

Benchmarking Reinforcement Learning Techniques for Autonomous Navigation

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Deep reinforcement learning (RL) has brought many successes for au-
2 tonomous robot navigation. However, there still exists important limitations that
3 prevent *real-world* use of RL-based navigation systems. For example, most learn-
4 ing approaches lack safety guarantees; and learned navigation systems may not
5 generalize well to unseen environments. Despite a variety of recent learning tech-
6 niques to tackle these challenges in general, a lack of an open-source benchmark
7 and reproducible learning methods specifically for autonomous navigation makes
8 it difficult for roboticists to choose what learning methods to use for their mobile
9 robots and for learning researchers to identify current shortcomings of general
10 learning methods for autonomous navigation. In this paper, we identify four ma-
11 jor desiderata of applying deep RL approaches for autonomous navigation: (D1)
12 *reasoning under uncertainty*, (D2) *safety*, (D3) *learning from limited trial-and-*
13 *error data*, and (D4) *generalization to diverse and novel environments*. Then, we
14 explore four major classes of learning techniques with the purpose of achieving
15 one or more of the four desiderata: *memory-based neural network architectures*
16 (D1), *safe RL* (D2), *model-based RL* (D2, D3), and *domain randomization* (D4).
17 By deploying these learning techniques in a new open-source large-scale naviga-
18 tion benchmark and real-world environments, we perform a comprehensive study
19 aimed at establishing to what extent can these techniques achieve these desiderata
20 for RL-based navigation systems.

21 **Keywords:** RL, Autonomous navigation, Benchmark

22 1 Introduction

23 Autonomous robot navigation, i.e., moving a robot from one point to another without colliding
24 with any obstacle, has been studied by the robotics community for decades. Classical navigation
25 systems [1, 2] can successfully solve such navigation problem in many real-world scenarios, e.g.,
26 handling noisy, partially observable sensory input but still providing verifiable collision-free safety
27 guarantees. However, these systems require extensive engineering effort, and can still be brittle
28 in challenging scenarios, e.g., in highly constrained environments [3]. Recently, data-driven ap-
29 proaches have also been used to tackle the navigation problem [4] thanks to advances in the machine
30 learning community. In particular, Reinforcement Learning (RL), i.e., learning from self-supervised
31 trial-and-error data, has achieved tremendous progress on multiple fronts, including safety [5, 6, 7],
32 generalizability [8, 9, 10, 11], sample efficiency [12, 13], and addressing temporal data [14, 15, 16].
33 For the problem of navigation, learned navigation systems from RL [17] have the potential to re-
34 lieve roboticists from extensive engineering efforts [18, 19, 20, 21, 22] spent on developing and
35 fine-tuning classical systems. Moreover, a simple case study conducted in five randomly generated
36 obstacle courses where classical navigation systems often fail shows that a RL-based navigation
37 has the potential to achieve superior behaviors in terms of both navigation efficiency and collision
38 avoidance (Fig. 1 left).

39 Despite such promising advantages, learning-based navigation systems are far from finding their
40 way into real-world robotics use cases, which currently still rely heavily on their classical counter-

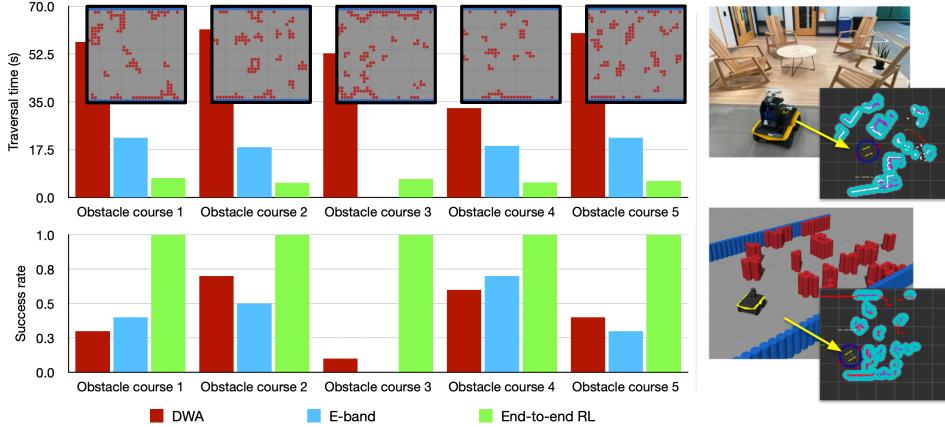


Figure 1: **Left:** Traversal time (top) and success rate (bottom) of two classical navigation systems, DWA [2] (red) and E-band [1] (blue), and vanilla end-to-end RL-based navigation systems (individually trained) in five randomly generated difficult obstacle courses (green). The insets at the top show top-down views of the five obstacle courses. **Right:** Navigation environments in the real world (top) and the proposed benchmark (bottom) are similar to the robot perception system (e.g., white/red laser scans and cyan/purple costmaps).

41 parts. Such reluctance in adopting learning-based systems in the real world stems from a series of
 42 fundamental limitations of learning methods, e.g., lack of safety, explainability, and generalizability.
 43 To make things even worse, a lack of well-established comparison metrics and reproducible learning
 44 methods further obfuscates the effects of different learning approaches on navigation across both the
 45 robotics and learning communities, making it difficult to assess the state of the art and therefore to
 46 adopt learned navigation systems in the real world.
 47 To facilitate research in developing RL-based navigation systems with the goal of deploying them
 48 in real-world scenarios, we introduce a new open-source large-scale navigation benchmark with a
 49 variety of challenging, highly constrained obstacle courses to evaluate different learning approaches,
 50 along with the implementation of several state-of-the-art RL algorithms. The obstacle courses re-
 51 semble highly-constraint real-world navigation environments (right of Fig. 1), and present major
 52 challenges to existing classical navigation systems, while RL-based navigation systems have the
 53 potential to perform well in them (left of Fig. 1).
 54 We identify four major desiderata that ought to be fulfilled by any learning-based system that is to
 55 be deployed: (D1) *reasoning under uncertainty of partially observed sensory inputs*, (D2) *safety*,
 56 (D3) *learning from limited trial-and-error data*, and (D4) *generalization to diverse and novel envi-*
 57 *ronments*. By deploying four major classes of learning techniques: *memory-based neural network*
 58 *architectures*, *safe RL*, *model-based RL*, and *domain randomization*, we perform extensive experi-
 59 ments and empirically compare a large range of RL-based methods based on the degree to which
 60 they achieve each of these desiderata. Each trial of the method runs for about 8 hours and uses
 61 100 CPU cores for distributed training of the RL agents. The whole set of experiments includes
 62 120 independent trials which account for about 100000 hours of CPU time in total. Moreover, by
 63 deploying 6 selected navigation systems in three qualitatively different real-world navigation envi-
 64 ronments, we investigate to what degree the conclusions drawn from the benchmark can be applied
 65 to the real world.

66 2 Desiderata for Learning-based Autonomous Navigation

67 In this section, we introduce four desiderata for learning-based autonomous navigation systems and
 68 briefly discuss the learning techniques as their corresponding solutions. More detailed descriptions
 69 of these techniques can be found in Appendix A.

70 **(D1) reasoning under uncertainty of partially observed sensory inputs.** Autonomous navigation
 71 without explicit mapping and localization is usually formalized as a Partially Observable Markov

72 Decision Process (POMDP), where the agent produces the motion of the robot only based on limited
73 sensory inputs that are usually not sufficient to recover the full state of the navigation environment.
74 Most RL approaches solve POMDPs by maintaining a history of past observations and actions
75 [14, 15]. Then, neural network architectures like Recurrent Neural Networks (RNNs) that process
76 sequential data are employed to encode history and address partial observability. In this study, we
77 investigate various design choices of history-dependent architectures and history lengths.

78 **(D2) safety.** Even though in some cases deep RL methods achieve comparable performance to
79 classical navigation, they still suffer from poor explainability and do not guarantee collision-free
80 navigation. The lack of safety guarantee is a major challenge preventing RL-based navigation from
81 being used in the real world. Prior works have addressed this challenge by formalizing the navigation
82 as a multi-objective problem that treats collision avoidance as a separate objective from reaching the
83 goal and solving it with Lagrangian or Lyapunov-based methods [5]. For simplicity, we only explore
84 Lagrangian method and investigate whether explicitly treat safety as a separate objective leads to
85 safer and smoother learned navigation behavior.

86 **(D3) learning from limited trial-and-error data.** Although deep RL approaches can alleviate
87 roboticists from extensive engineering effort, a large amount of data is still required to train a typical
88 deep RL agent. However, autonomous navigation data is usually expensive to collect in the real
89 world. Therefore, data collection is usually conducted in simulation, e.g., in the Robot Operating
90 System (ROS) Gazebo simulator, which provides an easy interface with real-world robots. However,
91 simulating a full navigation stack from perception to actuation is more computationally expensive
92 compared to other RL domains, e.g., MuJuCo or Atari games [23, 24], which presents a high re-
93 quirement for sample efficiency. Most prior works have used off-policy RL algorithms to improve
94 sample efficiency with experience replay [25, 26]. In addition, model-based RL methods can ex-
95 plicitly improve sample efficiency, and are widely used in robot control problems. In this study, we
96 compare two common classes of model-based RL method [12, 13] combined with an off-policy RL
97 algorithm, and empirically study to what extent model-based approaches improve sample efficiency
98 when provided with different amounts of data.

99 **(D4) generalization to diverse and novel environments.** The ultimate goal of deep RL approaches
100 for autonomous navigation is to learn a generalizable policy for all kinds of navigation environments
101 in the real world. A simple strategy is to train the agent in as many diverse navigation environments
102 as possible or domain randomization, but it is unclear what is the necessary amount of training envi-
103 ronments to efficiently achieve good generalization. Utilizing the large-scale navigation benchmark
104 proposed in this paper, we empirically study the dependence of generalization on the number of
105 training environments.

106 **3 Navigation Benchmark**

107 This section details the proposed navigation benchmark for RL-based navigation systems, which
108 aims to provide a unified and comprehensive testbed for future autonomous navigation research. In
109 Sec. 3.1 and 3.2, the navigation task is formally defined and formulated as a POMDP. More detailed
110 background of MDP and POMDP can be found in Appendix A.1. Finally, Sec. 3.3 introduces
111 simulated and real-world environments that benchmark different aspects of navigation performance.

112 **3.1 Navigation Problem Definition**

113 **Definition 1 (Robot Navigation Problem)** *Situated within a navigation environment e which in-
114 cludes information of all the obstacle locations at any time t , a start location (x_i, y_i) , a start orien-
115 tation θ_i , and a goal location (x_g, y_g) , the navigation problem T_e is to maximize the probability p of
116 a mobile robot reaching the goal location from the start location and orientation under a constraint
117 on the number of collisions with any obstacle $C < 1$ and a time limit $t < T_{max}$.*

118 A navigation problem can be formally defined as above. Given the current location (x_t, y_t) , the
119 robot is considered to have reached the goal location if and only if its distance to the goal location is
120 smaller than a threshold, $d_t < d_s$, where d_t is the Euclidean distance between (x_t, y_t) and (x_g, y_g) ,
121 and d_s is a constant threshold.

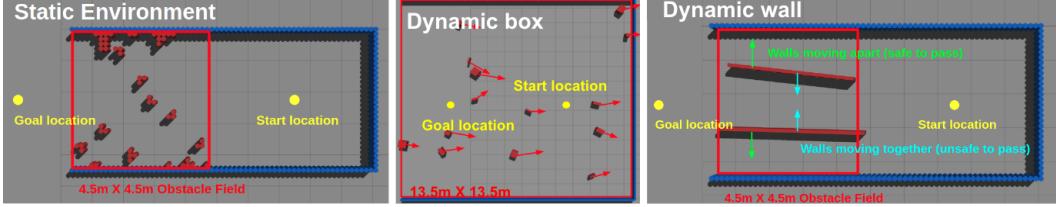


Figure 2: Three types of navigation environments: static, dynamic box, and dynamic-wall. The red squares mark the obstacle fields, and the yellow circles mark the start and goal locations. In dynamic-wall, the green (blue) arrows indicate the case when the two walls are moving apart (together). In dynamic box, the red arrows indicate the velocities of obstacles.

122 3.2 POMDP Formulation

123 A navigation task \mathcal{T}_e can be formulated as a POMDP conditioned on a navigation environment e ,
 124 which can be represented by a 7-tuple $(S_e, A_e, O_e, T_e, \gamma_e, R_e, Z_e)$. In this POMDP, the state $s_t \in$
 125 S_e is a 5-tuple $(x_t, y_t, \theta_t, c_t, e)$ with x_t, y_t, θ_t the two-dimensional coordinates and the orientation
 126 of the robot at time step t , c_t a binary indicator of whether a collision has occurred since the last
 127 time step $t - 1$, and e the navigation environment. The action $a_t = (v_t, \omega_t) \in A_e$ is a two-
 128 dimensional continuous vector that encodes the robot’s linear and angular velocity. The observation
 129 $o_t = (\chi_t, \bar{x}_t, \bar{y}_t) \in O_e$ is a 3-tuple composed of the sensory input χ_t from LiDAR scans and the
 130 relative goal position (\bar{x}_t, \bar{y}_t) in the robot frame. The observation model $Z : S \rightarrow O$ maps the state
 131 to the observation. The reward function for this POMDP is defined as follows:

$$R_e(s_t, a_t) = +b_f \cdot \mathbb{1}(d_t < d_s) + b_p \cdot (d_{t-1} - d_t) - b_c \cdot c_t, \quad (1)$$

132 where $\mathbb{1}(d_t < d_s)$ is the indicator function of reaching the goal location, d_t is the Euclidean distance
 133 to the goal location, and b_f, b_p, b_c are the coefficient constants. In this reward function, the first
 134 term is the true reward function that assigns a positive constant b_f for the success of an agent, which
 135 matches with the objective of the navigation task in Definition 1. The second and third terms are
 136 auxiliary rewards that facilitate the training by encouraging local progress and penalizing collisions.

137 We perform a grid search over different values of the coefficients in this reward function, and the
 138 result shows that the auxiliary reward term $(d_{t-1} - d_t)$ is necessary for successful training, and a
 139 much smaller coefficient b_p relative to b_f can lead to a better asymptotic performance. The agent
 140 can learn without the penalty reward for collision ($b_c = 0$), but a moderate value of b_c can improve
 141 the asymptotic performance and speed up training. For all the experiments in this paper, we fix the
 142 coefficients as $b_f = 20$, $b_p = 1$ and $b_c = 4$.

143 In our experiments, the RL algorithm solves a multi-task RL problem where the tasks
 144 are randomly sampled from a task distribution $\mathcal{T}_e \sim p(\mathcal{T}_e)$. Here the task distribution
 145 $p(\mathcal{T}_e) := U(\{e_i\}_{i=1}^N)$ is a uniform distribution on a set of N navigation environments $\{e_i\}_{i=1}^N$.
 146 The overall objective of this multi-task RL problem is to find an optimal policy $\pi^* =$
 147 $\max_{\pi} \mathbb{E}_{\mathcal{T}_e \sim p(\mathcal{T}_e), \tau_t \sim \pi} [\sum_{t=0}^{\infty} \gamma^t R_e(s_t, a_t)]$.

148 3.3 Navigation Environments

149 The navigation is performed by a ClearPath Jackal differential-drive ground robot in simulated
 150 by the Gazebo simulator. More details of the robot and simulation can be found in Appendix C.
 151 Each environment in this benchmark will have a navigation system navigating the robot through
 152 a 10m navigation path that passes through a highly constrained obstacle course. Walls are placed
 153 at three edges of a square so that passing through the obstacle field is the only path to the goal
 154 location (see Fig. 2). The benchmark includes 300 static environments, 100 dynamic-box
 155 environments, and 100 dynamic-wall environments. The static environments contains a di-
 156 verse set of obstacle course covering a large range of difficulty levels from easy to hard. A
 157 dynamic-box environment has small boxes with random shapes and velocities to test the sys-
 158 tem’s immediate reactions to small moving obstacles. A dynamic-wall has two walls mov-
 159 ing oppositely that requires the system to make a longer-term decision of whether to pass or
 160 wait. The detailed procedures of generating these environments can be found in the Appendix



Figure 3: Real-world benchmark-like (left), in-door highly-constrained (middle), and large-scale (right) environments. The yellow curves mark roughly the paths of navigation.

161 B. We randomly select 50 environments from each type as the test sets, which are denoted as
 162 static-test, dynamic-box-test, and dynamic-wall-test. The remaining environments
 163 are denoted as static-train, dynamic-box-train, and dynamic-wall-train respectively.
 164 To study the effect of randomization, static-train is further separated as static-train-5,
 165 static-train-10, static-train-50, static-train-100, and static-train-250 by ran-
 166 domly sampling 5, 10, 50, 100, and all 250 environments from static-train.
 167 To test the sim-to-real transferability of the policies learning with different techniques, the navi-
 168 gation systems are deployed in three qualitatively different static navigation environments includ-
 169 ing a benchmark-like environment (Fig. 3 left), an indoor highly-constrained environment (Fig. 3
 170 right), and a large-scale environment of 30 meters in length. We denote them as real-world-1,
 171 real-world-2, and real-world-3 respectively.

172 4 Experiments

173 In this section, we present experimental results of each studied technique to achieve the proposed
 174 desiderata in Sec. 2. Unless otherwise specified, all the experiments mentioned in this section use
 175 a distributed TD3 RL algorithm (similar to [21]) combined with the corresponding techniques, and
 176 all the data points presented are averaged over three independent runs.

	static env.			dynamic-box env.			dynamic-wall env.		
	H = 1	H = 4	H = 8	H = 1	H = 4	H = 8	H = 1	H = 4	H = 8
MLP	65 ± 4	57 ± 7	42 ± 2	50 ± 5	35 ± 2	46 ± 3	67 ± 7	72 ± 1	69 ± 4
GRU	-	51 ± 2	43 ± 4	-	48 ± 4	45 ± 1	-	82 ± 4	78 ± 5
CNN	-	55 ± 4	45 ± 5	-	42 ± 5	40 ± 1	-	63 ± 3	43 ± 3
Transformer	-	68 ± 2	46 ± 3	-	52 ± 1	44 ± 4	-	33 ± 28	15 ± 13

Table 1: (D1) Success rate (%) of policies trained with different neural network architectures and history lengths in static, dynamic-box, and dynamic-wall environments.

Methods	Baseline (model-free)	Lagrangian method	MPC (model-based)	DWA
Success rate (%)	65 ± 4	74 ± 2	70 ± 3	82
Survival time (s)	8.0 ± 1.5	16.2 ± 2.5	55.7 ± 4.9	62.7
Traversal time (s)	7.5 ± 0.3	8.6 ± 0.2	24.7 ± 2.0	35.6

Table 2: (D2) Success rate, survival time, and traversal time of policies trained with Lagrangian method, MPC with probabilistic transition model, and DWA.

177 4.1 Memory-based Neural Network Architectures (D1)

178 To benchmark the performance of different neural network (NN) architectures, deep RL policies rep-
 179 resented by architectures of Multilayer Perceptron (MLP), One-dimensional Convolutional Neural
 180 Network (CNN), Gated Recurrent Units (GRU), and Transformer with history length of 4 and 8 are
 181 trained in static-train-50, and the two types of dynamic environments dynamic-box-train
 182 and dynamic-wall-train from Sec. 3.3. After training, the policies are tested in their correspond-
 183 ing test sets. In addition, MLP with history length of one is added as a memory-less baseline. Table

184 1 shows the success rates of policies with different architectures and history lengths evaluated in
185 static-test (left), dynamic-wall-test (middle) and dynamic-box-test respectively.

186 **Memory-based NNs only marginally improve navigation performance in static environments.**
187 In Table 1, the policy represented by Transformer with a history length of 4 shows the best success
188 rate of 68%, with a slightly worse success rate of 65% achieved by the baseline MLP. Additionally,
189 a monotonic decrease in success rate with increasing history length is observed in each tested NN
190 architecture. For example, a 32% drop in the success rate of Transformer is shown by increasing the
191 history length from 4 to 8. One possible explanation is that, if only few past observations are useful
192 to make the decision, including more history will make it more difficult to learn a generalized policy
193 in this very diverse training set.

194 **Memory is essential when possible catastrophic failures will happen by making the wrong**
195 **long-term decisions.** Memory usually matters for dynamic environments when a single time frame
196 is not sufficient to estimate the motion of obstacles. Surprisingly, in dynamic-box where the dy-
197 namic obstacles are completely random, the memory-based NN architectures do not outperform the
198 memory-less baseline. On the other hand, in dynamic-wall with a manually designed dynamic
199 challenge, the best success rate of 82% is observed in GRU with a history length of 4, which im-
200 proves about 15% over the non-memory baseline. During our deployment of the policies, we observe
201 that, in dynamic-box even though the memory-less agent does not estimate the motion and adjust
202 its plan in advance, it tends to perform safely and avoids the obstacles when they get close enough.
203 This simple strategy works surprisingly well and achieves similar success rate as the memory-based
204 policies. However, this strategy does not work in the manually designed dynamic challenges like
205 dynamic-wall where the agent has to estimate the motion of the obstacles to pass safely.

206 4.2 Safe RL (D2)

207 To investigate to what extent safe RL methods can help to improve safety, a TD3 agent with the
208 Lagrangian-based safe RL method (Appendix A.3) is trained in static-train-50, and then tested
209 in static-test. The policy is represented by a MLP with its input containing only one history
210 length. Table 2 shows the success rate, average survival time, and average traversal time of the
211 safe RL agent trained with Lagrangian method and a baseline MLP agent tested in static-test.
212 We define survival time as the time cost of an unsuccessful episode (collision or exceeding a time
213 limit of 80s). Traversal time, instead, is the time cost of a successful episode. With the same level of
214 success rate, a longer survival time means that the agent tends to, at least, avoid collisions if it cannot
215 succeed. To compare the safe RL method with classical navigation systems which are believed to
216 have better safety, we also add evaluation metrics from a classical navigation stack with the Dynamic
217 Window Approach (DWA) [2] local planner.

218 **Lagrangian method reduces the gap between training and test environments.** When deployed
219 in the training environments, both the baseline MLP and the safe RL method achieves about 80%
220 success rate. However, in the test environments, the Lagrangian method has a better success rate
221 of 74% compare to 65% by the baseline MLP. We hypothesize that the safety constraint applied by
222 the safe RL methods forms a way of regularization, and therefore, improves the generalization to
223 unseen environments.

224 **Lagrangian method increases the average survival time in failed episodes.** As expected, the
225 Lagrangian method increases the average survival time by 8.2s compared to the baseline MLP at a
226 cost of 1.1s longer average traversal time. However, such improved safety are still worse than the
227 classical navigation systems given the best survival time of 88.6s achieved by DWA.

228 4.3 Model-based RL (D2 and D3)

229 To explore how the model-based approaches help with the autonomous navigation tasks, we im-
230 plement both Dyna-style and MPC methods (Appendix A.4), and evaluate the methods in static
231 environments. The transition models are either represented by a deterministic NN or a probabilistic
232 NN that predicts the mean and variance of the next state. During the training in static-train-50,

233 the policies are saved when 100k, 500k and 2000k transition samples are collected, then tested in
 234 static-test. The success rates of these policies are reported in Table 3.

235 **Model-based methods do not improve sample efficiency.** As shown in the second and third
 236 columns in Table 3, better success rates of 13% and 58% are achieved by the baseline MLP method
 237 provided by limited 100k and 500k transition samples respectively. In addition, Higher success
 238 rates at 500k transition samples are observed in probabilistic models compared to their deterministic
 239 counterparts, which indicates a more efficient learning with probabilistic transition models.

240 **Model-based methods with probabilistic dynamic models improve the asymptotic performance.** In the last column of Table 3, both Dyna-style and MPC with probabilistic dynamic models
 241 achieve slightly better success rates of 70% compared to 65% in the baseline MLP method when
 242 sufficient transition samples of 2000k are given to the learning agent.

243 **The MPC policy performs conservatively when deployed in unseen test environments and
 244 shows a better safety performance.** The safety performances of MPC policies with probabilistic
 245 dynamic models are also tested (see Table 2). We observe that the agents with MPC policies
 246 navigate very conservatively with an average traversal time of 24.7s, which is about two times more
 247 than the MLP baseline. In the meantime, MPC policies achieve improved safety with the best sur-
 248 vival time of 55.7s among the RL-based methods.

250 4.4 Domain Randomization (D4)

251 To explore how model generalization depends on the degree of randomness in the training environ-
 252 ments, baseline MLP policies with one history length are trained in the environment sets with five
 253 different sizes: static-train-5, static-train-10, static-train-50, static-train-100,
 254 and static-train-250. The trained policies are tested in the same static-test. To investigate
 255 the performance gap between training and test, the policies trained with 50, 100, and 250 environ-
 256 ments are also tested on static-train-50, which is part of their training sets. Fig 4 shows the
 257 success rate of policies trained with different number of training environments.

258 **The generalization to unseen environments monotonously improves with increasing number
 259 of training environments.** As shown in Fig. 4, the performances on the unseen test environments
 260 monotonously increase from 43% to 74% with the number of training environments increasing from
 261 5 to 250. Moreover, the gaps between training and test environments gradually shrink by adding
 262 more training environments provided by that the polices are robust enough to maintain similar per-
 263 formances of about 80% on the training environments.

Transition samples	100k	500k	2000k
MLP	13 ± 7	58 ± 2	65 ± 4
Dyna-style deterministic	8 ± 2	30 ± 10	66 ± 5
MPC deterministic	0 ± 0	21 ± 10	62 ± 3
Dyna-style probabilistic	0 ± 0	48 ± 4	70 ± 1
MPC probabilistic	0 ± 0	45 ± 4	70 ± 3

Table 3: (D3) Success rate (%) of policies trained with different model-based methods and different number of transition samples.

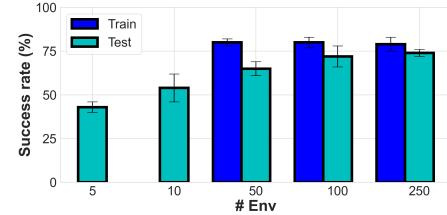


Figure 4: (D4) Success rate (%) of polices trained with different number of training environments.

264 4.5 Physical experiments

265 To study the consistency of the above observations in simulation and the real world, we deploy one
 266 baseline MLP policy, one best policy for each studied desideratum, and one classical navigation
 267 system (DWA [2]) in the three real-world environments introduced in Sec. 3.3. Each deployment is
 268 repeated three times, and the average traversal time and the number of successful trials are reported
 269 in Table. 4. Even though the best memory-based policy, transformer architecture with 4 history
 270 length, was only marginally better than the baseline MLP in simulation, in the real world it can nav-
 271 igate very smoothly and fails only once in real-world-2 and real-world-3, while baseline MLP
 272 fails most of the trials in all the environments including the benchmark-like environment. Similarly,

	history	# envs	real-world-1	real-world-2	real-world-3
MLP	1	50	6.9 (1/3)	10.6 (1/3)	N (0/3)
MLP	1	250	4.6 ± 0.8 (3/3)	6.6 ± 0.6 (3/3)	22.6 ± 0.5 (3/3)
Transformer	4	50	6.1 ± 0.4 (3/3)	6.1 ± 0.1 (2/3)	20.5 ± 2 (2/3)
Lagrangian	1	50	4.4 ± 0.6 (3/3)	7.1 ± 0.1 (2/3)	26.2 (1/3)
MPC	1	50	13.2 ± 0.7 (3/3)	24.8 ± 3.7 (3/3)	N (0/3)
DWA	-	-	16.2 ± 0.7 (3/3)	35.2 ± 8.2 (2/3)	66.9 ± 0.6 (3/3)

Table 4: Physical experiments. The table shows the traversal time (s) and the number of successful trials of 5 RL-based navigation systems and a classical navigation system (DWA) evaluated in `real-world-1`, `real-world-2`, and `real-world-3`.

273 MLP policy trained with 250 environments can successfully navigate in all the environments without
 274 any failures, while baseline MLP trained with 50 environments fails most of the trials. Safe RL
 275 improves the chances of success in all the environments and can navigate more safely by performing
 276 backups and small adjustments of robots’ poses. Similar to the simulation, MPC navigates very
 277 conservatively and succeeds in all the trials in `real-world-1` and `real-world-2`, but has much
 278 more difficulty generalizing to large-scale `real-world-3`.

279 5 Limitations

280 In this section, we discuss the limitations of this benchmark and the current stages of the studied
 281 learning techniques in solving the four desiderata. We also point out potential future directions to
 282 improve the benchmark and to better solve the desiderata.

283 **(D1) Reasoning under uncertainty of partially observed sensory inputs.** While our results do
 284 not show an obvious benefit by adding memory in static and very random dynamic (`dynamic-box`)
 285 environments, much more significant improvements have been shown in the real world and in more
 286 challenging dynamic environments (`dynamic-wall`). However, `dynamic-wall` is still very toy-
 287 like. Due to the lack of real-world dynamic environments, we cannot justify whether the proposed
 288 dynamic environments in the benchmark is instructive to the real-world deployment or not. To
 289 further evaluate the ability of RL-based systems with memory to reason under uncertainty, more
 290 realistic but also challenging dynamic environments will be needed in the future.

291 **(D2) safety** is improved by both safe RL and model-based MPC methods. However, classical navigation
 292 systems still achieve the best safety performance at a cost of very long traversal time. Whether
 293 RL-based navigation systems can achieve similar level of safety guarantee as classical navigation
 294 systems and whether safety can be improved without significantly sacrificing the traversal time are
 295 still open questions.

296 **(D3) the ability to learn from limited trial-and-error data** is not improved by the proposed model-
 297 based methods. However, our implementation is based on an off-policy TD3 algorithm which already
 298 ensures a relatively good sample efficiency. Since we did not include on-policy RL algorithms
 299 in this study, it is still unclear whether model-based methods combined with on-policy RL algo-
 300 rithms can achieve better sample efficiency or not.

301 **(D4) the generalization to diverse and novel environments** is improved by increasing the ran-
 302 domness of training environments. However, a noticeable gap of about 5% between training and
 303 test environments can not be eliminated by further increasing the randomness. While domain ran-
 304 domization is a simple but efficient method, open-ended curriculum learning [27] that leverages
 305 procedural content generation is reported to further improve the zero-shot generalization.

306 **(D5) Navigation is not solved yet!** Although the proposed benchmark cannot represent every real-
 307 world navigation scenario, we believe it serves as a simple yet comprehensive testbed for RL-based
 308 navigation methods. Because we observed that for every desideratum, no methods can achieve 100%
 309 success rate on all *training* environments. Moreover, we have made sure that every environment is
 310 indeed individually solvable even using the naive TD3 method with a simple MLP architecture. This
 311 alone indicates that there exists an optimization and generalization challenge when we have a large
 312 number of training environments as in our proposed benchmark.

313 **References**

- 314 [1] S. Quinlan and O. Khatib. Elastic bands: Connecting path planning and control. In *[1993] Pro-*
315 *ceedings IEEE International Conference on Robotics and Automation*, pages 802–807. IEEE,
316 1993.
- 317 [2] D. Fox, W. Burgard, and S. Thrun. The dynamic window approach to collision avoidance.
318 *IEEE Robotics & Automation Magazine*, 4(1):23–33, 1997.
- 319 [3] Icra 2022 benchmark autonomous robot navigation (barn) challenge. [https://www.cs.](https://www.cs.utexas.edu/~xiao/BARN_Challenge/BARN_Challenge.html)
320 [utexas.edu/~xiao/BARN_Challenge/BARN_Challenge.html](https://www.cs.utexas.edu/~xiao/BARN_Challenge/BARN_Challenge.html), 2022. Accessed: 2022-
321 06-09.
- 322 [4] X. Xiao, B. Liu, G. Warnell, and P. Stone. Motion planning and control for mobile robot
323 navigation using machine learning: a survey. *Autonomous Robots*, pages 1–29, 2022.
- 324 [5] Y. Chow, O. Nachum, A. Faust, M. Ghavamzadeh, and E. A. Duéñez-Guzmán. Lyapunov-
325 based safe policy optimization for continuous control. *CoRR*, abs/1901.10031, 2019. URL
326 <http://arxiv.org/abs/1901.10031>.
- 327 [6] G. Thomas, Y. Luo, and T. Ma. Safe reinforcement learning by imagining the near future,
328 2022.
- 329 [7] E. Rodríguez-Seda, D. Stipanovic, and M. Spong. Lyapunov-based cooperative avoidance
330 control for multiple lagrangian systems with bounded sensing uncertainties. In *2011 50th IEEE*
331 *Conference on Decision and Control and European Control Conference, CDC-ECC 2011*,
332 Proceedings of the IEEE Conference on Decision and Control, pages 4207–4213, Dec. 2011.
333 ISBN 9781612848006. doi:10.1109/CDC.2011.6160783. 2011 50th IEEE Conference on
334 Decision and Control and European Control Conference, CDC-ECC 2011 ; Conference date:
335 12-12-2011 Through 15-12-2011.
- 336 [8] K. Cobbe, O. Klimov, C. Hesse, T. Kim, and J. Schulman. Quantifying generalization in
337 reinforcement learning. In *ICML*, 2019.
- 338 [9] K. Cobbe, C. Hesse, J. Hilton, and J. Schulman. Leveraging procedural generation to bench-
339 mark reinforcement learning. *arXiv preprint arXiv:1912.01588*, 2019.
- 340 [10] N. Justesen, R. R. Torrado, P. Bontrager, A. Khalifa, J. Togelius, and S. Risi. Illuminating gen-
341 eralization in deep reinforcement learning through procedural level generation. *arXiv: Learn-*
342 *ing*, 2018.
- 343 [11] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel. Domain randomization
344 for transferring deep neural networks from simulation to the real world. In *2017 IEEE/RSJ*
345 *international conference on intelligent robots and systems (IROS)*, pages 23–30. IEEE, 2017.
- 346 [12] R. S. Sutton. Dyna, an integrated architecture for learning, planning, and reacting. *SIGART*
347 *Bull.*, 2(4):160–163, jul 1991. ISSN 0163-5719. doi:10.1145/122344.122377. URL <https://doi.org/10.1145/122344.122377>.
- 348 [13] A. Nagabandi, G. Kahn, R. S. Fearing, and S. Levine. Neural network dynamics for model-
349 based deep reinforcement learning with model-free fine-tuning. In *2018 IEEE International*
350 *Conference on Robotics and Automation (ICRA)*, pages 7559–7566. IEEE, 2018.
- 352 [14] M. J. Hausknecht and P. Stone. Deep recurrent q-learning for partially observable mdps. In
353 *AAAI Fall Symposia*, 2015.
- 354 [15] D. Wierstra, A. Förster, J. Peters, and J. Schmidhuber. Solving deep memory pomdps with
355 recurrent policy gradients. In *ICANN*, 2007.

- 356 [16] K. Chua, R. Calandra, R. McAllister, and S. Levine. Deep reinforcement learning in a handful
 357 of trials using probabilistic dynamics models. *Advances in neural information processing*
 358 *systems*, 31, 2018.
- 359 [17] H.-T. L. Chiang, A. Faust, M. Fiser, and A. Francis. Learning navigation behaviors end-to-end
 360 with autorl. *IEEE Robotics and Automation Letters*, 4(2):2007–2014, 2019.
- 361 [18] X. Xiao, B. Liu, G. Warnell, J. Fink, and P. Stone. Appld: Adaptive planner parameter learning
 362 from demonstration. *IEEE Robotics and Automation Letters*, 5(3):4541–4547, 2020.
- 363 [19] Z. Wang, X. Xiao, B. Liu, G. Warnell, and P. Stone. APPLI: Adaptive planner parameter learning-
 364 ing from interventions. In *2021 IEEE International Conference on Robotics and Automation*
 365 (*ICRA*). IEEE, 2021.
- 366 [20] Z. Wang, X. Xiao, G. Warnell, and P. Stone. Apple: Adaptive planner parameter learning from
 367 evaluative feedback. *IEEE Robotics and Automation Letters*, 6(4):7744–7749, 2021.
- 368 [21] Z. Xu, G. Dhamankar, A. Nair, X. Xiao, G. Warnell, B. Liu, Z. Wang, and P. Stone. AP-
 369 PLR: Adaptive planner parameter learning from reinforcement. In *2021 IEEE International*
 370 *Conference on Robotics and Automation (ICRA)*. IEEE, 2021.
- 371 [22] X. Xiao, Z. Wang, Z. Xu, B. Liu, G. Warnell, G. Dhamankar, A. Nair, and P. Stone. Appl:
 372 Adaptive planner parameter learning. *Robotics and Autonomous Systems*, 154:104132, 2022.
- 373 [23] E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. In
 374 *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–
 375 5033, 2012. doi:[10.1109/IROS.2012.6386109](https://doi.org/10.1109/IROS.2012.6386109).
- 376 [24] M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. The arcade learning environment:
 377 An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:
 378 253–279, jun 2013.
- 379 [25] H.-T. L. Chiang, A. Faust, M. Fiser, and A. Francis. Learning navigation behaviors end-to-end
 380 with autorl. *IEEE Robotics and Automation Letters*, 4:2007–2014, 2019.
- 381 [26] A. Wahid, A. Toshev, M. Fiser, and T.-W. E. Lee. Long range neural navigation policies for
 382 the real world. *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems*
 383 (*IROS*), pages 82–89, 2019.
- 384 [27] M. Dennis, N. Jaques, E. Vinitsky, A. Bayen, S. Russell, A. Critch, and S. Levine. Emergent
 385 complexity and zero-shot transfer via unsupervised environment design. *Advances in Neural*
 386 *Information Processing Systems*, 33:13049–13061, 2020.
- 387 [28] R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- 388 [29] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–
 389 1780, 1997.
- 390 [30] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio. Empirical evaluation of gated recurrent neural
 391 networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.
- 392 [31] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polo-
 393 sukhin. Attention is all you need. *Advances in neural information processing systems*, 30,
 394 2017.
- 395 [32] E. Altman. *Constrained Markov decision processes: stochastic modeling*. Routledge, 1999.
- 396 [33] S. Boyd, S. P. Boyd, and L. Vandenberghe. *Convex optimization*. Cambridge university press,
 397 2004.

- 398 [34] C. G. Atkeson and J. C. Santamaria. A comparison of direct and model-based reinforcement
399 learning. In *Proceedings of international conference on robotics and automation*, volume 4,
400 pages 3557–3564. IEEE, 1997.
- 401 [35] R. S. Sutton. Integrated architectures for learning, planning, and reacting based on approximat-
402 ing dynamic programming. In *Machine learning proceedings 1990*, pages 216–224. Elsevier,
403 1990.
- 404 [36] J. A. Rossiter. *Model-based predictive control: a practical approach*. CRC press, 2017.
- 405 [37] F. Sadeghi and S. Levine. CAD²RL: Real single-image flight without a single real image.
406 *ArXiv*, abs/1611.04201, 2017.
- 407 [38] X. B. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel. Sim-to-real transfer of robotic
408 control with dynamics randomization. *2018 IEEE International Conference on Robotics and*
409 *Automation (ICRA)*, pages 1–8, 2018.
- 410 [39] D. Perille, A. Truong, X. Xiao, and P. Stone. Benchmarking metric ground navigation. In
411 *2020 IEEE International Symposium on Safety, Security, and Rescue Robotics (SSRR)*, pages
412 116–121. IEEE, 2020.
- 413 [40] S. Wolfram. Statistical mechanics of cellular automata. *Reviews of modern physics*, 55(3):601,
414 1983.
- 415 [41] OSRF. Ros wiki move_base. http://wiki.ros.org/move_base, 2018.

416 **A Studied Techniques**

417 In Sec. A.1, a brief review of Reinforcement Learning (RL) and Markov decision processes (MDPs)
 418 is provided. Then, Sec. A.2 to A.5 describes in detail the studied techniques and how they can
 419 potentially achieve the desiderata.

420 **A.1 RL and MDPs**

421 In RL, an agent optimizes its discounted cumulative return through interactions with an environment,
 422 which is formulated as an MDP. Specifically, an MDP is a 5-tuple (S, A, T, γ, R) , where S, A are
 423 the state and action spaces, $T : S \times A \rightarrow S$ is the transition kernel that maps the agent's current
 424 state and its action to the next state, γ is a discount factor and $R : S \times A \rightarrow \mathbb{R}$ is the reward function.
 425 The overall objective is for the agent to find a policy function $\pi : S \rightarrow A$ such that its discounted
 426 cumulative return is maximized: $\pi^* = \arg \max_{\pi} \mathbb{E}_{s_t, a_t \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \right]$ [28].

427 MDPs assume the agent has access to the world state s which encapsulates sufficient information
 428 for making optimal decisions. However, in real applications, the agent often only perceives part of
 429 the state s at any moment. Such partial observability leads to uncertainty in the world state and
 430 the problem becomes a Partially Observable Markov decision process (POMDP). A POMDP is a
 431 7-tuple $(S, A, O, T, \gamma, R, Z)$. In addition to the elements of an MDP, O denotes the observation
 432 space and $Z : S \rightarrow O$ is an observation model that maps the world state to an observation. For
 433 instance, at each time step t , the agent receives an observation $o_t \sim Z(\cdot | s_t)$. In general, solving a
 434 POMDP optimally requires taking the entire history into consideration, which means the objective
 435 is then to find a policy that maps its past trajectory $\tau_t = (o_0, a_0, \dots, o_t, a_t)$ to an action a_t such that
 436 $\max_{\pi} \mathbb{E}_{\tau_t \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \right]$.

437 **A.2 Memory-based Neural Network Architectures (D1)**

438 Due to the uncertainty from partial observations (e.g., caused by dynamic obstacles or imperfect
 439 sensory inputs), a mobile agent often needs to aggregate the information along its trajectory history
 440 (as in many classical systems, a costmap is being continuously built based on previous perceptions).
 441 When using a parameterized model as the policy, recurrent neural networks (RNNs) such as those
 442 incorporating Long-Short Term Memory (LSTM) [29] or Gated Recurrent Units (GRUs) [30] are
 443 widely adopted for solving POMDPs [14, 15]. More recently, transformers [31], a type of deep
 444 neural architecture that use multiple layers of attention mechanism to process sequence data, have
 445 been proposed. Transformers have demonstrated superior performance over RNN-based models in
 446 vision and natural language applications [31]. In this work, we consider both GRU and transformers
 447 as the backbone model for the navigation policy. The reason for choosing GRU over LSTM is
 448 due to the fact that GRU has a simpler architecture than, but performs comparably with, LSTM in
 449 practice [30].

450 **A.3 Safe RL (D2)**

451 Successful navigation involves both navigating efficiently towards a user specified goal and avoiding
 452 collisions. While most prior RL-based navigation approaches design a *single* reward function that
 453 summarizes both objectives, it is not clear whether explicitly treating the two objectives separately
 454 will have any benefits. Specifically, assume the reward function R only rewards the agent for making
 455 progress to the goal. Additionally, a cost function $C : S \times A \rightarrow \mathbb{R}^+$ maps a state s and the
 456 agent's action a to a penalty c . Then the navigation problem can be transformed into a constrained
 457 optimization problem with 2 objectives [32]:

$$\max_{\pi} \mathbb{E}_{s_t, a_t \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \right] \quad \text{s.t.} \quad \mathbb{E}_{s_t, a_t \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t C(s_t, a_t) \right] \leq \epsilon. \quad (2)$$

458 Here $\epsilon \geq 0$ is a threshold that controls how tolerant we are of the risk of collision. By formulating
 459 the navigation problem in this way, existing constrained optimization techniques can be applied for
 460 solving (2). One of the most common approaches in constrained optimization is using a Lagrangian

461 multiplier, which transforms the constraint into a penalty term multiplied by a Lagrangian multiplier
 462 $\lambda \geq 0$. Specifically, the objective becomes

$$\max_{\pi} \mathbb{E}_{s_t, a_t \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \right] + \lambda \left(\mathbb{E}_{s_t, a_t \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t C(s_t, a_t) \right] - \epsilon \right). \quad (3)$$

463 To solve (3), one can either manually specify a λ based on prior knowledge or optimize λ simultaneously [33]. In this work, we perform a grid search over λ but fix it during learning.

465 A.4 Model-based RL (D2, D3)

466 Given an accurate model, model-based RL often benefits from better sample efficiency compared to
 467 model-free methods [34]. In addition to improved sample efficiency, prior work on model-based RL
 468 also reported that planning with a learned model can improve the agent's safety [6]. Therefore, we
 469 investigate whether these claims also hold for autonomous navigation if the agent learns a transition
 470 model \hat{T} from its interaction with the world. Although the reward model R can also be learned, as the
 471 structure of the reward function is usually designed manually, we let the agent take advantage of the
 472 known reward function. In this work, we consider two ways of using a learned model: a Dyna-style
 473 method [35] and model-predictive control (MPC) [36]. Specifically, assume the agent's rollout tra-
 474 jectories are saved into a replay buffer B in the form of transition tuples: $B = \{(s_t, a_t, s_{t+1}, r_t)\}_t$.¹
 475 Then model-based methods assume the agent also learns a model $\hat{T} : S \times A \rightarrow S$ that approximates
 476 T . A common learning objective of \hat{T} is to minimize the mean-square-error between \hat{T} 's and T 's
 477 predictions on transitions:

$$\hat{T} = \arg \min_{T'} \mathbb{E}_{(s, a, s', r) \sim B} \|T'(s, a) - s'\|_2^2. \quad (4)$$

478 Once \hat{T} is learned, given a transition pair $(s, a) \sim B$, the Dyna-style method samples additional
 479 $s' \sim \hat{T}(s, a)$ to enrich the replay buffer that can potentially benefit the learning of the value function.
 480 MPC, on the other hand, first uses \hat{T} to form samples of future trajectories, then it outputs the first
 481 action corresponding to the trajectory that has the highest return. MPC usually enables a more
 482 efficient exploration by selecting promising actions, and also potentially improves the asymptotic
 483 performance with the help of model predictions.

484 A.5 Domain randomization (D4)

485 A direct deployment of a policy trained in limited training environments to unseen target environ-
 486 ments is usually formalized as a zero-shot transfer problem, where no extra training of the policy is
 487 allowed in the target environments. One promising approach for zero-shot transfer has been Domain
 488 Randomization (DR) [11]. In DR, the environment parameters (i.e. obstacle configurations) in pre-
 489 defined ranges are randomly selected in each training environment. By randomizing everything that
 490 might vary in the target environments, the generalization can be improved by covering target envi-
 491 ronments as variations of random training environments. This simple yet strategy has been reported
 492 to be efficient in practice when solving many sim-to-real transfer problems [11, 37, 38].

¹Replay buffer is a commonly used technique to learn the value function (critic) in reinforcement learning.

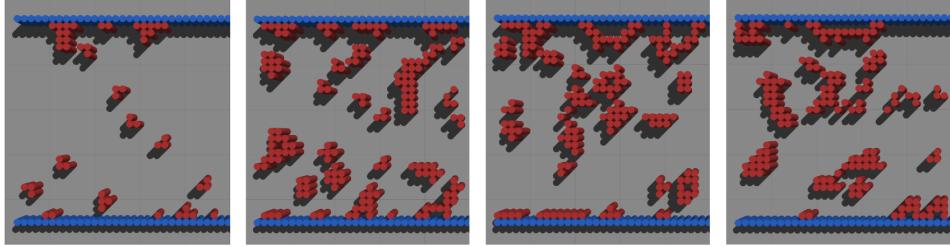


Figure 5: Examples of static environments.

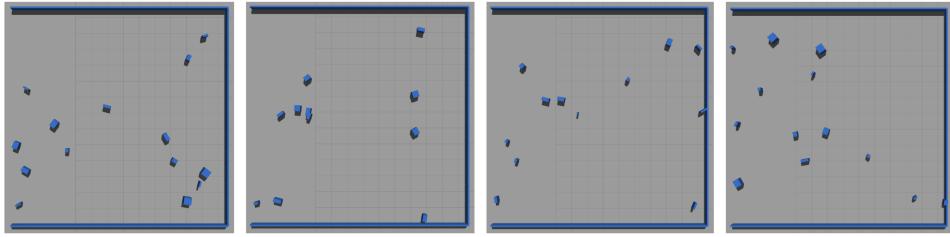


Figure 6: Examples of dynamic-box environments.

493 B Navigation environments

494 **Static environments.** We use the 300 static environments from the BARN dataset [39]. The obstacle
 495 fields in these environments are represented by a 30×30 black-white grid, which corresponds to
 496 an area of $4.5\text{m} \times 4.5\text{m}$ with black and white cells representing obstacle-occupied and free space
 497 respectively (see Fig. 5). The grid cells are generated by a method of cellular automation [40],
 498 which is originally designed to generate a collection of black cells on a white grid of specified shape
 499 that evolves through a number of discrete time steps according to a set of rules based on the states of
 500 neighboring cells. In BARN, such generation process begins with randomly filling the grid cells by
 501 a ratio of *initial fill percentage*, then performs *smoothing iterations* that either fill an empty cell if the
 502 number of its filled neighbors is larger than *fill threshold* or empty a filled cell if its filled neighbors
 503 is smaller than *clear threshold*. The grids that are not navigable will be discarded. Due to the cellular
 504 automation, the resulting grid resembles real-world obstacles more than the initial randomly filled
 505 grid does. To generate BARN dataset, 12 different sets of hyper-parameters are used with *initial*
 506 *fill percentage* chosen from $\{0.15, 0.2, 0.25, 0.3\}$ and *smoothing iterations* ranges from 2 to 4. The
 507 *fill threshold* and *clear threshold* are kept at 5 and 1 respectively. Each set of hyper-parameters
 508 generates 25 environments which constitute 300 environments in total.

509 **Dynamic box environments.** These environments are $13.5\text{m} \times 13.5\text{m}$ obstacle fields (larger than
 510 the static environments) that give the agent more time to respond to the moving obstacles (see Fig.
 511 6). The obstacles are randomly generated without any manually designed challenging scenarios.
 512 Each obstacle is a $w \times l \times h$ box with its width w and length l randomly sampled from a range of
 513 $[0.1\text{m}, 0.5\text{m}]$ and a height $h = 1\text{m}$. The obstacles start from a random position on the left edge of
 514 the obstacle field with a random orientation and a constant linear velocity. The magnitude of the
 515 velocity is randomly sampled from a range of $[1\text{m/s}, 1.5\text{m/s}]$, and the direction of the velocity is
 516 randomly sampled from all the possible directions pointing into the obstacle field. Each obstacle
 517 repeats its motion once it moves out of the obstacle field. Each dynamic box environment has 10 to
 518 15 such randomly generated obstacles.

519 **Dynamic wall environments.** These environments are $4.5\text{m} \times 4.5\text{m}$, which have two long parallel
 520 walls moving in opposite directions with their velocities perpendicular to the start-goal direction
 521 (see Fig. 7). The walls are long enough so that the robot can only pass when the two walls are
 522 moving apart. This manually designed navigation scenario requires the agent to maintain a memory
 523 of past observations and actions, and to estimate the motion of obstacles. To challenge the agent,
 524 we add small variances so that each wall's length, tilting angle, and magnitude of the velocity are
 525 randomly sampled from the ranges of $[3.5\text{m}, 4.5\text{m}]$, $[-10^\circ, 10^\circ]$ and $[1\text{m/s}, 1.4\text{m/s}]$ respectively.

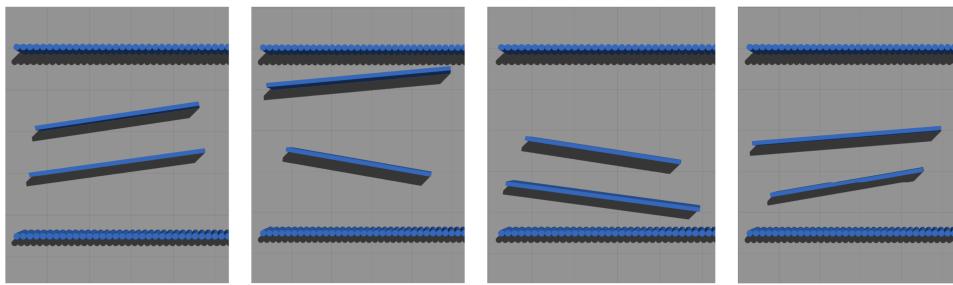


Figure 7: Examples of dynamic-wall environments.

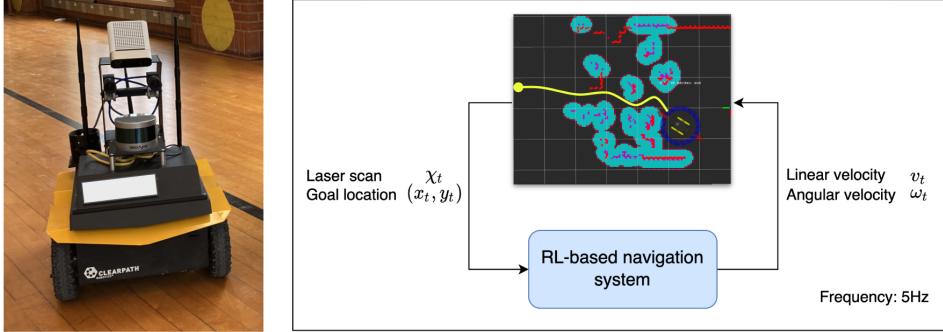


Figure 8: The picture of a Jackal Robot (left) and the diagram of the RL-based navigation system.

526 C Jackal Robot and Navigation system

527 As shown in Fig. 8 on the left, the navigation is performed by a ClearPath Jackal differential-drive
 528 ground robot in navigation environments simulated by the Gazebo simulator. The robot is equipped
 529 with a 720-dimensional planar laser scan with a 270° field of view, which is used as our sensory
 530 input χ_t . We preprocess the LiDAR scans by capping the maximum range to 5m which covers
 531 the entire obstacle field. The goal location (\bar{x}_t, \bar{y}_t) is queried directly from the Gazebo simulator.
 532 The robot receives velocity commands at a frequency of 5 Hz. A time limit of 80s is used which
 533 corresponds to a maximum of 400 time steps in an episode. In the real world, the robot is deployed
 534 in the same way as simulation except for that the goal location is queried from the on-board ROS
 535 move_base stack [41] that localizes the robot simultaneously. Fig. 8 on the right shows a diagram
 536 of the RL-based navigation system.