

Improving the Resilience of Quadrotors in Underground Environments by Combining Learning-based and Safety Controllers

Isaac R. Ward, Mark Paral, Kristopher Riordan, and Mykel J. Kochenderfer

Abstract—Autonomously controlling quadrotors in large-scale subterranean environments is applicable to many areas such as environmental surveying, mining operations, and search and rescue. Learning-based controllers represent an appealing approach to autonomy, but are known to not generalize well to ‘out-of-distribution’ environments not encountered during training. In this work, we train a normalizing flow-based prior over the environment, which provides a measure of how far out-of-distribution the quadrotor is at any given time. We use this measure as a runtime monitor, allowing us to switch between a learning-based controller and a safe controller when we are sufficiently out-of-distribution. Our methods are benchmarked on a point-to-point navigation task in a simulated 3D cave environment based on real-world point cloud data from the DARPA Subterranean Challenge Final Event Dataset. Our experimental results show that our combined controller simultaneously possesses the liveness of the learning-based controller (completing the task quickly) and the safety of the safety controller (avoiding collision).

I. Introduction

Reliable methods for autonomously navigating large-scale, subterranean environments with quadrotors are becoming increasingly relevant, with applications in search and rescue [1], [2], mining operations [3], [4], terrestrial cave exploration [5], [6], and space exploration [7], [8], [9].

Learning-based controllers represent an attractive approach to autonomy; they can be learned from data, exhibit extreme maneuverability, and account for non-linear dynamics [10]. However, these methods are known to generalize poorly to scenarios not encountered during training [11]. Note that we refer to the environment that an algorithm is tuned or trained on as in-distribution (InD); otherwise, it is out-of-distribution (OOD) [12].

Another branch of potential solutions to autonomy are well-established control theoretic approaches, which can guarantee safety [13]. Such optimal control algorithms allow designers to mathematically optimize desirable performance metrics in the presence of constraints (e.g. collision avoidance and dynamic feasibility).

There is an inherent tradeoff between safety (something bad will not happen, e.g. avoiding colliding into a wall) and liveness (something good must happen, e.g. reaching the goal state quickly) [14]. In this work, we balance both considerations by using learning-based

controllers when InD, and traditional safety controllers when OOD, relying on a runtime monitor to classify InD from OOD.

The contributions of this work are threefold:

- 1) We train a normalizing flow-based optimal control policy (FlowMPPI [15]) within a Bayesian model-based reinforcement learning paradigm [16]. To our knowledge, this is the largest 3D environment on which FlowMPPI has been trained on to date (3D extents of $41 \times 62 \times 11$ meters, and an internal volume of 11492 cubic metres).
- 2) We design a safety controller that solves for dynamically feasible, obstacle-avoiding trajectories [17], [18].
- 3) We switch between the two methods at test time using an OOD runtime monitor that considers the environment, current state, and goal state. We experimentally demonstrate that our combined controller simultaneously possesses the liveness of the learning-based controller (completing the task quickly) and the safety of the safety controller (avoiding collisions).

II. Related Literature

The learning-based controller in this work is based on FlowMPPI [15]. This method uses a sampling-based Model Predictive Control (MPC) [19] framework called Model Predictive Path Integral Control (MPPI) [20]. MPPI operates by sampling optimal control sequences from a Gaussian prior, simulating them with respect to some known dynamics, assigning a reward to each outcome, and computing an optimal control sequence based on a reward-weighted average of the samples.

The key innovation behind FlowMPPI is the use of a conditional normalizing flow [21] in place of the Gaussian prior when representing the optimal control distribution. Normalizing flows can represent arbitrarily complex distributions, yet can be learned from data. The normalizing flow in FlowMPPI is trained in a Bayesian model-based reinforcement learning paradigm [16], and represents the optimal control distribution. Numerous optimal control samples are drawn from the normalizing flow and a loss is computed based on 1) how goal-directed each sample is, and 2) the breadth of the samples (i.e. exploitation balanced with exploration).

The safety controller used in this work is an Augmented-Lagrangian Iterative Linear Quadratic Regulator [18] (AL-iLQR) tracking the trajectory output

I. R. Ward, M. Paral, K. Riordan, and M. J. Kochenderfer are with the Stanford Intelligent Systems Laboratory in the Department of Aeronautics and Astronautics at Stanford University, Stanford, CA 94305, USA. I. R. Ward is the corresponding author (irward@stanford.edu).

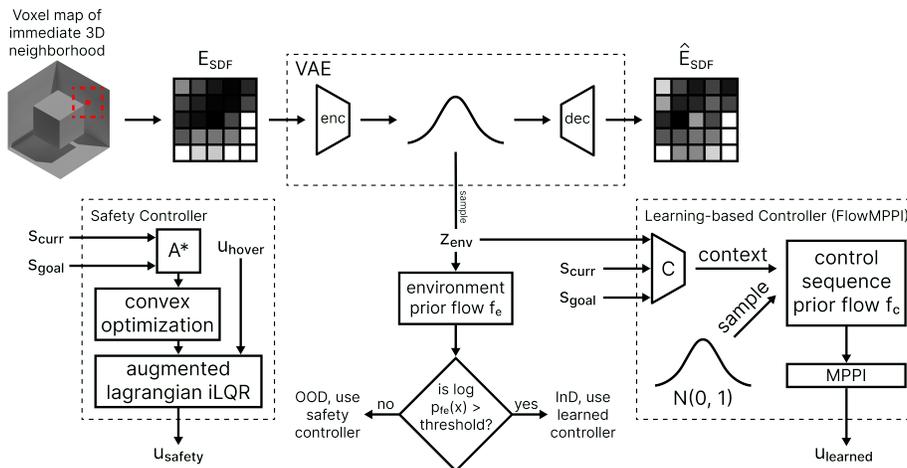


Fig. 1: A block diagram explanation of the method introduced in this work. A variational autoencoder is trained to encode signed distance fields of the environment in an online fashion. A prior over the environment encodings is trained in-tandem and probabilistically estimates if any new environmental encoding is from the prior distribution (i.e. ‘is the current environment similar to what was encountered during training?’). This probability determines which controller is more appropriate to use. The learning-based controller reaches the goal quickly, but is sensitive to out-of-distribution inputs. The safety controller reaches the goal slowly, but is robust to out-of-distribution inputs. By switching between the two, our combined controller achieves superior liveness (reaching the goal quickly) and safety (reaching the goal without collision) properties than either of the constituent methods do on their own.

from a sequential convex programming (SCP) algorithm. Prior works have demonstrated the value of iLQR-based feedback controllers in minimizing tracking errors [22], [23], and Malyuta et al. [17] presents a comprehensive guide to generating cost-minimizing, multi-objective, constrained, dynamically feasible, collision avoiding trajectories.

This work also takes advantage of the rich field of OOD detection methods [12], [24]. Many methods for OOD detection exist, including statistical methods [25], comparing latent space representations [26], and estimating the likelihood of encodings [27], the latter of which is the focus of this work.

We use these methods to create a combined controller than draws on learning-based methods where appropriate, and falls back to conservative SCP methods when needed. Prior works have investigated the construction of meta-controllers that combine 1) non-pedigreed nominal controllers, and 2) recovery controllers, and demonstrated that although safe, recovery controllers may be disruptive to operational objectives [28].

III. Methods

Our task is simulated point-to-point navigation within enclosed 3D cave environments. We consider two simple handcrafted cave environments and two complex cave environments from DARPA’s Subterranean Challenge Final Event dataset, as illustrated in Figure 2. Start and goal states are selected to be > 25 meters away from each other, and the quadrotor is initialized at the start position with zero velocity, rotation, and angular velocity. The task is considered complete if the quadrotor

stays within one half of its radius to the goal position for at least 3 seconds. The dynamics that govern the system are defined in Equation 15.

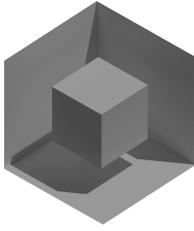
A. Model Predictive Path Integral Control (MPPI)

We implement MPPI (Algorithm 1) as a baseline, using a multivariate normal distribution as the optimal control distribution as suggested in Williams et al. [20]. This distribution is updated at each timestep to be centered about the previous optimal control u^* , with a covariance matrix that has zero off-diagonal elements. Diagonal elements represent variances, and are computed to be one-quarter of the range of the allowed control inputs. Increasing or decreasing these variances has the effect of widening or tightening the control distribution.

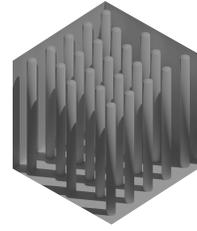
B. Normalizing Flows with Model Predictive Path Integral Control (FlowMPPI)

We implement a variation of FlowMPPI [15] and train it within a Bayesian model-based reinforcement learning paradigm [16]. The implementation follows Algorithm 1 with a conditional normalizing flow representing the optimal control distribution.

This flow is conditioned on the task variables (the start state and goal state), and the encoding of the quadrotor’s immediate environment (using a variational autoencoder [29]). We refer to the concatenated task variables and environment encoding as the context vector C . In practice, this ensures that the samples drawn from the optimal control distribution are contextually informed; they are goal directed (move from the start state to the goal state) and avoid collisions (by accounting for the immediate environment).



(a) The Block environment features a block obstacle with 14 meter side lengths suspended in the center of an enclosed $30 \times 30 \times 30$ cubic meter space.



(b) The Pillars environment features 23 cylindrical pillar obstacles (each 2 meters in diameter) in an enclosed $30 \times 30 \times 30$ cubic meter space.



(c) The Tunnels environment features interconnecting mineshafts with cross-sectional areas of approximately 3×3 meters and miscellaneous debris (e.g. rubble piles, barrels). The wall surfaces vary between naturally jagged and smooth concrete. This environment is Section 6 of DARPA’s Subterranean Challenge Final Event dataset, has extents of $31 \times 59 \times 4$ meters, and an internal volume of 536 cubic meters.



(d) The Chamber environment is a large cavernous space with piles of rubble, two mineshafts (geometrically similar to the Tunnels environment), and a ramp leading to an elevated area. This environment is Section 4 and Section 5 of DARPA’s Subterranean Challenge Final Event dataset, has extents of $41 \times 62 \times 11$ meters, and an internal volume of 11492 cubic metres.

Fig. 2: The four environments used in this work. Each represents a simulated (Figure 2a and 2b) or real-world (Figure 2c and 2d) cave environment. Learning-based methods trained in one environment may not necessarily generalize to different environments during testing, allowing us to test out-of-distribution performance.

C. Sequential Convex Programming for Trajectory Optimization

We determine non-colliding, dynamically feasible trajectories using sequential convex programming (SCP). The method iteratively refines an initial trajectory by solving a sequence of convex optimization problems [17].

We begin by computing a volume that 1) contains the start state s_{init} and goal state s_{goal} , and 2) contains no environmental obstacles. We compute this volume as a set of V spheres by first solving for the shortest path between s_{init} and s_{goal} using the A* search algorithm [30]. We dilate the voxel map representation of the environment by the quadrotor’s radius to ensure that there would be no collision even if the quadrotor were centered at any point along the resulting path.

The first sphere is set to be centered at s_{init} and is iteratively grown from a zero radius until it intersects with the environment. The next sphere is centered at the point that is both on the A* path and on the previous sphere’s surface. Again we grow the sphere until collision, and this process repeats until a sphere contains s_{goal} . The union of these V spheres is a contiguous volume that contains s_{init} and s_{goal} .

We initialize the solver with a state trajectory of shape $(T, 12)$, where the positional elements (first three elements

of the 12-dimensional state vector) are the A* path solutions, and all other values (velocity, rotation, and angular velocity) are zero. This initial state trajectory is dynamically infeasible but will be iteratively improved upon. We initialize the action trajectory with a set of controls that causes the quadcopter to hover in place. We find that this trajectory initialization warm-starts the solver and leads to faster solve times [17].

The optimization is a multiobjective cost function that simultaneously minimizes deviation from the goal, control effort, navigation slack (staying within the non-collision volume V), and dynamics slack (discussed further below):

$$J = w_{\text{dist}} \sum_{t=1}^T |s_t - s_{\text{goal}}|^2 + w_{\text{ctrl}} \sum_{t=1}^{T-1} |u_t|^2 + w_{\text{nav}} \sum_{t=1}^T \sum_{v=1}^V \xi_{t,v}^{\text{nav}} + w_{\text{dyn}} \max |\xi^{\text{dyn}}| \quad (1)$$

The trajectory is constrained to start and end at fixed states

$$s_1 = s_{\text{init}} \quad \text{and} \quad s_T = s_{\text{goal}} \quad (2)$$

Control inputs are constrained to dynamically feasible limits

$$u_t \in [u_{\text{min}}, u_{\text{max}}], \quad \forall t \in [1, T-1] \quad (3)$$

Algorithm 1 Model Predictive Path Integral Control

Require: initial state s_1 , samples K , horizon H ,
forward dynamics function Dyn ,
optimal control distribution \mathcal{P}^* ,
cost function $C(\cdot)$, weighting sharpness parameter λ .

- 1 // Store state and control trajectories and their costs
- 2 $\tau_s \leftarrow []$; $\tau_u \leftarrow []$; $C \leftarrow []$
- 3 for $k = 1$ to K do
- 4 $\tau_s^k \leftarrow [s_1]$; $\tau_u^k \leftarrow []$
- 5 for $h = 1$ to H do
- 6 // Execute the sampled control and get the next state
- 7 $u_h^k \leftarrow \text{sample}(\mathcal{P}^*)$
- 8 $s_h^k \leftarrow \text{Dyn}(\tau_s^k[-1], u_h^k)$
- 9 // Append the state and control to the trajectories
- 10 append(τ_s^k, s_h^k); append(τ_u^k, u_h^k)
- 11 end for
- 12 // Store state and control trajectory k
- 13 append(τ_s, τ_s^k); append(τ_u, τ_u^k)
- 14 // Relate the trajectory to a cost and store
- 15 append($C, \sum_{h=1}^H C(\tau_s^k[h], \tau_u^k[h])$)
- 16 end for
- 17 // Compute weights for each trajectory based on costs
- 18 $w \leftarrow \text{softmax}(-\lambda \cdot C)$
- 19 // Update the optimal control distribution (as needed)
- 20 $\mathcal{P}^* \leftarrow \text{update}(\cdot)$
- 21 // Combine optimal cost-weighted control sequence,
- 22 // and return the immediate (first) optimal control
- 23 $u^* \leftarrow (\sum_{k=1}^K w_k \tau_u^k[1])$
- 24 return u^*

Successive states are constrained to satisfy the affined dynamics with some allowed slack at each timestep ξ^{dyn}

$$s_{t+1} = A_t s_t + B_t u_t + C_t + \xi_t^{\text{dyn}}, \forall t \in [1, T-1] \quad (4)$$

A_t , B_t , and C_t represent the dynamics affined about the previous iteration's states s_t and controls u_t . To maintain the accuracy of this affine approximation we enforce a trust region

$$|s_t - s_t^{\text{prev}}| \leq r_s, \forall t \in [1, T] \quad (5)$$

$$|u_t - u_t^{\text{prev}}| \leq r_u, \forall t \in [1, T-1] \quad (6)$$

and the maximum dynamics slack must stay below some defined threshold $\xi_{\text{max}}^{\text{dyn}}$

$$\max |\xi_t^{\text{dyn}}| \leq \xi_{\text{max}}^{\text{dyn}} \quad (7)$$

Note that the trust region and maximum allowed dynamics slack are geometrically decayed at each iteration according to a decay schedule (Table II). We define a function $D(\cdot)$ which takes a state trajectory and returns a (T, V) -shaped matrix where the element t, v represents the signed distance field value of the t^{th} state's position

with respect to the v^{th} collision-free sphere [31]. We also introduce a navigation slack variable ξ^{nav} which is solved for at each iteration, and represents how much we deviated from the non-collision volume in that iteration. We introduce the first navigation constraint over every element of the navigation slack:

$$\xi_t^{\text{nav}} \leq D(s_t), \forall t \in [1, T] \quad (8)$$

which states that the navigation slack ξ^{nav} in this iteration must be less than the true signed distance field values from the current trajectory solution. We then introduce the second navigation constraint:

$$\text{diag} \left(\frac{e^{\sigma D(s^{\text{prev}})}}{\sum e^{\sigma D(s^{\text{prev}})}} (\xi^{\text{nav}} - D(s^{\text{prev}}))^T \right) + \frac{\log \sum e^{\sigma D(s^{\text{prev}})}}{\sigma} \geq 0 \quad (9)$$

which uses a softmax with temperature σ to ensure that each point in the trajectory is moving more towards the center of a non-colliding sphere in successive iterations [17].

This problem is solved iteratively until the solution does not meaningfully improve with respect to the objective J (Equation 1).

D. Augmented Lagrangian Iterative Linear Quadratic Regulator (AL-iLQR)

The Iterative Linear Quadratic Regulator (iLQR) is an optimization framework which leverages the Linear Quadratic Regulator (LQR) control policy and applies it to systems with nonlinear dynamics, using iterative dynamics linearization about the nominal trajectory s^{nom} , u^{nom} . The optimal controller follows the control law:

$$u_t = u_t^{\text{nom}} + K_t (s_t - s_t^{\text{nom}}) + d_t \quad (10)$$

where K_t is the feedback gain, d_t is the feedforward gain, s_t is the state for timestep $t \in [1, T]$, and u_t is the control input for timestep $t \in [1, T-1]$.

There are many practical constraints for a quadrotor controller that must be enforced, such as control bounds on rotor inputs. Traditional iLQR does not allow for this, so we use Augmented Lagrangians (AL), which reposes a standard optimization problem with explicit constraints into a problem with soft constraints enforced in the cost function [18]. The new optimization cost after applying the AL technique is referred to as the Augmented Lagrangian \mathcal{L}_A , and includes a Lagrange multiplier λ and penalty matrix I_μ . Using AL within the iLQR tracking framework results in the AL-iLQR tracking problem:

$$\min_{u_{1:T-1}} \mathcal{L}_A(s_{1:T}, u_{1:T-1}) \quad (11a)$$

$$\text{s.t.} \quad s_{t+1} = A_t s_t + B_t u_t, \quad \forall t = [1, T-1] \quad (11b)$$

which states that a control sequence that minimizes the AL:

$$\mathcal{L}_A(s_{1:T}, u_{1:T-1}) = \ell_T + \left(\lambda_T^\top + \frac{1}{2} c_T^\top I_{\mu, T} \right) c_T + \sum_{i=1}^{T-1} \left(\ell_i + \left(\lambda_i^\top + \frac{1}{2} c_i^\top I_{\mu, i} \right) c_i \right) \quad (12)$$

whilst maintaining dynamic feasibility, must be found. We further define bounds on the control inputs

$$[c_t(s_t, u_t), c_T(s_T)] = \left[\begin{bmatrix} u_t - u_{\max} \\ u_{\min} - u_t \end{bmatrix}, 0 \right] \quad (13)$$

and define the standard iLQR tracking costs, which penalize divergence from the desired state trajectory s^{track} using state cost, control cost, and terminal state cost matrices Q , R , and Q_T :

$$\begin{aligned} \ell_t(s_t, u_t) &= (s_t - s_t^{\text{track}})^\top Q (s_t - s_t^{\text{track}}) + u_t^\top R u_t \\ \ell_T(s_T) &= \frac{1}{2} (s_T - s_T^{\text{track}})^\top Q_T (s_T - s_T^{\text{track}}) \end{aligned} \quad (14)$$

Which are minimized.

IV. Results

A. The Learning-based Controller Completes the Task Faster but is More Sensitive to OOD Scenarios

In Table I we observe that FlowMPPI completes tasks the fastest out of all our controllers; in all environments apart from Block, FlowMPPI has the fastest task completion time for successful tasks, demonstrating the value of model-based reinforcement learning for InD inputs.

However, we observe a large drop in success rate when going from InD to OOD when using FlowMPPI. In the small environments, success rate drops from 100% to 71%, and in the large environments, success rate drops from 93% to 76%, illustrating the method’s OOD-sensitivity. We note that the failure of learning-based methods to generalize to new environments is well-observed [11].

B. The Safety Controller Completes the Task Slower but is More Robust to OOD Scenarios

Compared to the learning-based controller, the safety controller completes tasks much slower; in all four of the tested environments, the safety controller has the slowest completion times by far (see Table I). However, the controller exhibits a less severe drop in success rate compared to FlowMPPI when going from InD to OOD. In the small environments, success rate drops from 100% to 94%, and in the large environments, success rate drops from 88% to 86%. We also note that overall the safety controller produces smoother, lower control effort, shorter trajectories. We expect that this is due to SCP globally optimizing the entire path, whereas FlowMPPI and MPPI are myopic due to their local horizons.

These results experimentally show that the safety controller is more resilient to OOD than the learning-based controller, despite being much slower. This illustrates the tradeoff between liveness and safety properties [14].

C. Combining both Controllers with an OOD Detector Improves Task Completion Speed and OOD Resilience

Switching between the learning-based and safety controller based on an OOD score results in a controller that exhibits the best properties of its constituents. In

the small environments (Block and Pillars) the combined controller achieves comparable success rate (within 1%) to AL-iLQR, yet the average task completion speed is greatly reduced on account of FlowMPPI being used while InD (see Table I).

We note that the combined method having a lower success rate than the safety controller suggests that either 1) the OOD scenarios may not be perfectly classified, or 2) that the safety controller has a higher failure rate when switched to. Both scenarios are plausible; OOD classification is imperfect and the combined controller often switches to the safety controller when in the presence of close obstacles, or dynamically challenging states (i.e. states not encountered during training).

Table I shows a similar pattern on the large environments (Tunnels and Chamber). The combined controller achieves comparable success rates to the safety controller (within 2%), yet greatly reduced completion times due to the InD performance of the learning-based constituent controller.

V. Conclusion

In this work we have used out-of-distribution detection to determine when a learning-based controller and a safety controller should be used in the context of a point-to-point quadrotor navigation task in an underground 3D environment.

We find that learning-based methods (specifically, FlowMPPI) complete tasks quickly when in-distribution, but exhibit lower success rates when out-of-distribution. We find that our safety controller (based on sequential convex programming and AL-iLQR) is more robust to out-of-distribution scenarios, but exhibits overall slower task completion. By combining both, we create a controller that completes task quickly when in-distribution and maintains safety when out-of-distribution.

VI. Acknowledgments

Toyota Research Institute (TRI) provided funds to assist the authors with their research, but this article solely reflects the opinions and conclusions of its authors and not TRI or any other Toyota entity.

References

- [1] M. Petrlik, P. Petracek, V. Kratky, T. Musil, Y. Stasinchuk, M. Vrba, T. Baca, D. Hert, M. Pecka, T. Svoboda et al., “Uavs beneath the surface: Cooperative autonomy for subterranean search and rescue in darpa subt,” *Field Robotics Special Issue: DARPA Subterranean Challenge, Advancement and Lessons Learned from the Finals*, vol. 3, no. 1, pp. 1–68, January 2023.
- [2] R. Bogue, “Disaster relief, and search and rescue robots: the way forward,” *Industrial Robot: The International Journal of Robotics Research and Application*, vol. 46, no. 2, pp. 181–187, 2019.
- [3] C. Papachristos, S. Khattak, F. Mascarich, and K. Alexis, “Autonomous navigation and mapping in underground mines using aerial robots,” in *IEEE Aerospace Conference*, 2019.
- [4] G. U. Sikakwe, “Mineral exploration employing drones, contemporary geological satellite remote sensing and geographical information system (gis) procedures: A review,” *Remote Sensing Applications: Society and Environment*, vol. 31, p. 100988, 2023.

TABLE I: Test results across all environments. Methods are trained and tuned on the in-distribution environment before evaluation on the corresponding out-of-distribution environment (meaning that we train on Block before evaluating on Block and Pillars, and we train on Tunnels before evaluating on Tunnels and Chamber). We consider cases where the training environment matches the test environment to be in-distribution, and cases where the training environment does not match the test environment to be out-of-distribution. $N = 100$ randomized start and end goals are pre-generated for each environment so that each method can be tested consistently. Abbreviations: success rate SR, total time to complete task \bar{T}_{done} , average quadrotor velocity \bar{v} , average flight trajectory length \bar{d} , and average total control effort throughout trajectory \bar{u} .

Methods	Block (small, in-distribution)					Pillars (small, out-of-distribution)				
	SR (%)	\bar{T}_{done} (s)	\bar{v} (ms^{-1})	\bar{d} (m)	\bar{u} (N)	SR (%)	\bar{T}_{done} (s)	\bar{v} (ms^{-1})	\bar{d} (m)	\bar{u} (N)
AL-iLQR	100	40.32	2.24	43.22	2.62	94	60.29	2.85	54.39	3.27
MPPI	98	34.58	5.03	58.50	4.19	93	37.10	5.47	42.99	4.02
FlowMPPI	100	36.15	4.98	54.57	4.76	71	28.16	4.70	35.32	4.96
Combined (ours)	99	39.33	4.74	51.22	4.83	92	51.18	3.01	49.68	3.99
Methods	Tunnels (large, in-distribution)					Chamber (large, out-of-distribution)				
	SR (%)	\bar{T}_{done} (s)	\bar{v} (ms^{-1})	\bar{d} (m)	\bar{u} (N)	SR (%)	\bar{T}_{done} (s)	\bar{v} (ms^{-1})	\bar{d} (m)	\bar{u} (N)
AL-iLQR	93	144.81	2.18	64.12	2.14	86	133.47	2.06	39.97	2.63
MPPI	90	45.60	3.54	91.97	5.62	78	40.54	6.59	53.96	4.14
FlowMPPI	88	43.11	3.89	78.38	5.11	76	38.43	6.04	48.13	4.91
Combined (ours)	92	46.44	3.75	80.20	5.02	84	50.52	5.11	45.76	4.22

- [5] O. Peltzer, A. Bouman, S.-K. Kim, R. Senanayake, J. Ott, H. Delecki, M. Sobue, M. J. Kochenderfer, M. Schwager, J. Burdick et al., “Fig-op: Exploring large-scale unknown environments on a fixed time budget,” in IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), 2022.
- [6] P. Petráček, V. Krátký, M. Petrlik, T. Báča, R. Kratochvíl, and M. Saska, “Large-scale exploration of cave environments by unmanned aerial vehicles,” IEEE Robotics and Automation Letters, vol. 6, no. 4, pp. 7596–7603, 2021.
- [7] T. N. Titus, J. J. Wynne, M. J. Malaska, A.-a. Agha-Mohammadi, P. B. Buhler, E. C. Alexander, J. W. Ashley, A. Azua-Bustos, P. J. Boston, D. L. Buczkowski et al., “A roadmap for planetary caves science and exploration,” Nature Astronomy, vol. 5, no. 6, pp. 524–525, 2021.
- [8] B. J. Morrell, M. S. da Silva, M. Kaufmann, S. Moon, T. Kim, X. Lei, C. Patterson, J. Uribe, T. S. Vaquero, G. J. Correa et al., “Robotic exploration of martian caves: Evaluating operational concepts through analog experiments in lava tubes,” Acta Astronautica, vol. 223, pp. 741–758, 2024.
- [9] J. J. Wynne, T. N. Titus, A.-a. Agha-Mohammadi, A. Azua-Bustos, P. J. Boston, P. de León, C. Demirel-Floyd, J. De Waele, H. Jones, M. J. Malaska et al., “Fundamental science and engineering questions in planetary cave exploration,” Journal of Geophysical Research: Planets, vol. 127, no. 11, p. e2022JE007194, 2022.
- [10] E. Kaufmann, L. Bauersfeld, A. Loquercio, M. Müller, V. Koltun, and D. Scaramuzza, “Champion-level drone racing using deep reinforcement learning,” Nature, vol. 620, no. 7976, pp. 982–987, 2023.
- [11] D. Ghosh, J. Rahme, A. Kumar, A. Zhang, R. P. Adams, and S. Levine, “Why generalization in RL is difficult: Epistemic POMDPs and implicit partial observability,” in Advances in Neural Information Processing Systems (NeurIPS), 2021.
- [12] J. Guérin, K. Delmas, R. Ferreira, and J. Guiochet, “Out-of-distribution detection is not all you need,” in AAAI Conference on Artificial Intelligence (AAAI), 2023.
- [13] C. Dawson, A. Jasour, A. Hofmann, and B. Williams, “Provably safe trajectory optimization in the presence of uncertain convex obstacles,” in IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), 2020.
- [14] E. Kindler, “Safety and liveness properties: A survey,” Bulletin of the European Association for Theoretical Computer Science, vol. 53, no. 268-272, p. 30, 1994.
- [15] T. Power and D. Berenson, “Variational inference mpc using normalizing flows and out-of-distribution projection,” Robotics: Science and Systems (RSS), 2022.
- [16] M. Okada and T. Taniguchi, “Variational inference mpc for bayesian model-based reinforcement learning,” in Conference on Robot Learning (CoRL), 2020, pp. 258–272.
- [17] D. Malyuta, T. P. Reynolds, M. Szmuk, T. Lew, R. Bonalli, M. Pavone, and B. Açıkmeşe, “Convex optimization for trajectory generation: A tutorial on generating dynamically feasible trajectories reliably and efficiently,” IEEE Control Systems Magazine, vol. 42, no. 5, pp. 40–113, 2022.
- [18] B. Jackson, “AL-iLQR Tutorial,” 2022, accessed: 6 November 2024. [Online]. Available: {https://bjack205.github.io/papers/AL_iLQR_Tutorial.pdf}
- [19] M. Morari and J. H. Lee, “Model predictive control: past, present and future,” Computers & Chemical Engineering, vol. 23, no. 4-5, pp. 667–682, 1999.
- [20] G. Williams, A. Aldrich, and E. A. Theodorou, “Model predictive path integral control: From theory to parallel computation,” Journal of Guidance, Control, and Dynamics, vol. 40, no. 2, pp. 344–357, 2017.
- [21] D. Rezende and S. Mohamed, “Variational inference with normalizing flows,” in International Conference on Machine Learning (ICML), 2015.
- [22] Y. Althman and D. Gu, “Quadrotor transporting cable-suspended load using iterative linear quadratic regulator (iLQR) optimal control,” in IEEE Computer Science and Electronic Engineering (CEEC), 2016, pp. 168–173.
- [23] V. Sumathy and D. Ghose, “Adaptive augmented torque control of a quadcopter with an aerial manipulator,” Drone Systems and Applications, vol. 10, no. 1, pp. 26–50, 2021.
- [24] J. Yang, K. Zhou, Y. Li, and Z. Liu, “Generalized out-of-distribution detection: A survey,” International Journal of Computer Vision, vol. 132, no. 12, pp. 5635–5662, 2024.
- [25] C. Agia, R. Sinha, J. Yang, Z.-a. Cao, R. Antonova, M. Pavone, and J. Bohg, “Unpacking failure modes of generative policies: Runtime monitoring of consistency and progress,” in Conference on Robot Learning (CoRL), 2024.
- [26] S. Ramakrishna, Z. Rahiminasab, G. Karsai, A. Easwaran, and A. Dubey, “Efficient out-of-distribution detection using latent space of β -vae for cyber-physical systems,” ACM Transactions on Cyber-Physical Systems (TCPS), vol. 6, no. 2, pp. 1–34, 2022.
- [27] D. Gudovskiy, S. Ishizaka, and K. Kozuka, “Cflow-ad: Real-time unsupervised anomaly detection with localization via

conditional normalizing flows,” in IEEE/CVF Winter Conference on Applications of Computer Vision, 2022, pp. 98–107.

- [28] C. Lazarus, J. G. Lopez, and M. J. Kochenderfer, “Runtime safety assurance using reinforcement learning,” in Digital Avionics Systems Conference (DASC), 2020.
- [29] D. P. Kingma, M. Welling et al., “Auto-encoding variational bayes,” in International Conference on Learning Representations (ICLR). Banff, Canada, 2014.
- [30] P. E. Hart, N. J. Nilsson, and B. Raphael, “A formal basis for the heuristic determination of minimum cost paths,” IEEE Transactions on Systems Science and Cybernetics, vol. 4, no. 2, pp. 100–107, 1968.
- [31] H. Oleynikova, A. Millane, Z. Taylor, E. Galceran, J. Nieto, and R. Siegwart, “Signed distance fields: A natural representation for both mapping and planning,” in RSS Workshop: Geometry and Beyond—Representations, Physics, and Scene Understanding for Robotics. University of Michigan, 2016.
- [32] F. Sabatino, “Quadrotor control: Modeling, nonlinear control design, and simulation,” Master’s thesis, KTH, School of Electrical Engineering (EES), Automatic Control, October 2015.

Appendix

A. Dynamics

The system has a 12-dimensional state-space representing the 3D position, velocity, rotation (Euler angles), and angular velocity, and a 4-dimensional control-space representing the angular velocities of the rotors. Our discretized dynamics model follows the implementation outlined in [32], and assumes negligible motor resistance (i.e. motor power is assumed to be proportional to angular velocity), stationary air, linear fluid friction, negligible rotational drag, and does not account for blade flapping:

$$\begin{bmatrix} x \\ y \\ z \\ r_z \\ r_y \\ r_x \\ \dot{x} \\ \dot{y} \\ \dot{z} \\ \omega_x \\ \omega_y \\ \omega_z \end{bmatrix}_{t+1} = \begin{bmatrix} x \\ y \\ z \\ r_z \\ r_y \\ r_x \\ \dot{x} \\ \dot{y} \\ \dot{z} \\ \omega_x \\ \omega_y \\ \omega_z \end{bmatrix}_t + \Delta t \begin{bmatrix} \dot{x} \\ \dot{y} \\ \dot{z} \\ \omega_y \frac{\sin(r_x)}{\cos(r_y)} + \omega_z \frac{\cos(r_x)}{\cos(r_y)} \\ \omega_y \cos(r_x) - \omega_z \sin(r_x) \\ \omega_x + \omega_y \sin(r_x) \tan(r_y) + \omega_z \cos(r_x) \tan(r_y) \\ -(\sin(r_x) \sin(r_z) + \cos(r_z) \cos(r_x) \sin(r_y)) k \frac{\sum_{i=1}^4 \omega_i}{m} \\ -(\cos(r_x) \sin(r_z) \sin(r_y) - \cos(r_z) \sin(r_x)) k \frac{\sum_{i=1}^4 \omega_i}{m} \\ g - \cos(r_x) \cos(r_y) k \frac{\sum_{i=1}^4 \omega_i}{m} \\ \frac{(I_y - I_z) \omega_y \omega_z + k r (\omega_3^2 - \omega_1^2)}{I_x} \\ \frac{(I_z - I_x) \omega_x \omega_z + k r (\omega_4^2 - \omega_2^2)}{I_y} \\ \frac{(I_x - I_y) \omega_x \omega_y + b((\omega_3^2 + \omega_4^2) - (\omega_1^2 + \omega_2^2))}{I_z} \end{bmatrix} \quad (15)$$

where (x, y, z) is the position, $(\dot{x}, \dot{y}, \dot{z})$ is the velocity, (r_z, r_y, r_x) is the 3D rotation (Euler angles), $(\omega_x, \omega_y, \omega_z)$ are the body’s angular velocities, $(\omega_1, \omega_2, \omega_3, \omega_4)$ are the rotor’s angular velocities, k is the thrust coefficient, m is the quadrotor’s mass, b is the yaw-drag coefficient, r is the quadrotor’s radius, and (I_x, I_y, I_z) are the quadrotor’s moments of inertia about the $(\hat{x}, \hat{y}, \hat{z})$ axes. The gravitational acceleration is $g = 9.81 \text{m/s}^2$, and the time step is $\Delta t = 0.025 \text{s}$. Control bounds are $u_{\min} = 0$ and $u_{\max} = 4$.

B. Reproducibility

All experiments were executed on a workstation with an AMD Ryzen Threadripper 3970X 32-Core 64-Thread Processor (3.69 GHz) CPU and 128 GB of usable RAM. The training of neural networks was executed

on NVIDIA GeForce RTX 3090 GPUs using PyTorch Lightning. The code and data required to reproduce these results is available at <https://github.com/sisl/underground>. The hyperparameters used in our algorithms are specified in Table II.

TABLE II: Descriptions of hyperparameters and their values as used in this work.

Sequential Convex Programming Hyperparameters	
Initial state trust region radius r_s	1
State trust region radius decay γ_s	0.95
Initial control trust region radius r_c	1
Control trust region radius decay γ_c	0.95
Maximum dynamics slack ξ_{\max}^{dyn}	1
Maximum dynamics slack decay γ_{dyn}	0.5
Navigation sigma σ	50
Minimum J improvement	0.01
Goal directedness term weight w_{dist}	0.1
Control effort term weight w_{ctrl}	0.1
Navigation slack term weight w_{nav}	0.01
Dynamics slack term weight w_{dyn}	100
AL-iLQR hyperparameters	
State cost matrix Q_k	$Q_{k,ij} = \begin{cases} 5, & i = j \wedge i \leq 3 \\ 1, & i = j \wedge i > 3 \\ 0, & i \neq j \end{cases}$
Final state cost matrix Q_K	$Q_{K,ij} = \begin{cases} 100, & i = j \wedge i \leq 3 \\ 10, & i = j \wedge i > 3 \\ 0, & i \neq j \end{cases}$
Action cost matrix R_k	I_4
Maximum # of path segments	3
# points per meter resampled to	10
MPPI hyperparameters	
Samples K	1024
Horizon H	50
Temperature λ	20
FlowMPPI hyperparameters	
Samples K	1024
Horizon H	50
Temperature λ	20
Neighborhood size	$6 \times 6 \times 6 \text{ m}^3$
VAE input size	864
VAE latent size	128
CNF # flow layers	8
CNF hidden layers’ size	256
CNF hidden layer type	AffineCouplingTransform
Context vector input size	152
Context vector embedding size	64
# training episodes	256
Combined method’s hyperparameters	
Normalized OOD threshold (log space)	3.12