance  $\Sigma$ , that is,

$$q_{\mathbf{v}}(\mathbf{u}) = Z_q^{-1} \exp \left( -\frac{1}{2} (\mathbf{u} - \mathbf{v})^\top \Sigma^{-1} (\mathbf{u} - \mathbf{v}) \right), \quad (\text{B.23})$$

By choosing  $p_0(\mathbf{u}) = q_0(\mathbf{u})$  and  $r(\mathbf{u}) = q_{\bar{\mathbf{v}}}(\mathbf{u})$ , we have that

$$\frac{p_0(\mathbf{u})}{r(\mathbf{u})} = \frac{Z_q^{-1} \exp \left( -\frac{1}{2} \mathbf{u}^\top \Sigma^{-1} \mathbf{u} \right)}{Z_q^{-1} \exp \left( -\frac{1}{2} (\mathbf{u} - \bar{\mathbf{v}})^\top \Sigma^{-1} (\mathbf{u} - \bar{\mathbf{v}}) \right)} \quad (\text{B.24})$$

$$= \exp \left( -\frac{1}{2} \mathbf{u}^\top \Sigma^{-1} \mathbf{u} + \frac{1}{2} (\mathbf{u} - \bar{\mathbf{v}})^\top \Sigma^{-1} (\mathbf{u} - \bar{\mathbf{v}}) \right) \quad (\text{B.25})$$

$$= \exp \left( -\frac{1}{2} \mathbf{u}^\top \Sigma^{-1} \mathbf{u} + \frac{1}{2} \mathbf{u}^\top \Sigma^{-1} \mathbf{u} - \mathbf{u}^\top \Sigma^{-1} \bar{\mathbf{v}} + \frac{1}{2} \bar{\mathbf{v}}^\top \Sigma^{-1} \bar{\mathbf{v}} \right) \quad (\text{B.26})$$

$$= \exp \left( -\mathbf{u}^\top \Sigma^{-1} \bar{\mathbf{v}} + \frac{1}{2} \bar{\mathbf{v}}^\top \Sigma^{-1} \bar{\mathbf{v}} \right). \quad (\text{B.27})$$

Hence,

$$\omega(\mathbf{u}) = \frac{Z^{-1} \exp(-J(\mathbf{u})) p_0(\mathbf{u})}{r(\mathbf{u})} \quad (\text{B.28})$$

$$= Z^{-1} \exp(-[J(\mathbf{u}) + \mathbf{u}^\top \Sigma^{-1} \bar{\mathbf{v}}]) \exp\left(\frac{1}{2} \bar{\mathbf{v}}^\top \Sigma^{-1} \bar{\mathbf{v}}\right) \quad (\text{B.29})$$

$$\propto \exp(-[J(\mathbf{u}) + \mathbf{u}^\top \Sigma^{-1} \bar{\mathbf{v}}]), \quad (\text{B.30})$$

where the last line follows since terms not dependent on  $\mathbf{u}$  can be canceled out between the numerator and denominator in (10).  $\blacksquare$

## APPENDIX C NCBF TRAINING DETAILS

The definition of the avoid set is given by (2),

$$\mathcal{A} = \{x|h(x) > 0\},$$

where the avoidance heuristic  $h(x)$  is used to define the avoid set. Furthermore, the definition of the policy value function  $V^{h,\pi}(x_0)$  given by (19), and the ensuing training losses (IV-A), (25) are all dependent on  $h(x)$ . A suitable selection of the avoidance heuristic  $h(x)$  must ensure that its corresponding avoidance set  $\mathcal{A}$  encompasses all system states that overlap with any obstacles. Additionally, the heuristic should improve the quality of information captured by the loss function (25), thereby making it more effective for neural network training. To this end, a reasonable choice of the avoidance heuristic utilizes coordinates in (A.2),

$$h_0(x) = w_I^2 - e_y^2, \quad (\text{C.1})$$

where  $w_I = 1.5$  m is half of the track width, and  $e_y$  is the lateral deviation of the vehicle CoM to the track centerline. The visualization of  $h_0(x)$  is demonstrated by the orange curve in Fig. 16.

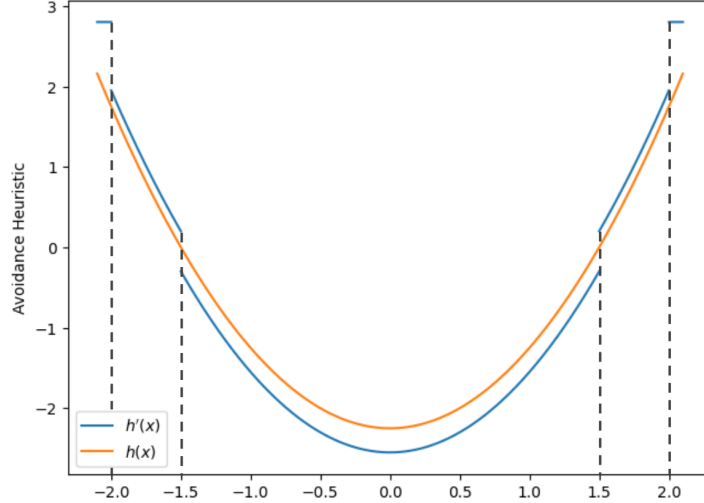


Fig. 16: Avoidance heuristic visualization. The orange curve shows the original avoidance heuristic  $h_0$ , while the blue curve demonstrates the modified avoidance heuristic  $h$  used for training the NCBF.

However, the neural network approximation  $V_\theta^{h_0,\pi}(x)$  of the policy value function  $V^{h_0,\pi}$  given by (19) may not accurately distinguish the unsafe states in the avoid set  $\mathcal{A}_0 = \{x|h_0(x) > 0\}$  from the rest of the state space, due to the fact that trained neural networks can suffer from insufficient accuracy and precision [35, 17]. This is demonstrated by the orange line and its error margins in Fig.17. To alleviate this problem, we introduce discontinuity to  $h_0(x)$  for enhanced tolerance of policy value function modelling errors while maintaining the same avoid set (2), such that,

$$\mathcal{A} = \{x|h(x) > 0\} = \{x|h_0(x) > 0\}. \quad (\text{C.2})$$

where our choice of the avoidance heuristic is given by,

$$h(x) = \begin{cases} h_0(x) - 0.3, & \text{if } e_y < w_I, \\ h_0(x) + 0.2, & \text{if } w_I \leq e_y \leq w_O, \\ 2.8, & \text{if } w_O \leq e_y, \end{cases} \quad (\text{C.3})$$

where  $w_O$  is the “crash width”. If  $w_O < e_y$ , the vehicle crashes and needs to be reset. If  $w_I \leq e_y \leq w_O$ , the AutoRally collides with the soft drainage pipes but can still continue driving. As shown by the blue curve in Fig. 16,  $h(x)$  is designed to reduce the error of the avoid set of the resulting NCBF by introducing a discontinuity around zero. This is further demonstrated in Fig. 17.



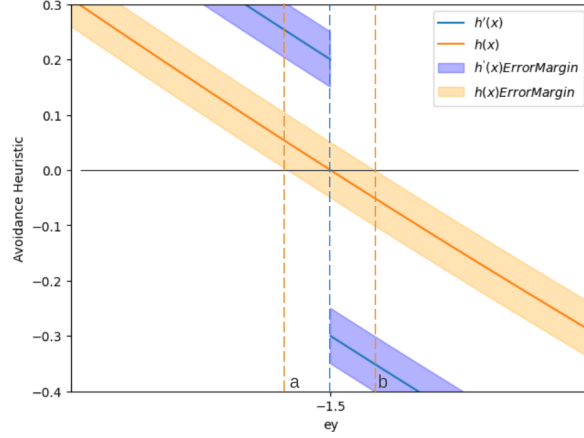


Fig. 17: Improving accuracy and precision of NCBF modeled avoid set boundary by using the modified heuristic  $h(x)$ . When using the original avoidance heuristic  $h_0(x)$  to supervise the NCBF training process, the resulting modeled avoid set boundary can be anywhere within  $e_y \in (a, b)$  due to model errors. Instead, using  $h(x)$  to supervise the training process results in a modeled avoid set with an accurate boundary shown by the dashed blue line, despite model errors.

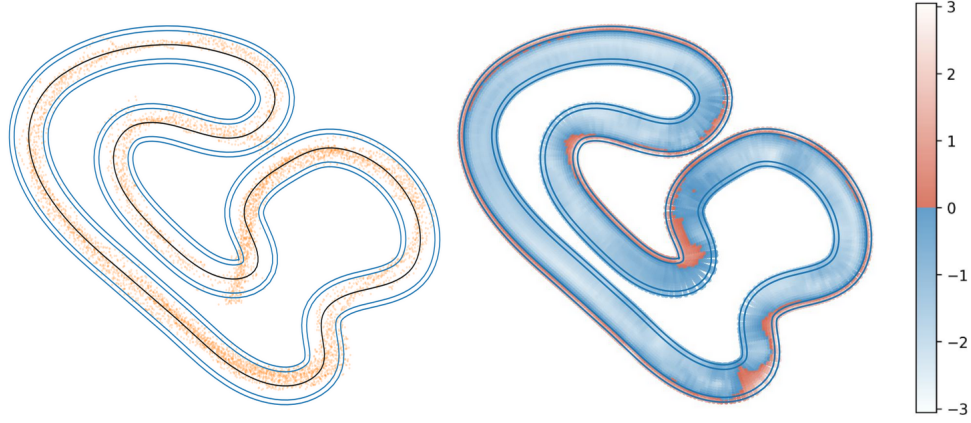


Fig. 18: NCBF Training. The left figure shows collected system states (yellow dots) in AutoRally simulations. The right figure visualizes the resulting Neural CBF. The NCBF  $B(x)$  takes the 8-dimensional state  $x$  as input. The red color indicates an unsafe region where  $B(x) > 0$ , while the blue color indicates a safe region where  $B(x) \leq 0$ .

The orange dots in the left plot in Fig. 18 represent the vehicle states collected using the training data. The dots outside of the track show that the autonomous car went off the track at several turns. These data are used to train a neural DCBF, leading to a visualization as demonstrated by the plot on the right in Fig. 18, in which the NCBF  $B(x)$  is mapped onto the 2D Cartesian plane by utilizing mean values from the training dataset for the remaining state dimensions. The red color shows the area where  $B(x) > 0$ , and as a result, the neural CBF slows down the Autorally when it approaches sharp turns to ensure safety.

From Fig. 19, we see that the learned NCBF  $B(x)$  has a smaller safe set (negative level set) than the avoid set heuristic  $h(x)$ , as it turns positive earlier than  $h(x)$  when the vehicle approaches the track boundary. When used as a safety filter,  $B(x)$  can slow down the autonomous vehicle much earlier, making its negative level set control-invariant.

## APPENDIX D

### PROOFS FOR RESAMPLING-BASED ROLLOUTS (RBR) SUPERVISED BY A CBF

#### D1 Proof of Theorem 1

*Proof:* Condition (14b) implies that  $B(x_k) \leq (\text{Id} - \alpha) \circ B(x_{k-1})$ , where  $\circ$  denotes function composition and  $\text{Id}$  denotes the identity function [2]. Since  $B(x_1) \leq (\text{Id} - \alpha) \circ B(x_0)$ , it follows that,

$$B(x_k) \leq (\text{Id} - \alpha)^k \circ B(x_0). \quad (\text{D.1})$$

Since  $(\text{Id} - \alpha)$  is a class- $\kappa$  function for  $a \in (0, 1)$ , it follows from  $B(x_0) \leq 0$  that  $B(x_k) \leq 0$ . Hence, the set  $\mathcal{S}$  is forward invariant for system (1). ■

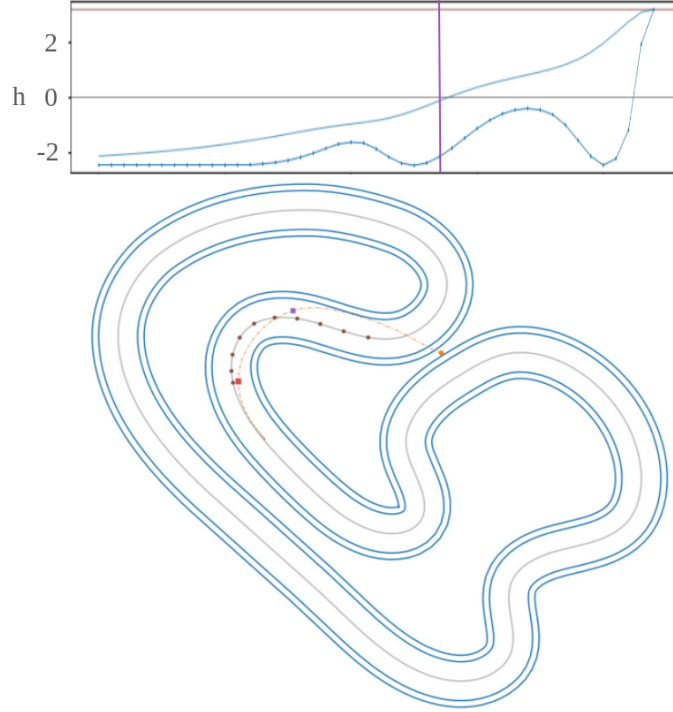


Fig. 19: A simulation example of the value change of the learned Neural CBF  $B(x)$  compared to the avoid set heuristic  $h(x)$ . The orange trajectory produced by the standard MPPI results in a collision with the track boundary. The square purple dot shows the point where  $B(x)$  shifts from negative to positive values, corresponding to the purple vertical line in the curve plot above, which shows the value change of the heuristics  $B(x)$ ,  $h(x)$  along the MPPI trajectory. The dots along the track centerline are observations collected starting at the square red dot, used to augment the training data.

## D2 Proof of Theorem 3

*Proof:* Note that

$$f(\tilde{b} \mid a, b) = \mathbb{1}_{b \in S} \delta(\tilde{b} - b) + \mathbb{1}_{b \notin S} \delta(\tilde{b} - a). \quad (D.2)$$

Let  $\tilde{f}$  denote the marginal density of  $\tilde{b}$ . Then, using the independence of  $a$  and  $b$ ,

$$\tilde{f}(\tilde{b}) = \int \left( \int f(\tilde{b} \mid a, b) f(a \mid a \in S) da \right) f(b) db \quad (D.3)$$

Simplifying the inner integral first using properties of the Dirac delta function gives

$$\int f(\tilde{b} \mid a, b) f(a \mid a \in S) da \quad (D.4)$$

$$= \mathbb{1}_{b \in S} \delta(\tilde{b} - b) + \mathbb{1}_{b \notin S} \int \delta(\tilde{b} - a) f(a \mid a \in S) da \quad (D.5)$$

$$= \mathbb{1}_{b \in S} \delta(\tilde{b} - b) + \mathbb{1}_{b \notin S} f(\tilde{b} \mid \tilde{b} \in S). \quad (D.6)$$

Hence,

$$\tilde{f}(\tilde{b}) = \int \left( \mathbb{1}_{b \in S} \delta(\tilde{b} - b) + \mathbb{1}_{b \notin S} f(\tilde{b} \mid \tilde{b} \in S) \right) f(b) db \quad (D.7)$$

$$= f(\tilde{b}) \mathbb{1}_{\tilde{b} \in S} + f(\tilde{b} \mid \tilde{b} \in S) \int \mathbb{1}_{b \notin S} f(b) db \quad (D.8)$$

$$= f(\tilde{b} \mid \tilde{b} \in S) P(b \in S) + f(\tilde{b} \mid \tilde{b} \in S) P(b \notin S) \quad (D.9)$$

$$= f(\tilde{b} \mid \tilde{b} \in S). \quad (D.10)$$

■

### D3 Proof of Corollary 1

*Proof:* Applying Theorem 3,  $\mathbf{x}$ ,  $\mathbf{a}$  and  $\tilde{\mathbf{b}}$  have the same distribution,

$$\mathbf{x} \stackrel{d}{=} \mathbf{a} \stackrel{d}{=} \tilde{\mathbf{b}}. \quad (\text{D.11})$$

Hence,

$$\mathbb{E}\left[\frac{1}{2}w(\mathbf{a}) + \frac{1}{2}w(\tilde{\mathbf{b}})\right] = \frac{1}{2}\mathbb{E}[w(\mathbf{a})] + \frac{1}{2}\mathbb{E}[w(\tilde{\mathbf{b}})], \quad (\text{D.12})$$

$$= \frac{1}{2}\mathbb{E}[w(\mathbf{x})] + \frac{1}{2}\mathbb{E}[w(\mathbf{x})], \quad (\text{D.13})$$

$$= \mathbb{E}[w(\mathbf{x})]. \quad (\text{D.14})$$

■

### D4 Proof of Theorem 4

We prove the two claims separately.

**Lemma 7.** *The variance of the Monte Carlo estimator of the optimal control law is*

$$\text{Var}[\hat{v}_k] = \frac{1}{N} \left( \frac{1}{3}2^K - \frac{1}{4} \right), \quad k = 1, \dots, K. \quad (\text{D.15})$$

*Proof:* First, note that the optimal control  $\mathbf{v}^*$  is given by

$$v_k^* = \int_{[-1,1]^K} p(\mathbf{u} \mid o = 1) u_k \, d\mathbf{u} \quad (\text{D.16})$$

$$= \int_{[-1,1]^K} \mathbb{1}_{\mathbf{u} \in [0,1]} u_k \, d\mathbf{u} \quad (\text{D.17})$$

$$= \int_{[0,1]^K} u_k \, d\mathbf{u} \quad (\text{D.18})$$

$$= \frac{1}{2}. \quad (\text{D.19})$$

Using the formula for computing the variance of an importance sampled Monte Carlo estimator [53], we then have that

$$\text{Var}[\hat{v}_k] = \frac{1}{N} \int_{[-1,1]^K} \frac{(u_k p(\mathbf{u} \mid o = 1) - v_k^* r(\mathbf{u}))^2}{r(\mathbf{u})} \, d\mathbf{u} \quad (\text{D.20})$$

$$= \frac{1}{N} \int_{[-1,1]^K} \frac{(u_k p(\mathbf{u} \mid o = 1) - \frac{1}{2}2^{-K})^2}{2^{-K}} \, d\mathbf{u} \quad (\text{D.21})$$

$$= \frac{1}{N} 2^K \left( \int_{[0,1]^K} (u_k)^2 \, d\mathbf{u} - 2^{-K} \int_{[0,1]^K} u_k \, d\mathbf{u} + \frac{1}{4} (2^{-K})^2 \int_{[-1,1]^K} 1 \, d\mathbf{u} \right) \quad (\text{D.22})$$

$$= \frac{1}{N} 2^K \left( \frac{1}{3} - \frac{1}{2} 2^{-K} + \frac{1}{4} 2^{-K} \right) \quad (\text{D.23})$$

$$= \frac{1}{N} \left( \frac{1}{3} 2^K - \frac{1}{4} \right). \quad (\text{D.24})$$

■

**Proof of second claim.** We now prove the second claim. When performing resampling-based rollouts, *only* if  $u_k^i$  lies in  $[-1, 0]$  for all  $k = 0, \dots, K-2$  does resampling not occur, and the output  $\tilde{u}_k^i \in [-1, 0]$  for  $k = 0, \dots, K-2$ . Otherwise, we have that  $\tilde{u}_k^i \in [0, 1]$  for  $k = 0, \dots, K-2$ . The last control  $u_{K-1}^i = \tilde{u}_{K-1}^i$  is never resampled. Hence, the probability that all  $N$  samples lie in  $[-1, 0]^{K-1} \times [-1, 1]$  is equal to  $2^{-N(K-1)}$ .

For all timesteps except the last, we write the joint probability density function (over all  $N$  samples) as

$$p(\tilde{u}_k^1, \dots, \tilde{u}_k^N) = \begin{cases} 2^{-N}, & \text{if } \tilde{u}_k^i \in [-1, 0] \text{ for all } i = 1, \dots, N, \\ 1 - 2^{-N}, & \text{if } \tilde{u}_k^i \in [0, 1] \text{ for all } i = 1, \dots, N, \\ 0, & \text{otherwise,} \end{cases} \quad \text{for } k = 0, \dots, K-2, \quad (\text{D.25})$$

and

$$p(\tilde{u}_{K-1}^i) = 1, \quad p(\tilde{u}_{K-1}^1, \dots, \tilde{u}_{K-1}^N) = \prod_{i=1}^N p(\tilde{u}_{K-1}^i) = 1. \quad (\text{D.26})$$

For convenience, let  $t := 2^{-N}$  such that  $P(\tilde{u}_k^{1:N} \in [0, 1]) = 1 - t$  for  $k < K - 1$ , and

$$P(\tilde{u}_{0:K-2}^{1:N} \in [0, 1]) = \prod_{k=0}^{K-2} P(\tilde{u}_k^{1:N} \in [0, 1]) = (1 - t)^{K-1}. \quad (\text{D.27})$$

Next, consider the Monte Carlo estimator  $\hat{v}_k$  using the resampled controls  $\tilde{u}_k^i$ , defined by

$$\hat{v}_k = \frac{1}{N} \sum_{i=1}^N \frac{\mathbb{1}_{\tilde{u}_{0:K-1}^i \in [0, 1]}}{\frac{1}{2}(1-t)^{K-1}} \tilde{u}_k^i. \quad (\text{D.28})$$

We will compute the expectation of  $\hat{v}_k$  using the law of total expectation. To this end, we have that

$$\mathbb{E}_{\tilde{u}_{0:K-1}^{1:N}}[\hat{v}_k] = \frac{1}{N} \mathbb{E}_{\tilde{u}_{0:K-1}^{1:N}} \left[ \sum_{i=1}^N \frac{\mathbb{1}_{\tilde{u}_{0:K-1}^i \in [0, 1]}}{\frac{1}{2}(1-t)^{K-1}} \tilde{u}_k^i \right] \quad (\text{D.29})$$

$$= \frac{1}{N} \left( P(\tilde{u}_{0:K-2}^{1:N} \geq 0) \mathbb{E}_{\tilde{u}_{0:K-1}^{1:N}} \left[ \sum_{i=1}^N \frac{\mathbb{1}_{\tilde{u}_{0:K-1}^i \in [0, 1]}}{\frac{1}{2}(1-t)^{K-1}} \tilde{u}_k^i \mid \tilde{u}_{0:K-2}^{1:N} \geq 0 \right] \right. \\ \left. + P(\tilde{u}_{0:K-2}^{1:N} < 0) \mathbb{E}_{\tilde{u}_{0:K-1}^{1:N}} \left[ \sum_{i=1}^N \frac{\mathbb{1}_{\tilde{u}_{0:K-1}^i \in [0, 1]}}{\frac{1}{2}(1-t)^{K-1}} \tilde{u}_k^i \mid \tilde{u}_{0:K-2}^{1:N} < 0 \right] \right) \quad (\text{D.30})$$

$$= \frac{1}{N} \left( (1-t)^{K-1} \mathbb{E}_{\tilde{u}_{0:K-1}^{1:N}} \left[ \sum_{i=1}^N \frac{\mathbb{1}_{\tilde{u}_{0:K-1}^i \in [0, 1]}}{\frac{1}{2}(1-t)^{K-1}} \tilde{u}_k^i \mid \tilde{u}_{0:K-2}^{1:N} \geq 0 \right] \right) \quad (\text{D.31})$$

$$= \frac{2}{N} \left( \mathbb{E}_{\tilde{u}_{0:K-1}^{1:N}} \left[ \sum_{i=1}^N \mathbb{1}_{\tilde{u}_{0:K-1}^i \in [0, 1]} \tilde{u}_k^i \mid \tilde{u}_{0:K-2}^{1:N} \geq 0 \right] \right). \quad (\text{D.32})$$

We now split into two cases. When  $k = K - 1$ ,

$$\mathbb{E}_{\tilde{u}_{0:K-1}^{1:N}}[\hat{v}_{K-1}] = \frac{2}{N} \left( \mathbb{E}_{\tilde{u}_{K-1}^{1:N}} \left[ \mathbb{E}_{\tilde{u}_{0:K-2}^{1:N}} \left[ \sum_{i=1}^N \mathbb{1}_{\tilde{u}_{0:K-1}^i \in [0, 1]} \tilde{u}_{K-1}^i \mid \tilde{u}_{0:K-2}^{1:N} \geq 0 \right] \right] \right) \quad (\text{D.33})$$

$$= \frac{2}{N} \left( \mathbb{E}_{\tilde{u}_{K-1}^{1:N}} \left[ \sum_{i=1}^N \mathbb{1}_{\tilde{u}_{K-1}^i \in [0, 1]} \tilde{u}_{K-1}^i \right] \right) \quad (\text{D.34})$$

$$= \frac{2}{N} \left( \sum_{i=1}^N \mathbb{E}_{\tilde{u}_{K-1}^i} \left[ \mathbb{1}_{\tilde{u}_{K-1}^i \in [0, 1]} \tilde{u}_{K-1}^i \right] \right) \quad (\text{D.35})$$

$$= \frac{2}{N} \left( \sum_{i=1}^N \frac{1}{4} \right) \quad (\text{D.36})$$

$$= \frac{1}{2}. \quad (\text{D.37})$$

Similarly, when  $k < K - 1$ ,

$$\mathbb{E}_{\tilde{u}_{0:K-1}^{1:N}}[\hat{v}_k] = \frac{2}{N} \left( \mathbb{E}_{\tilde{u}_{K-1}^{1:N}} \left[ \mathbb{E}_{\tilde{u}_{0:K-2}^{1:N}} \left[ \sum_{i=1}^N \mathbb{1}_{\tilde{u}_{K-1}^i \in [0, 1]} \mathbb{1}_{\tilde{u}_{0:K-2}^i \in [0, 1]} \tilde{u}_k^i \mid \tilde{u}_{0:K-2}^{1:N} \geq 0 \right] \right] \right) \quad (\text{D.38})$$

$$= \frac{1}{N} \left( \mathbb{E}_{\tilde{u}_{0:K-2}^{1:N}} \left[ \sum_{i=1}^N \tilde{u}_k^i \mid \tilde{u}_{0:K-2}^{1:N} \geq 0 \right] \right) \quad (\text{D.39})$$

$$= \frac{1}{2}. \quad (\text{D.40})$$

Thus,  $\hat{v}_k$  is an unbiased estimator.

Next, we look at the variance. Summarizing the above computation, we have that

$$\mathbb{E}_{\tilde{u}_{0:K-1}^{1:N}}[\hat{v}_k \mid \tilde{u}_{0:K-2}^{1:N} \geq 0] = \frac{1}{2(1-t)^{K-1}}, \quad \mathbb{E}_{\tilde{u}_{0:K-1}^{1:N}}[\hat{v}_k \mid \tilde{u}_{0:K-2}^{1:N} < 0] = 0. \quad (\text{D.41})$$

Let  $c := \mathbb{E}_{\tilde{u}_{0:K-1}^{1:N}}[\hat{v}_k \mid \tilde{u}_{0:K-2}^{1:N} \geq 0] = \frac{1}{2(1-t)^{K-1}}$  for convenience. Computing the conditional variances, for  $\tilde{u}_{0:K-2}^{1:N} < 0$ , we have that

$$\text{Var}_{\tilde{u}_{0:K-1}^{1:N}}[\hat{v}_k \mid \tilde{u}_{0:K-2}^{1:N} < 0] = 0. \quad (\text{D.42})$$

For  $\tilde{u}_{0:K-2}^{1:N} \geq 0$  we have,

$$\text{Var}_{\tilde{u}_{0:K-1}^{1:N}}[\hat{v}_k \mid \tilde{u}_{0:K-2}^{1:N} \geq 0] = \mathbb{E}_{\tilde{u}_{0:K-1}^{1:N}}[\hat{v}_k^2 \mid \tilde{u}_{0:K-2}^{1:N} \geq 0] - c^2 \quad (\text{D.43})$$

$$= \mathbb{E}_{\tilde{u}_{0:K-1}^{1:N}} \left[ \left( \frac{1}{N} \sum_{i=1}^N \frac{\mathbb{1}_{\tilde{u}_{0:K-1}^i \in [0,1]}}{\frac{1}{2}(1-t)^{K-1}} \tilde{u}_k^i \right)^2 \mid \tilde{u}_{0:K-2}^{1:N} \geq 0 \right] - c^2 \quad (\text{D.44})$$

$$= \left( \frac{1}{N} \frac{1}{\frac{1}{2}(1-t)^{K-1}} \right)^2 \underbrace{\mathbb{E}_{\tilde{u}_{0:K-1}^{1:N}} \left[ \left( \sum_{i=1}^N \mathbb{1}_{\tilde{u}_{0:K-1}^i \in [0,1]} \tilde{u}_k^i \right)^2 \mid \tilde{u}_{0:K-2}^{1:N} \geq 0 \right]}_{:= (\star)} - c^2 \quad (\text{D.45})$$

Expanding  $(\star)$ , yields

$$(\star) = \mathbb{E}_{\tilde{u}_{0:K-1}^{1:N}} \left[ \sum_{i=1}^N \mathbb{1}_{\tilde{u}_{0:K-1}^i \in [0,1]} (\tilde{u}_k^i)^2 + \sum_{i=1}^N \sum_{j=1, j \neq i}^N \mathbb{1}_{\tilde{u}_{0:K-1}^i \in [0,1]} \mathbb{1}_{\tilde{u}_{0:K-1}^j \in [0,1]} \tilde{u}_k^i \tilde{u}_k^j \mid \tilde{u}_{0:K-2}^{1:N} \geq 0 \right]. \quad (\text{D.46})$$

We now split into two cases.

**Case 1:**  $k = K - 1$ . Looking at the first term of  $(\star)$  in (D.46),

$$\mathbb{E}_{\tilde{u}_{0:K-1}^{1:N}} \left[ \sum_{i=1}^N \mathbb{1}_{\tilde{u}_{0:K-1}^i \in [0,1]} (\tilde{u}_{K-1}^i)^2 \mid \tilde{u}_{0:K-2}^{1:N} \geq 0 \right] = \mathbb{E}_{\tilde{u}_{0:K-1}^{1:N}} \left[ \sum_{i=1}^N \mathbb{E}_{\tilde{u}_{K-1}^i} [\mathbb{1}_{\tilde{u}_{K-1}^i \in [0,1]} (\tilde{u}_{K-1}^i)^2] \mid \tilde{u}_{0:K-2}^{1:N} \geq 0 \right] \quad (\text{D.47})$$

$$= \sum_{i=1}^N \mathbb{E}_{\tilde{u}_{K-1}^i} [\mathbb{1}_{\tilde{u}_{K-1}^i \in [0,1]} (\tilde{u}_{K-1}^i)^2] \quad (\text{D.48})$$

$$= \frac{N}{6}. \quad (\text{D.49})$$

For the second term in (D.46), we have

$$\mathbb{E}_{\tilde{u}_{0:K-1}^{1:N}} \left[ \sum_{i=1}^N \sum_{j=1, j \neq i}^N \mathbb{1}_{\tilde{u}_{0:K-1}^i \in [0,1]} \mathbb{1}_{\tilde{u}_{0:K-1}^j \in [0,1]} \tilde{u}_{K-1}^i \tilde{u}_{K-1}^j \mid \tilde{u}_{0:K-2}^{1:N} \geq 0 \right] \quad (\text{D.50})$$

$$= \sum_{i=1}^N \sum_{j=1, j \neq i}^N \mathbb{E}_{\tilde{u}_{K-1}^{i,j}} \left[ \mathbb{1}_{\tilde{u}_{K-1}^i \in [0,1]} \mathbb{1}_{\tilde{u}_{K-1}^j \in [0,1]} \tilde{u}_{K-1}^i \tilde{u}_{K-1}^j \right] \quad (\text{D.51})$$

$$= \sum_{i=1}^N \sum_{j=1, j \neq i}^N \mathbb{E}_{\tilde{u}_{K-1}^{i,j}} \left[ \mathbb{1}_{\tilde{u}_{K-1}^i \in [0,1]} \tilde{u}_{K-1}^i \right]^2 \quad (\text{D.52})$$

$$= \sum_{i=1}^N \sum_{j=1, j \neq i}^N \frac{1}{4} \quad (\text{D.53})$$

$$= \frac{N(N-1)}{4}. \quad (\text{D.54})$$



Substituting the two terms (D.49) and (D.54) into (D.45), yields

$$\text{Var}_{\tilde{u}_{0:K-1}^{1:N}}[\hat{v}_{K-1} \mid \tilde{u}_{0:K-2}^{1:N} \geq 0] = \left( \frac{1}{N} \frac{1}{\frac{1}{2}(1-t)^{K-1}} \right)^2 \left( \frac{N}{6} + \frac{N(N-1)}{4} \right) - c^2 \quad (\text{D.55})$$

$$\leq \left( \frac{1}{N} \frac{1}{\frac{1}{2}(1-t)^{K-1}} \right)^2 \frac{N^2}{4} - c^2 \quad (\text{D.56})$$

$$= \left( \frac{1}{(1-t)^2} \right)^{K-1} - \frac{1}{4} \left( \frac{1}{(1-t)^2} \right)^{K-1}, \quad (\text{D.57})$$

$$= \frac{3}{4} \left( \frac{1}{(1-t)^2} \right)^{K-1}. \quad (\text{D.58})$$

**Case 2:**  $k < K - 1$ . Looking at the first term of  $(\star)$  in (D.46),

$$\mathbb{E}_{\tilde{u}_{0:K-1}^{1:N}} \left[ \sum_{i=1}^N \mathbb{1}_{\tilde{u}_{0:K-1}^i \in [0,1]} (\tilde{u}_k^i)^2 \mid \tilde{u}_{0:K-2}^{1:N} \geq 0 \right] = \mathbb{E}_{\tilde{u}_{0:K-1}^{1:N}} \left[ \sum_{i=1}^N \mathbb{E}_{\tilde{u}_{K-1}^i} [\mathbb{1}_{\tilde{u}_{K-1}^i \in [0,1]}] (\tilde{u}_k^i)^2 \mid \tilde{u}_{0:K-2}^{1:N} \geq 0 \right] \quad (\text{D.59})$$

$$= \frac{1}{2} \mathbb{E}_{\tilde{u}_{0:K-1}^{1:N}} \left[ \sum_{i=1}^N (\tilde{u}_k^i)^2 \mid \tilde{u}_{0:K-2}^{1:N} \geq 0 \right]. \quad (\text{D.60})$$

From Theorem 3, the marginal of any resampled controls will all be uniform over  $[0, 1]$ . Hence,

$$\frac{1}{2} \mathbb{E}_{\tilde{u}_{0:K-1}^{1:N}} \left[ \sum_{i=1}^N (\tilde{u}_k^i)^2 \mid \tilde{u}_{0:K-2}^{1:N} \geq 0 \right] = \frac{1}{2} \sum_{i=1}^N \mathbb{E}_{\tilde{u}_k^i} [(\tilde{u}_k^i)^2 \mid \tilde{u}_k^i \geq 0] = \frac{N}{6}. \quad (\text{D.61})$$

For the second term in (D.46),

$$\mathbb{E}_{\tilde{u}_{0:K-1}^{1:N}} \left[ \sum_{i=1}^N \sum_{j=1, j \neq i}^N \mathbb{1}_{\tilde{u}_{0:K-1}^i \in [0,1]} \mathbb{1}_{\tilde{u}_{0:K-1}^j \in [0,1]} \tilde{u}_k^i \tilde{u}_k^j \mid \tilde{u}_{0:K-2}^{1:N} \geq 0 \right] \quad (\text{D.62})$$

$$= \mathbb{E}_{\tilde{u}_{0:K-2}^{1:N}} \left[ \sum_{i=1}^N \sum_{j=1, j \neq i}^N \mathbb{E}_{\tilde{u}_{K-1}^{1:N}} [\mathbb{1}_{\tilde{u}_{K-1}^i \in [0,1]} \mathbb{1}_{\tilde{u}_{K-1}^j \in [0,1]}] \tilde{u}_k^i \tilde{u}_k^j \mid \tilde{u}_{0:K-2}^{1:N} \geq 0 \right] \quad (\text{D.63})$$

$$= \frac{1}{4} \mathbb{E}_{\tilde{u}_{0:K-2}^{1:N}} \left[ \sum_{i=1}^N \sum_{j=1, j \neq i}^N \tilde{u}_k^i \tilde{u}_k^j \mid \tilde{u}_{0:K-2}^{1:N} \geq 0 \right]. \quad (\text{D.64})$$

To make this computation easier, let  $\Xi_{ij} = 1$  denote the event that  $\tilde{u}_k^i$  and  $\tilde{u}_k^j$  for  $i \neq j$  were resampled from the same control, i.e.,  $\tilde{u}_k^i = \tilde{u}_k^j$ , and let  $\beta = \mathbb{P}(\Xi_{ij} = 1)$ . When  $\Xi_{ij} = 1$ ,  $\tilde{u}_k^i$  and  $\tilde{u}_k^j$  are the same random variable, and therefore it follows that

$$\mathbb{E}_{\tilde{u}_{0:K-2}^{1:N}} \left[ \tilde{u}_k^i \tilde{u}_k^j \mid \Xi_{ij} = 1 \mid \tilde{u}_{0:K-2}^{1:N} \geq 0, \Xi_{ij} = 1 \right] = \mathbb{E}_{\tilde{u}_k^i} [(\tilde{u}_k^i)^2 \mid \tilde{u}_k^i \geq 0] = \frac{1}{3}. \quad (\text{D.65})$$

Otherwise, when  $\Xi_{ij} = 0$ ,  $\tilde{u}_k^i$  and  $\tilde{u}_k^j$  are independent. Hence,

$$\mathbb{E}_{\tilde{u}_{0:K-2}^{1:N}} \left[ \tilde{u}_k^i \tilde{u}_k^j \mid \Xi_{ij} = 1 \mid \tilde{u}_{0:K-2}^{1:N} \geq 0, \Xi_{ij} = 0 \right] = \mathbb{E}_{\tilde{u}_k^i} [(\tilde{u}_k^i)^2 \mid \tilde{u}_k^i \geq 0] = \frac{1}{4}. \quad (\text{D.66})$$

Hence, the expectation becomes

$$(\text{D.64}) = \frac{1}{4} \sum_{i=1}^N \sum_{j=1, j \neq i}^N \beta \frac{1}{3} + (1 - \beta) \frac{1}{4} \quad (\text{D.67})$$

$$= \frac{1}{4} N(N-1) \left( \beta \frac{1}{3} + (1 - \beta) \frac{1}{4} \right) \quad (\text{D.68})$$

$$\leq \frac{1}{12} N(N-1). \quad (\text{D.69})$$

Combining the two terms (D.61) and (D.69), we thus have that

$$\text{Var}_{\tilde{u}_{0:K-1}^{1:N}}[\hat{v}_k \mid \tilde{u}_{0:K-2}^{1:N} \geq 0] \leq \left( \frac{1}{N} \frac{1}{\frac{1}{2}(1-t)^{K-1}} \right)^2 \left( \frac{N}{6} + \frac{N(N-1)}{12} \right) - c^2 \quad (\text{D.70})$$

$$\leq \left( \frac{1}{N} \frac{1}{\frac{1}{2}(1-t)^{K-1}} \right)^2 \frac{N^2}{6} - c^2 \quad (\text{D.71})$$

$$= \frac{2}{3} \left( \frac{1}{(1-t)^2} \right)^{K-1} - \frac{1}{4} \left( \frac{1}{(1-t)^2} \right)^{K-1} \quad (\text{D.72})$$

$$= \frac{5}{12} \left( \frac{1}{(1-t)^2} \right)^{K-1}. \quad (\text{D.73})$$

Hence, taking both cases into account, we have that

$$\text{Var}_{\tilde{u}_{0:K-1}^{1:N}}[\hat{v}_k \mid \tilde{u}_{0:K-2}^{1:N} \geq 0] \leq \frac{3}{4} \left( \frac{1}{(1-t)^2} \right)^{K-1}. \quad (\text{D.74})$$

Finally, using the law of total variance, we have that

$$\text{Var}[\hat{v}_k] = \underbrace{\mathbb{E}_{\text{sign of } \tilde{u}_{0:K-2}^{1:N}}[\text{Var}[\hat{v}_k \mid \tilde{u}_{0:K-2}^{1:N}]]}_{\textcircled{1}} + \underbrace{\text{Var}_{\text{sign of } \tilde{u}_{0:K-2}^{1:N}}[\mathbb{E}[\hat{v}_k \mid \tilde{u}_{0:K-2}^{1:N}]]}_{\textcircled{2}}. \quad (\text{D.75})$$

The first term gives

$$\textcircled{1} = \mathbb{P}(\tilde{u}_{0:K-2}^{1:N} \geq 0) \text{Var}[\hat{v}_k \mid \tilde{u}_{0:K-2}^{1:N} \geq 0] + \mathbb{P}(\tilde{u}_{0:K-2}^{1:N} < 0) \text{Var}[\hat{v}_k \mid \tilde{u}_{0:K-2}^{1:N} < 0] \quad (\text{D.76})$$

$$\leq (1-t)^{K-1} \frac{3}{4} \left( \frac{1}{(1-t)^2} \right)^{K-1} + 0 \quad (\text{D.77})$$

$$= \frac{3}{4} \left( \frac{1}{1-t} \right)^{K-1}. \quad (\text{D.78})$$

The second term gives

$$\textcircled{2} = \mathbb{P}(\tilde{u}_{0:K-2}^{1:N} \geq 0)(\mathbb{E}[\hat{v}_k \mid \tilde{u}_{0:K-2}^{1:N} \geq 0] - \mathbb{E}[\hat{v}_k])^2 + \mathbb{P}(\tilde{u}_{0:K-2}^{1:N} < 0)(\mathbb{E}[\hat{v}_k \mid \tilde{u}_{0:K-2}^{1:N} < 0] - \mathbb{E}[\hat{v}_k])^2 \quad (\text{D.79})$$

$$= (1-t)^{K-1} \left( \frac{1}{2} \frac{1}{(1-t)^{K-1}} - \frac{1}{2} \right)^2 + (1 - (1-t)^{K-1}) \frac{1}{4} \quad (\text{D.80})$$

$$= \frac{1}{2}(1-t)^{K-1} - 1 + \frac{1}{2} \frac{1}{(1-t)^{K-1}} + \frac{1}{4} - \frac{1}{4}(1-t)^{K-1} \quad (\text{D.81})$$

$$= \frac{1}{4}(1-t)^{K-1} - \frac{3}{4} + \frac{1}{2} \frac{1}{(1-t)^{K-1}}. \quad (\text{D.82})$$

Hence, combining both terms gives us

$$\text{Var}[\hat{v}_k] \leq \frac{3}{4} \left( \frac{1}{1-t} \right)^{K-1} + \frac{1}{4}(1-t)^{K-1} - \frac{3}{4} + \frac{1}{2} \left( \frac{1}{1-t} \right)^{K-1} \quad (\text{D.83})$$

$$= \frac{5}{4} \left( \frac{1}{1-t} \right)^{K-1} + \frac{1}{4}(1-t)^{K-1} - \frac{3}{4} \quad (\text{D.84})$$

$$= O \left( \left( \frac{1}{1-2^{-N}} \right)^K + (1-2^{-N})^K \right). \quad (\text{D.85})$$

#### D5 Proof of Theorem 5

*Proof:* We will prove the equivalent statement that

$$\|\tilde{\mathbf{w}}\|_2^2 - \|\tilde{\mathbf{w}}'\|_2^2 \geq 0. \quad (\text{D.86})$$

To begin, we expand  $\tilde{\mathbf{w}}'$  to obtain

$$\|\tilde{\mathbf{w}}'\|_2^2 = \frac{\|\mathbf{w}\|_2^2 + \|\mathbf{c}\|_2^2}{\|\mathbf{w}'\|_1^2}. \quad (\text{D.87})$$

Then,

$$\|\tilde{\mathbf{w}}\|_2^2 - \|\tilde{\mathbf{w}}'\|_2^2 = \frac{\|\mathbf{w}\|_2^2}{\|\mathbf{w}\|_1^2} + \frac{-\|\mathbf{w}\|_2^2 - \|\mathbf{c}\|_2^2}{\|\mathbf{w}'\|_1^2} \quad (\text{D.88})$$

$$= \frac{1}{\|\mathbf{w}'\|_1^2} \left( \|\mathbf{w}'\|_1^2 \frac{\|\mathbf{w}\|_2^2}{\|\mathbf{w}\|_1^2} - \|\mathbf{w}\|_2^2 - \|\mathbf{c}\|_2^2 \right). \quad (\text{D.89})$$

Simplifying the first term, yields,

$$\|\mathbf{w}'\|_1^2 \frac{\|\mathbf{w}\|_2^2}{\|\mathbf{w}\|_1^2} = \left( \|\mathbf{w}\|_1^2 + 2\|\mathbf{w}\|_1\|\mathbf{c}\|_1 + \|\mathbf{c}\|_1^2 \right) \frac{\|\mathbf{w}\|_2^2}{\|\mathbf{w}\|_1^2} \quad (\text{D.90})$$

$$= \|\mathbf{w}\|_2^2 + 2\frac{\|\mathbf{w}\|_2^2\|\mathbf{c}\|_1}{\|\mathbf{w}\|_1} + \frac{\|\mathbf{w}\|_2^2\|\mathbf{c}\|_1^2}{\|\mathbf{w}\|_1^2}. \quad (\text{D.91})$$

Substituting (D.91) back into the parenthesis (D.89) gives us that

$$\|\mathbf{w}'\|_1^2 \frac{\|\mathbf{w}\|_2^2}{\|\mathbf{w}\|_1^2} - \|\mathbf{w}\|_2^2 - \|\mathbf{c}\|_2^2 = 2\frac{\|\mathbf{w}\|_2^2\|\mathbf{c}\|_1}{\|\mathbf{w}\|_1} + \frac{\|\mathbf{w}\|_2^2\|\mathbf{c}\|_1^2}{\|\mathbf{w}\|_1^2} - \|\mathbf{c}\|_2^2 \quad (\text{D.92})$$

$$\geq 2\frac{\|\mathbf{w}\|_2^2\|\mathbf{c}\|_1}{\|\mathbf{w}\|_1} + \frac{\|\mathbf{w}\|_2^2\|\mathbf{c}\|_1^2}{\|\mathbf{w}\|_1^2} - \|\mathbf{c}\|_1^2 \quad (\text{D.93})$$

$$= \|\mathbf{c}\|_1 \left( 2\frac{\|\mathbf{w}\|_2^2}{\|\mathbf{w}\|_1} + \frac{\|\mathbf{w}\|_2^2\|\mathbf{c}\|_1}{\|\mathbf{w}\|_1^2} - \|\mathbf{c}\|_1 \right). \quad (\text{D.94})$$

Simplifying the parenthesis in (D.94), we obtain

$$2\frac{\|\mathbf{w}\|_2^2}{\|\mathbf{w}\|_1} + \frac{\|\mathbf{w}\|_2^2\|\mathbf{c}\|_1}{\|\mathbf{w}\|_1^2} - \|\mathbf{c}\|_1 = 2\frac{\|\mathbf{w}\|_2^2}{\|\mathbf{w}\|_1} + \|\mathbf{c}\|_1 \left( \frac{\|\mathbf{w}\|_2^2}{\|\mathbf{w}\|_1^2} - 1 \right) \quad (\text{D.95})$$

$$\geq 2\frac{\|\mathbf{w}\|_2^2}{\|\mathbf{w}\|_1} + \|\mathbf{c}\|_1 \left( \frac{1}{N} - 1 \right). \quad (\text{D.96})$$

Using our assumption that  $\mathbf{c}$  is not “drastically larger” than  $\mathbf{w}$  in (37) then gives us the desired result.  $\blacksquare$

## APPENDIX E DETAILS ON VIMPC COST DESIGNS

In both simulations and experiments, MPPI variants using DCBF to ensure safety utilized a state-dependent cost function:

$$q(x_k^m) = (x_k^m - x_g)^\top Q (x_k^m - x_g), \quad (\text{E.1})$$

whereas the standard MPPI employed the cost,

$$q(x_k^m) = (x_k^m - x_g)^\top Q (x_k^m - x_g) + \mathbf{1}(x_k^m), \quad (\text{E.2})$$

where  $Q = \text{diag}(q_{v_x}, q_{v_y}, q_{\dot{\psi}}, q_{\omega_F}, q_{\omega_R}, q_{e_\psi}, q_{e_y}, q_s)$  represents the cost weights, and  $x_g = \text{diag}(v_g, 0, \dots, 0)$  specifies the target velocity. The collision cost function is defined as:

$$\mathbf{1}(x_k^m) := \begin{cases} 0, & \text{if } x_k^m \text{ is within the track,} \\ C_{\text{obs}}, & \text{otherwise.} \end{cases} \quad (\text{E.3})$$

Unlike standard MPPI, Shield-MPPI and NS-MPPI incorporate a DCBF constraint violation penalty but do not include an explicit collision cost.

APPENDIX F  
CONSTRAINT SATISFACTION ON THE VARIATIONAL DISTRIBUTION

One way of guaranteeing that the variational distribution  $q(\mathbf{u})$  satisfies the constraints  $x_k \in \bar{\mathcal{X}} := \mathcal{X} \setminus \mathcal{A}$  is by making the strong (but unrealistic) assumption that the set of controls that satisfy the constraints is a convex set. Specifically, for generic state constraints  $x \in \bar{\mathcal{X}}$ , define the set of safe control trajectories  $\mathcal{U}_{\text{safe}}^K \subseteq \mathcal{U}^K$  as

$$\mathcal{U}_{\text{safe}}^K(x_0) = \{\mathbf{u} \in \mathcal{U}^K \mid x_k \in \bar{\mathcal{X}}, \ k = 0, 1, \dots, K\}. \quad (\text{F.1})$$

We then have the following lemma.

**Lemma 8.** *Assume  $\mathcal{U}_{\text{safe}}^K(x_0)$  is convex for all states  $x_0 \in \bar{\mathcal{X}}$ , and suppose that  $p(\mathbf{u} \mid o = 1)$  has zero density outside of  $\mathcal{U}_{\text{safe}}^K(x_0)$ . Then, the mean  $\mathbf{v}^*$  (8) of the variational distribution  $q_{\mathbf{v}}$  (and the state trajectory resulting from following  $\mathbf{v}^*$ ) will also satisfy the constraints, i.e.,*

$$\mathbf{v}^* \in \mathcal{U}_{\text{safe}}^K(x_0). \quad (\text{F.2})$$

*Proof:* By definition (6),  $\mathbf{v}^*$  is the mean of the optimal distribution  $p(\mathbf{u} \mid o = 1)$ . Since  $p(\mathbf{u} \mid o = 1)$  has zero density outside  $\mathcal{U}_{\text{safe}}^K(x_0)$ , its support is contained within  $\mathcal{U}_{\text{safe}}^K(x_0)$ . Since  $\mathcal{U}_{\text{safe}}^K(x_0)$  is convex, the integral of  $\mathbf{u}$  over  $p(\mathbf{u} \mid o = 1)$  will also be contained within  $\mathcal{U}_{\text{safe}}^K(x_0)$ . Hence,  $\mathbf{v}^* \in \mathcal{U}_{\text{safe}}^K(x_0)$ . ■