# Bridging the Gap: Regularized Reinforcement Learning for Improved Classical Motion Planning with Safety Modules

Elias Goldsztejn[1], Ronen I. Brafman[2]

*Abstract*— Classical navigation planners can provide safe navigation, albeit often suboptimally and with hindered human norm compliance. ML-based, contemporary autonomous navigation algorithms can imitate more natural and human-compliant navigation, but usually require large and realistic datasets and do not always provide safety guarantees. We present an approach that leverages a classical algorithm to guide reinforcement learning. This greatly improves the results and convergence rate of the underlying RL algorithm and requires no human-expert demonstrations to jump-start the process. Additionally, we incorporate a practical fallback system that can switch back to a classical planner to ensure safety. The outcome is a sample efficient ML approach for mobile navigation that builds on classical algorithms, improves them to ensure human compliance, and guarantees safety.

Fig. 1: A robot navigating autonomously using a system that integrates an RL-based policy, regularized with a classical planner, alongside a safety switching mechanism.

## I. INTRODUCTION

The proliferation of self-driving cars by an increasing number of companies and the ability of robots to efficiently deliver food and supplies are clear manifestations of the vast improvements in autonomous navigation. Nonetheless, various unresolved challenges persist across diverse domains.

Classical planners use analytic optimization techniques and reactive rules for collision avoidance and for finding safe paths [1]–[3]. These methods can be successful in specific domains, require no or little data, are well understood and can lead to safe and interpretable behavior. However, they often exhibit unnatural and inefficient behaviors and poor social norm compliance [4].

Machine learning methodologies build on the latest advancements in imitation learning (e.g., [5], [6]) and deep reinforcement learning (e.g., [7]–[11]) through which they can capture the intricacy of human actions and provide enhanced environmental awareness and human-aware behavior. However, they require large and realistic datasets, and often lack safety guarantees [12]. Furthermore, because they are often based on end-to-end learning, they lack interpretability and transparency [13]. For these reasons, they can fail badly in unexpected ways on particular inputs [14], [15], making it difficult to rely on them.

This work seeks to alleviate the shortcomings of classical and learning-based methods by combining suitable components of each, building on and modifying various existing methodologies, reviewed in the next section. In particular, we exploit classical algorithms to improve the sample efficiency of a learning algorithm and the performance and safety of its resulting navigation policy.

Our main contribution is a sample-efficient learning strategy for improving classical planners and a fallback system with a trained supervisor that guarantees safety. More specifically, we suggest the following approach:

1) Train a planner using DRL with policy guidance derived from a classical planner: We seed the replay buffer with experiences generated by the classical planner and regularize the actor in an actor-critic algorithm using the classical planner's policy.
2) Use a classical rule-based navigation policy as a fallback system and train a supervisor that performs minimal switching between the neural and classical planner to ensure safety.

While our focus is navigation planning, the above offers a general recipe for using learning to improve classical algorithms while retaining their respective benefits. Our approach provides safety and offers transparency at the supervisor level. Its reliance on good, existing classical algorithms helps jump-start the learning algorithm, leading to faster and better convergence, as our empirical evaluation clearly demonstrates. The regularization term further stabilizes the learning process and ensures greater transparency, forcing it to remain in the vicinity of the well-understood classical algorithm. And unlike methods that rely on human demonstrations to achieve some of these effects, no human involvement is needed – as it is already implicitly present in the formulation of the classical algorithm and the reward function.

Relevant code and videos for this paper are available at `https://github.com/eliasgoldsztejn95`

## II. RELATED WORK

We review previous work on classical planners for autonomous driving in mobile robots, reinforcement learning,

[1]Elias Goldsztejn and [2]Ronen Brafman are with the Department of Computer Science at Ben-Gurion University of the Negev. eliasgol@post.bgu.ac.il, brafman@bgu.ac.il
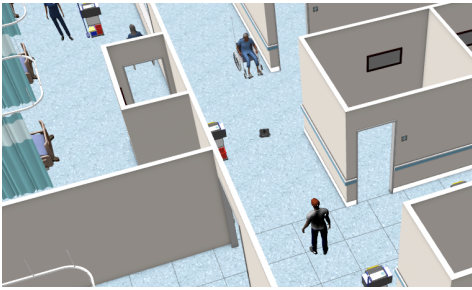
Fig. 2: All Learning is carried out in this realistic simulated hospital with people moving according to social forces.

especially using learning from demonstrations, as well as techniques for safe planning.

### A. Classical Planning for Mobile Robots

Early planning algorithms for mobile robots used reactive rules and analytic optimization to avoid collision and reach targets. Important approaches include Gaussian processes [16], decentralized and scalable multi-agent reciprocal optimizations [2], reactive techniques [1], spatial and temporal constraint-aware elastic-band methods [3], and more. Classical planners serve as foundational benchmarks in mobile navigation and are commonly integrated into critical robotic systems, exemplified by their inclusion in the widely adopted ROS [17] navigation stack. Although in some environments, these methods can provide effective navigation, they usually require fine-tuning from experts and are limited by their lack of human compliance and social norms, non-smooth behavior, and freezing robot problems, among others.

### B. Imitation and Reinforcement Learning

Imitation learning (IL) and reinforcement learning (RL) are effective and popular techniques for autonomous driving that can create flexible, human-compliant navigation systems. IL uses supervised learning to map states to actions so as to mimic the behavior of a demonstrator, such as a human expert. Some pioneering works include using synthesized data in the form of perturbations [18], and applying command-conditions to represent expert's intentions [19]. IL has well-known limitations, such as dataset bias and overfitting, new generalization issues, covariate shifts, and requires large amounts of data [20]. Mitigating these effects is an ongoing area of interest. Proposed methods include utilizing space-time cost volumes [21], which enable interpretability and faster learning. Others include combining reinforcement learning to address the covariate shift problem [22].

In RL, agents learn to make decisions by interacting with an environment. Feedback is obtained in the form of rewards from state-action pairs, and the primary objective is to derive a policy that maximizes the cumulative reward over time. In contrast with IL, RL has a closed-loop exploration-exploitation approach, which addresses the stability and overfitting challenges of IL.

RL is used extensively for autonomous driving and navigation [7]–[11], lane following and urban driving [23],

[24]. While RL is a promising approach it faces numerous challenges, including the difficulty of defining a suitable reward function, the need for extensive training datasets, and concerns related to safety and interpretability [25], [26]. Techniques like inverse and safe RL and RL with demonstrations and human feedback have been developed aimed at addressing these challenges more effectively. However, in many domains, solutions for these issues are still elusive.

### C. Reinforcement Learning with Demonstrations

RL with demonstrations combines elements of IL into the RL procedure offering complementary strengths. IL enhances realism, and alleviates the challenge of reward design and extensive datasets, while RL enhances safety and robustness. IL can be used for the initialization of policies [27], [28] and expert demonstrations can be included in the experience replay buffer [29], [30] to sample both expert demonstrations and interactions. However, these methods can fail when randomness of exploration is needed at the beginning of training or when expert rewards are not accessible.

We are specifically interested in approaches that directly incorporate expert demonstrations during the training stage, explicitly integrating them into the learning process (e.g. [31]–[34]). This approach is especially useful when we seek a learned policy similar to that of the expert.

Particularly relevant is the work of [33], which uses an RL strategy that incorporates the KL divergence between imitative expert priors and the agent policy in the reward function. This approach regularizes the agent's behavior with the expert policy while maintaining its exploration ability.

We propose to exploit the same strategy but use classical planners as the "expert" prior policy. As explained before, classical planners for mobile robots are competent enough to be employed as baselines or, in this case, as expert priors. The added benefit of using available planners is that the collection of data-sets containing human-expert demonstrations is no longer required!

### D. Safe Planning

Classical planners for mobile robots prioritize safety, and use well-understood and, therefore, trustworthy safety mechanisms, but often behave sub-optimally. While RL algorithms often provide enhanced efficiency, complete safety assurance remains challenging, especially when using Deep RL.

We focus on rule-based systems with fallback layers (e.g., [35], [36]) that assess ML planner outputs against predefined checks. If a plan is deemed unsafe, they can modify it, opt for alternative policies, or employ other strategies. These methods ensure predictable and reliable behavior through explicit rules. We employ a similar approach: In critical scenarios, the policy is seamlessly switched from the RL policy to the classical planner. The switching module is trained to guarantee safety while minimizing transitions between the ML policy and the classical planner. This allows for the utilization of an optimal, human-compliant planning system when safety is assured while seamlessly transitioning to a practical system that prioritizes safety in critical settings.
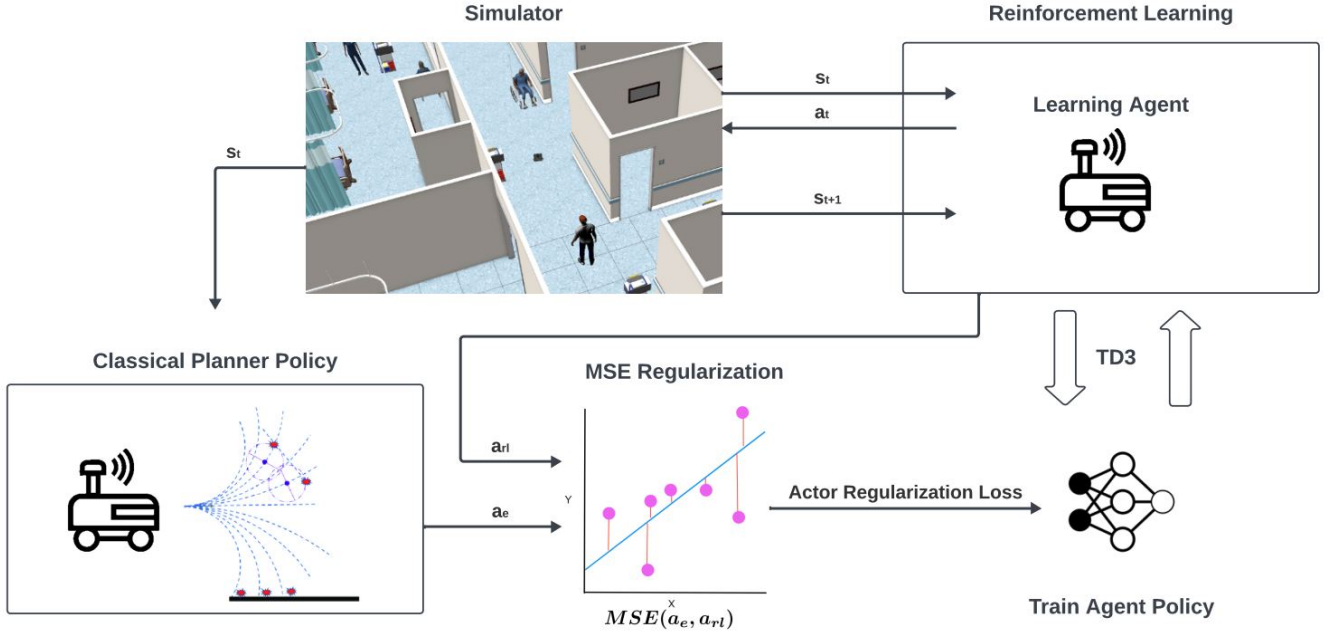
Fig. 3: Framework of the learning strategy with a classical planner regularization.

Unlike [37], which trains the RL algorithm and the switching module in unison, We first train an RL agent, on top of which we train a fuzzy-control supervisor. This makes learning easier, prevents the RL policy from exploiting the fallback algorithm too much, and makes for more transparent behavior and hence, better safety guarantees.

## III. BACKGROUND

### A. Reinforcement Learning

Reinforcement learning (RL) addresses the challenge of mastering the control of a dynamic system, characterized by a Markov decision process (MDP) denoted as $M = (S, A, Tr, r, \gamma)$, where $S$ is the state space, $A$ is the action space, $Tr : S \times A \rightarrow \Pi(S)$ is the transition function, $r : S \times A \rightarrow \mathbb{R}$ is the reward function, and $\gamma \in (0, 1]$ serves as a discount factor. The objective in RL is to develop a policy, described as a distribution over actions conditioned on states $\pi(a_t|s_t)$, aiming to maximize the long-term discounted cumulative reward:

$$\max_{\pi} \mathbb{E}_{\tau \sim p_\pi(\tau)} \left[ \sum_{t=0}^{T} \gamma^t r(s_t, a_t) \right] \tag{1}$$

Here, $\tau$ denotes a trajectory, $p_\pi(\tau)$ is the distribution of the trajectory under policy $\pi$, and $T$ represents the time horizon.

The Actor-Critic (AC) method is a prominent paradigm within RL, combining an actor (policy) and a critic (value function) for improved decision-making and evaluation:

$$\text{Actor: } \pi(a_t|s_t; \theta^\pi) \quad \text{Critic: } Q(s_t, a_t; \theta^Q)$$

Here, $\pi$ represents the policy function with parameters $\theta^\pi$, and $Q$ denotes the action-value function with parameters $\theta^Q$.
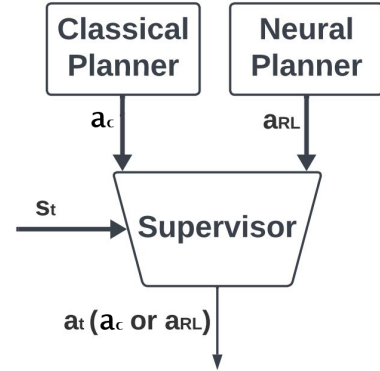


Fig. 4: A supervisor selects whether to switch from the RL to the Classical planner in critical situations.

The actor aims to determine the optimal policy, while the critic evaluates the chosen actions.

Drawing inspiration from [33], we impose a constraint ensuring similarity between the actor and the expert. This has two advantages: (1) A simplified exploration space concentrated around the expert policy, leading to decreased sampling complexity. (2) The policy mimics the expert, reducing the need for a finely tuned reward function. To implement this, we introduce a regularization term in the actor loss update, defining the *regularized actor loss* as:

$$\mathcal{L}_{\text{actor}}(\theta^\pi) = \mathbb{E}_{s_t \sim \mathcal{D}} \bigg[ -Q(s_t, \pi(s_t; \theta^\pi); \theta^Q)$$
$$+ \lambda \cdot \text{MSE}(\pi(s_t; \theta^\pi), \pi_{\text{expert}}(s_t)) \bigg]$$

Here, $\mathcal{L}_{\text{actor}}(\theta^\pi)$ is the actor loss. $Q(s_t, a_t; \theta^Q)$ is the

estimated action value by the critic. $\pi_{\text{expert}}(s_t)$ represents expert actions. $\lambda$ is the regularization strength. $\text{MSE}(\cdot, \cdot)$ is the mean squared error.

This formulation encourages the actor to generate actions that maximize expected rewards while regularizing them towards expert actions through the MSE term. This is particularly useful given sparse reward functions since the expert guides the learned policy. In our work, the "expert" is the classical planner. Unlike other approaches, such as [34], which rely on a specific classical planner to guide the search, our method can seamlessly integrate with any expert.

### B. "Expert" Prior

Our method assumes access to a good analytic algorithm for the task at hand. In robotics, the ROS platform supplies multiple such algorithms. In particular, for navigation, we employ the *move_base* framework alongside the classical local planner DWA [1], from which we can derive actions at each time step that we use to guide an AC algorithm. The expert actions are remapped to an unused ROS *topic*, while the RL actions are mapped to control the robot.

Unlike approaches relying on behavioral cloning, as seen in [22], [33], our method eliminates the need for a preceding learning step to imitate an expert. This makes our approach more cost-effective and less time-consuming.

### C. Safety Protocol

Drawing inspiration from various hybrid navigation approaches such as [26], [35]–[37], we formulate a practical protocol that can switch between a learning-based planner and a classical planner. Following [35], [36], we design this protocol with a focus on ensuring safety in critical scenarios, recognizing that, in many instances, reliability and safety take precedence over swift navigation.

Similar to the approach of [36], we employ deterministic rules to determine when a switch to a safer navigation mode is imperative. However, we introduce a strategy to fine-tune these rules, aiming to optimize safety while minimizing the frequency of transitions to the safe planner. More specifically, we utilize fuzzy rules and refine the membership functions through a genetic algorithm with the dual objectives of minimizing both transitions and critical situations.

### IV. ARCHITECTURE

The high-level architecture of the navigation algorithm is described in Fig. 4: A supervisor module dynamically switches between a neural-net based policy and a classical planner, effectively maintaining a balanced overall policy that maintains safety, robustness and efficiency in real-world scenarios. We now describe its components in more detail. Section V describes the methods used to train them.

*a) Supervisor:* Generally, machine learning techniques perform well in navigation, but may fail badly in certain scenarios. To address this, we adopt a transparent and interpretable strategy that incorporates a rule-based safety mechanism directly during deployment. First, like many local planners, if the distance of the robot to the nearest obstacle
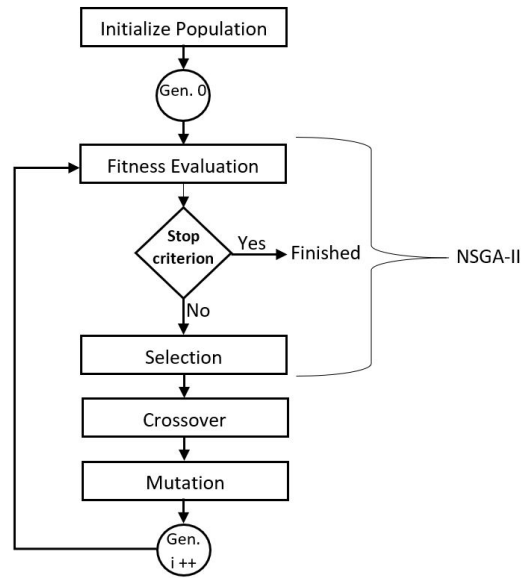


Fig. 5: Evolutionary Training Diagram for Supervisor. **Initialization:** Generate diverse populations with varying fuzzy membership values. **Fitness Evaluation:** Assess fitness through robot episodes, tracking critical situations and supervisor activations. **Selection:** Select values minimizing supervisor activations and critical situation occurrences.

$< 30cm$ then it moves backward. In addition, our *Supervisor* module switches to a default safe policy if the distance to the nearest obstacle is below a parameter called *radius* below, which is determined using the following simple fuzzy rules:

---
**Algorithm 1** Fuzzy Rules for Supervisor Radius
---
**if** robot velocity is high **then**
    supervisor radius is big
**else if** robot velocity is low **then**
    supervisor radius is small
**end if**

---

*b) Neural Navigation Policy:* The neural network receives as input an egocentric 2D bird-eye view grayscale image of the local costmap (6x6 Meters). This image is created with a lidar sensor mounted on the robot. We also incorporate the current velocity of the robot and the location of the next waypoint (which is 2 meters away) as indicated by the global planner (a modified version of A*). While its output is the linear and angular velocity command to the robot (continuous action space).

We use a deep convolutional NN followed by fully connected layers for the actor and critic networks (See Fig. 6). We process the image and concatenate it with velocity and waypoint vectors and, in the case of the critic, the action.

*c) Safe Policy:* The policy used when switching is the pure pursuit algorithm [38] with a very low, fixed linear speed ($0.5m/s$) towards the next waypoint, $30cm$ away. This straightforward algorithm is safe because of the following: (1) The $A^*$ variation used in this work creates a collision-
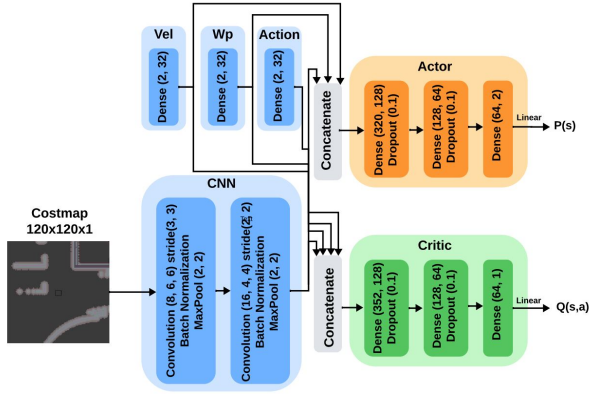
Fig. 6: Actor Critic neural networks. The Costmap is encoded using CNN's. The velocity (vel), waypoint, and action are concatenated. Two fully connected neural networks output the next action (actor) and the Q function (critic).

free path given static obstacles. (2) The robot's low velocity affords ample time for individuals (dynamic obstacles) to adapt to its presence, and for $A^*$ to adapt to these obstacles.

## V. METHODS

We train the neural policy and the genetic algorithm in simulation. We use a simulation platform of a realistic hospital [39] integrated with humans moving according to Social Forces, as described in [40]. A computer with AMD® Ryzen 9 5900x12 CPU and GeForce RTX 4070 was used for simulation and training. We used Pytorch [41] for the implementation of the neural networks.

The robots used in simulation and real life are differential drive robots equipped with 2D Lidar sensors with a $360^0$ field of view. With this sensor, a 6x6-meter egocentric local costmap is generated. We use *move_base* to generate the global plan (which uses a variant of A*) and, from there extract intermediate waypoints, as well as the DWA, recommended actions.

Episodes involved the robot navigating between designated hospital rooms, with simulations sped up about five-fold, completing each episode in 8 to 14 seconds. The training approach, inspired by [26], separated Simulation, Robot, Task environments, and the Training algorithm.

### A. Neural Policy Training

We use RL to develop a policy, training it for $60,000$ episodes (around 6 hours). The training process, illustrated in Fig. 3, refines the Twin Delayed DDPG (TD3) algorithm [42], an online deep RL algorithm, through regularization from an expert, specifically the DWA planner. TD3 uses 2 critics to mitigate overestimation bias and uses delayed updates to stabilize training by averting rapid policy changes and reducing divergence risk. It is known to provide improved sample efficiency compared to other AC algorithms.

Our modified algorithm, Expert-Enhanced TD3 (E2TD3), uses the classical local planner DWA to modify two elements of TD3: its replay buffer initialization and its policy update step. The pseudo-code is shown in Algorithm 2. The changes are highlighted by underlining them, and are explained

below. We note that similar modifications using an expert policy can be applied to many other RL algorithms.

**Buffer Initialization.** TD3 uses a policy-gradient method to update the policy. The gradient is computed by sampling a mini-batch of experiences from the replay buffer. E2TD3 initializes this buffer by recording experiences obtained by executing the expert (DWA) policy with some added normally distributed noise. This data pushes the algorithm in the general direction of the expert policy, while the use of an initially random policy and the noise added to the expert data provide some exploration bias.

**Regularization.** During learning, we augment new experiences $(s, a, r, s')$ with the expert recommended action $a_e$, obtaining the quintuple $(s, a, \underline{a_e}, r, s')$. In the *Actor*'s learning stage, we add to the standard gradient update a regularizer in the form of the MSE between the current policy and the expert policy: $\lambda \cdot \text{MSE}(\pi_\phi(s), \pi_{\text{expert}}(s))$ (while we observed comparable results with alternative values, we opted for $\lambda = 1$). This has two advantages: (1) Simplified exploration space, reducing sampling complexity. (2) The actor mimics the expert policy, reducing the need for a finely tuned reward function.

---

**Algorithm 2** The Expert-Enhanced TD3 (E2TD3) Algorithm

Initialize critic networks $Q_{\theta_1}, Q_{\theta_2}$, and actor network $\pi_\phi$ with random parameters $\theta_1, \theta_2, \phi$.
Initialize target networks: $\theta_1' \leftarrow \theta_1, \theta_2' \leftarrow \theta_2, \phi' \leftarrow \phi$.
Initialize replay buffer $B$ <u>by simulating the expert policy with Gaussian noise: $a = a_e + \epsilon, \ \epsilon \sim \mathcal{N}(0, \sigma)$.</u>
**for** $t = 1$ to $T$ **do**
    Select action with exploration noise: $a \sim \pi_\phi(s) + \epsilon$, $\epsilon \sim \mathcal{N}(0, \sigma)$, and observe <u>expert action $a_e$</u>, reward $r$ and new state $s'$.
    Store transition tuple $(s, a, \underline{a_e}, r, s')$ in $B$.
    Sample mini-batch of $N$ transitions $(s, a, \underline{a_e}, r, s')$ from $B$.
    $\tilde{a} \leftarrow \pi_{\phi'}(s') + \epsilon, \ \epsilon \sim \text{clip}(\mathcal{N}(0, \tilde{\sigma}), -c, c)$.
    $y \leftarrow r + \gamma \min_{i=1,2} Q_{\theta_i'}(s', \tilde{a})$.
    Update critics: $\theta_i \leftarrow \text{argmin}_{\theta_i} \frac{1}{N} \sum_{i=1}^{N} (y - Q_{\theta_i}(s, a))^2$.
    **if** $t \mod d$ **then**
        Update $\phi$ by the deterministic policy gradient:

$$\nabla_\phi J(\phi) = \frac{1}{N} \sum_{i=1}^{N} \nabla_a Q_{\theta_1}(s, a) \big|_{a = \pi_\phi(s)} \nabla_\phi \pi_\phi(s)$$
$$+ \underline{\lambda \cdot \text{MSE}(\pi_\phi(s), \pi_{\text{expert}}(s))}$$

        Update target networks:

$$\theta_i' \leftarrow \tau \theta_i + (1 - \tau)\theta_i'$$
$$\phi' \leftarrow \tau \phi + (1 - \tau)\phi'$$

    **end if**
**end for**

---

**Reward Function.** The RL algorithm attempts to find a policy maximizing the expected sum of discounted rewards (Eq. 1) minus the regularizing term. The reward function
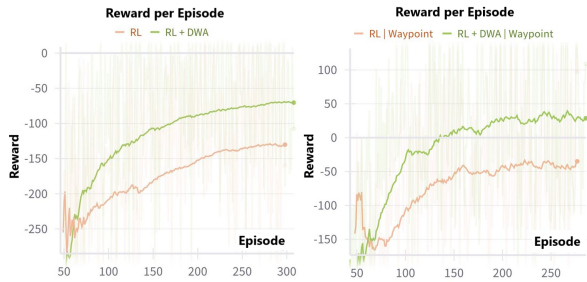
Fig. 7: Total reward per episode in simulations. We compare a plain RL policy with a regularized RL policy using the DWA classical planner, for two reward functions. As can be seen, RL + DWA outperforms the plain RL policy.
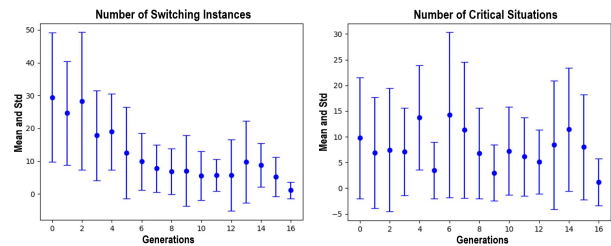


Fig. 8: Mean and STD of number of switching instances and critical situations across generations. We see that the Genetic Algorithm finds parameters that lower these values.

is supposed to quantify desirable properties of the selected path. Our reward function penalizes time and collisions, and rewards progress towards the waypoint. We tested the learning with two variations. In the first one, the reward function is more dense and greedy towards the waypoint: (1) $r(s,a) = r_{timestep} + r_{collision} + r_{progress\,waypoint}$ where $r_{timestep} = -0.5$ is a penalty for each time step, $r_{collision} = -vel - 1$ *if dist* $< 0.5$ is a penalty for the robot being close and fast towards obstacles, and $r_{progress\,waypoint} = |dist(waypoint, postion_t) - dist(waypoint, postion_{t+1})|$ is the distance traveled towards the current waypoint at each time step.

The second reward function is more sparse: the robot is rewarded only when it reaches the waypoint (with a margin of error): (2) $r(s,a) = r_{timestep} + r_{collision} + r_{waypoint}$ where $r_{waypoint} = 10$ is a reward given when the robot reaches each specific waypoint. The total reward during training can be seen in Fig. 7. In both cases our approach shows better performance than plain RL.

*B. Learning Supervisor Parameters*

The *Supervisor* is a ruled-based mechanism that switches to a safe navigation policy as the robot nears obstacles, adjusting its radius dynamically based on velocity, increasing it with higher speeds. We define its parameters using fuzzy rules that assign degrees of membership to variables for more flexible and nuanced decision-making compared to deterministic rules.

We optimize performance by fine-tuning membership function parameters of rule variables using a genetic algorithm. Variables include (1) *Robot velocity* (low/high) and (2) *Supervisor radius* (small/big). Membership functions for these variables are Gaussian, centered around $0m/s$, $1.5m/s$, and $0m$, $1.3m$, with variable standard deviations, learned using the NonDominated Sorting Genetic Algorithm (NSGA-II), Fig. 5. We optimze two objectives: (1) Minimum switching. (2) Minimal critical scenarios.

The Genetic Algorithm runs episodes of the robot moving in the hospital according to the learned RL policy while using the *Supervisor* with a specific set of parameters for the fuzzy rules. For each episode, the fitness function contains the number of alternations and the number of critical scenarios.

A critical scenario is defined as the robot being closer to obstacles than a fixed threshold (0.3 meters).

The supervisor module was trained for 16 generations, each generation with 16 individuals. Individuals were episodes of the robot moving through hospital rooms with the learned RL + DWA policy, augmented with randomized actions. The robot would switch to a low-velocity pure pursuit algorithm as a function of its velocity and proximity to obstacles, determined by the current values of fuzzy rules. The fitness function was evaluated based on the number of switching instances and critical situations (extreme closeness to obstacles) in an episode. As can be seen in Fig. 8, NSGA-II learned to reduce the number of switching instances while also reducing the number of critical situations.

## VI. EXPERIMENTS

In our study, we aim to show the efficacy of our algorithm, Regularized RL with a Supervisor, by demonstrating that it can: (1) be easily implemented with an accessible expert policy (2) improve the expert policy while maintaining proximity to it (3) learn faster, and possibly reach better cumulative rewards than plain RL (4) effectively integrate a supervisor module to diminish unexpected critical situations, thereby rendering it suitable for practical applications.

*A. Simulated Experiments*

We first test and compare DWA, RL, RL + DWA (regularized RL), and RL + DWA + Supervisor (regularized RL with a supervisor) in a simulated realistic hospital with dynamic people moving according to social forces, see Table I. Each algorithm was run for $1e4$ timesteps (around 40 episodes). The metrics for comparison were: Total Reward, Total $r_{collision}$, Timesteps, MSE between DWA and algorithm actions (MSE(dwa,$\pi$)), and the number of critical situations (Critical). A bigger Total Reward and Total $r_{collision}$ are better. Smaller Timesteps, MSE(dwa,$\pi$), and critical situations are preferred.

RL learns a decent navigation policy capable of reaching distinct rooms in the hospital. However, it can get stuck in local minima in complex situations (narrow passages or small rooms). The expert (DWA) is better than it in most metrics.

RL + DWA already shows an improvement (Total Reward, $r_{collision}$) or similar performance (Timesteps) both over RL and DWA. Furthermore, RL + DWA stays closer to DWA

TABLE I: Performance statistics on $1e4$ steps ($\sim 40$ episodes) per each algorithm in simulation. Average and $95\%$ confidence interval for: Total Reward, Total $r_{collision}$, Timesteps, MSE between DWA and algorithm actions (MSE(dwa,$\pi$)), and the number of critical situations (Critical) are presented. The algorithms tested are DWA, plain RL (RL), RL with DWA regularization (RL + DWA), and RL with regularization and supervisor module (RL + DWA + Supervisor).

| Performance Statistics **Simulation**. ((2) reward function) | | | | |
|---|---|---|---|---|
| | DWA | RL | RL + DWA | RL + DWA + Supervisor |
| Total Reward | $M = 15.8 \pm 40.7$ | $M = -3.2 \pm 45.8$ | $M = 53.7 \pm 37.2$ | $\mathbf{M = 64.8 \pm 30.9}$ |
| Total $r_{collision}$ | $M = -175.7 \pm 23.3$ | $M = -94.0 \pm 18.3$ | $\mathbf{M = -83.5 \pm 14.0}$ | $M = -89.5 \pm 18.1$ |
| Timesteps | $M = 225.5 \pm 20.9$ | $M = 296.3 \pm 40.0$ | $M = 226.1 \pm 27.0$ | $\mathbf{M = 221.9 \pm 25.2}$ |
| MSE(dwa, $\pi$) | $M = 0.0\%$ | $M = 55.3\% \pm 12.4\%$ | $M = 23.2\% \pm 3.0\%$ | $\mathbf{M = 20.8\% \pm 1.0\%}$ |
| Critical | $M = 4.8 \pm 3.3\%$ | $M = 20.4\% \pm 10.9\%$ | $M = 11.4\% \pm 6.2\%$ | $\mathbf{M = 3.5\% \pm 0.7\%}$ |



Fig. 9: Rover Zero 3 robot navigating autonomously in different environments.

TABLE II: Performance statistics in 18 real test cases

| Performance Statistics **Real-life** | | |
|---|---|---|
| | DWA | RL + DWA + Superv. |
| Time | $\mathbf{M = 33.3 \pm 7.4}$ | $M = 36.0 \pm 10.2$ |
| Negative Score | $M = 13.6\% \pm 3.7\%$ | $\mathbf{M = 4.7\% \pm 2.2\%}$ |
| Interventions | $M = 1.3 \pm 0.4$ | $\mathbf{M = 0.2 \pm 0.3\%}$ |

($MSE = 23.2\%$) than RL ($MSE = 55.3\%$). However, RL + DWA has a high percentage of critical situations ($11.4\%$), which makes it unreliable.

RL + DWA + Supervisor has the best or similar performance in all metrics. Most importantly, it lowers the number of critical situations ($3.5\%$) to a lower threshold than RL + DWA ($11.4\%$) and even of DWA ($4.5\%$).

This shows the advantage of guiding the search space of machine learning algorithms when a proficient yet suboptimal expert is available. Our algorithm can be straightforwardly implemented in other AC implementations whether in online or offline RL, with different experts. The inclusion of the Supervisor module enhances its practicality and reliability. This is important in real-life applications.

*B. Physical Experiments*

We performed physical experiments using a Roverrobotics Rover Zero robot wth a SLAMTEC RPLidar S1 and an Intel RealSense T265 Tracking Camera. The robot's computer (Intel NUCi7) executed the algorithms, and external control of the robot was established through the Master-Slave ROS protocol from a separate computer. A Dualshock 3 joystick was connected remotely to the robot.

We tested our algorithm, RL + DWA + Supervisor, using the policy learned in simulation with no additional fine-tuning, against DWA in various university areas (corridors, corners, entrances, hall) with static obstacles and people

moving around the robot. The robot navigated autonomously using the *move_base* plugin to generate a global plan with waypoints and DWA or our algorithm to control its actions. At each moment, a person could press the joystick's $X$ button to indicate dissatisfaction and could also control the robot's movement in critical situations (*Intervention*).

We carried 18 trials, 9 for each algorithm, with 3 different people recording Time, Negative Score, and Interventions. The results are shown in Table II. Negative Score is the percentage of time of dissatisfaction with the robot's performance reflecting suboptimal driving situations like proximity to obstacles, slow movement in open areas, or significant deviation from the goal.

While DWA exhibited quicker completion times, our algorithm obtained a lower negative score and fewer interventions. It demonstrated faster performance in open areas, efficiently navigated in clustered spaces, maintained a greater distance from people, passed through doors more effectively, and avoided collisions thanks to the Supervisor module. However, it tends to oscillate more due to policy switching.

## VII. CONCLUSION

This work introduces a framework that uses a well-known, trustworthy algorithm to help regularize the actor within an AC-based deep RL algorithm for the autonomous navigation of mobile robots. It then learns fuzzy rules for optimizing a safety module that switches between the learned policy and a safe but suboptimal one in critical situations. This pragmatic approach is able to improve a classical algorithms through reinforcement learning while preserving practicality and safety. Unlike most imitation learning algorithms, it requires no expert human input.

Our algorithm was able to enhance a prior algorithm (DWA, here) with low training time, adopting some of

the key features of DWA while introducing new elements to address specific scenarios demonstrating superior performance on key measures. Additionally, the inclusion of the Supervisor module helped ensure safety at all times, making it a practical and reliable algorithm.

## REFERENCES

[1] D. Fox, W. Burgard, and S. Thrun, "The dynamic window approach to collision avoidance," *IEEE Robotics & Automation Magazine*, vol. 4, no. 1, pp. 23–33, 1997.

[2] J. van den Berg, S. Guy, M. Lin, and D. Manocha, *Reciprocal n-Body Collision Avoidance*. Springer, Berlin, Heidelberg, 04 2011, vol. 70, pp. 3–19.

[3] C. Rösmann, F. Hoffmann, and T. Bertram, "Timed-elastic-bands for time-optimal point-to-point nonlinear model predictive control," in *ECC*. IEEE, 2015, pp. 3352–3357.

[4] C. I. Mavrogiannis, F. Baldini, A. Wang, D. Zhao, P. Trautman, A. Steinfeld, and J. Oh, "Core challenges of social robot navigation: A survey," *CoRR*, vol. abs/2103.05668, 2021.

[5] Y. Pan, C.-A. Cheng, K. Saigol, K. Lee, X. Yan, E. Theodorou, and B. Boots, "Agile autonomous driving using end-to-end deep imitation learning," *arXiv preprint arXiv:1709.07174*, 2017.

[6] F. Codevilla, M. Müller, A. López, V. Koltun, and A. Dosovitskiy, "End-to-end driving via conditional imitation learning," in *ICRA*, 2018, pp. 4693–4700.

[7] L. Tai, G. Paolo, and M. Liu, "Virtual-to-real deep reinforcement learning: Continuous control of mobile robots for mapless navigation," in *IEEE/RSJ*, 2017, pp. 31–36.

[8] M. Bojarski, D. D. Testa, D. Dworakowski, B. Firner, B. Flepp, P. Goyal, L. D. Jackel, M. Monfort, U. Muller, J. Zhang, X. Zhang, J. Zhao, and K. Zieba, "End to end learning for self-driving cars," *CoRR*, vol. abs/1604.07316, 2016.

[9] H.-T. L. Chiang, A. Faust, M. Fiser, and A. Francis, "Learning navigation behaviors end-to-end with autorl," *IEEE Robotics and Automation Letters*, vol. 4, no. 2, pp. 2007–2014, 2019.

[10] Y. F. Chen, M. Everett, M. Liu, and J. P. How, "Socially aware motion planning with deep reinforcement learning," in *IEEE/RSJ*, 2017, pp. 1343–1350.

[11] C. Chen, Y. Liu, S. Kreiss, and A. Alahi, "Crowd-robot interaction: Crowd-aware robot navigation with attention-based deep reinforcement learning," in *ICRA*, 2019, pp. 6015–6022.

[12] E. Yurtsever, J. Lambert, A. Carballo, and K. Takeda, "A survey of autonomous driving: Common practices and emerging technologies," *IEEE Access*, vol. 8, pp. 58 443–58 469, 2020.

[13] T. Räuker, A. Ho, S. Casper, and D. Hadfield-Menell, "Toward transparent ai: A survey on interpreting the inner structures of deep neural networks," 2023.

[14] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. J. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," in *ICLR 2014*, 2014.

[15] J. Uesato, A. Kumar, C. Szepesvári, T. Erez, A. Ruderman, K. Anderson, K. Dvijotham, N. Heess, and P. Kohli, "Rigorous agent evaluation: An adversarial approach to uncover catastrophic failures," *CoRR*, vol. abs/1812.01647, 2018.

[16] P. Trautman, J. Ma, R. M. Murray, and A. Krause, "Robot navigation in dense human crowds: the case for cooperation," in *ICRA*. IEEE, 2013, pp. 2153–2160.

[17] M. Quigley, K. Conley, B. Gerkey, J. Faust, T. Foote, J. Leibs, R. Wheeler, and A. Ng, "Ros: an open-source robot operating system," in *ICRA Workshop on Open Source Software*, vol. 3, 01 2009.

[18] M. Bansal, A. Krizhevsky, and A. Ogale, "Chauffeurnet: Learning to drive by imitating the best and synthesizing the worst," *arXiv preprint arXiv:1812.03079*, 2018.

[19] F. Codevilla, M. Müller, A. López, V. Koltun, and A. Dosovitskiy, "End-to-end driving via conditional imitation learning," in *ICRA*, 2018, pp. 4693–4700.

[20] F. Codevilla, E. Santana, A. M. López, and A. Gaidon, "Exploring the limitations of behavior cloning for autonomous driving," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 9329–9338.

[21] W. Zeng, W. Luo, S. Suo, A. Sadat, B. Yang, S. Casas, and R. Urtasun, "End-to-end interpretable neural motion planner," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8660–8669.

[22] Y. Lu, J. Fu, G. Tucker, X. Pan, E. Bronstein, B. Roelofs, B. Sapp, B. White, A. Faust, S. Whiteson, *et al.*, "Imitation is not enough: Robustifying imitation with reinforcement learning for challenging driving scenarios," *arXiv preprint arXiv:2212.11419*, 2022.

[23] A. Kendall, J. Hawke, D. Janz, P. Mazur, D. Reda, J.-M. Allen, V.-D. Lam, A. Bewley, and A. Shah, "Learning to drive in a day," in *ICRA*, 2019, pp. 8248–8254.

[24] J. Chen, S. E. Li, and M. Tomizuka, "Interpretable end-to-end urban autonomous driving with latent deep reinforcement learning," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 6, pp. 5068–5078, 2021.

[25] B. R. Kiran, I. Sobh, V. Talpaert, P. Mannion, A. A. Al Sallab, S. Yogamani, and P. Pérez, "Deep reinforcement learning for autonomous driving: A survey," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 6, pp. 4909–4926, 2021.

[26] X. Xiao, B. Liu, G. Warnell, and P. Stone, "Motion planning and control for mobile robot navigation using machine learning: a survey," *Autonomous Robots*, vol. 46, no. 5, pp. 569–597, 2022.

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

[28] X. Liang, T. Wang, L. Yang, and E. Xing, "Cirl: Controllable imitative reinforcement learning for vision-based self-driving," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 584–599.

[29] M. Vecerik, T. Hester, J. Scholz, F. Wang, O. Pietquin, B. Piot, N. Heess, T. Rothörl, T. Lampe, and M. Riedmiller, "Leveraging demonstrations for deep reinforcement learning on robotics problems with sparse rewards," *arXiv preprint arXiv:1707.08817*, 2017.

[30] A. Nair, B. McGrew, M. Andrychowicz, W. Zaremba, and P. Abbeel, "Overcoming exploration in reinforcement learning with demonstrations," in *ICRA*, 2018, pp. 6292–6299.

[31] B. Kang, Z. Jie, and J. Feng, "Policy optimization with demonstrations," in *International conference on machine learning*. PMLR, 2018, pp. 2469–2478.

[32] Y. Lu, J. Fu, G. Tucker, X. Pan, E. Bronstein, B. Roelofs, B. Sapp, B. White, A. Faust, S. Whiteson, *et al.*, "Imitation is not enough: Robustifying imitation with reinforcement learning for challenging driving scenarios," *arXiv preprint arXiv:2212.11419*, 2022.

[33] Z. Huang, J. Wu, and C. Lv, "Efficient deep reinforcement learning with imitative expert priors for autonomous driving," *IEEE Transactions on Neural Networks and Learning Systems*, 2022.

[34] U. Patel, N. K. S. Kumar, A. J. Sathyamoorthy, and D. Manocha, "Dwa-rl: Dynamically feasible deep reinforcement learning policy for robot navigation among mobile obstacles," in *ICRA*, 2021, pp. 6057–6063.

[35] M. Vitelli, Y. Chang, Y. Ye, A. Ferreira, M. Wołczyk, B. Osiński, M. Niendorf, H. Grimmett, Q. Huang, A. Jain, *et al.*, "Safetynet: Safe planning for real-world self-driving vehicles using machine-learned policies," in *ICRA*, 2022, pp. 897–904.

[36] T. Fan, P. Long, W. Liu, and J. Pan, "Distributed multi-robot collision avoidance via deep reinforcement learning for navigation in complex scenarios," *The International Journal of Robotics Research*, vol. 39, no. 7, pp. 856–892, 2020.

[37] S. Dey, A. Sadek, G. Monaci, B. Chidlovskii, and C. Wolf, "Learning whom to trust in navigation: dynamically switching between classical and neural planning," 2023.

[38] R. C. Coulter *et al.*, *Implementation of the pure pursuit path tracking algorithm*. Carnegie Mellon University, The Robotics Institute, 1992.

[39] AWS RoboMaker, "AWS RoboMaker Hospital World ROS package," GitHub repository, jul 2021. [Online]. Available: https://github.com/aws-robotics/aws-robomaker-hospital-world

[40] D. Helbing and P. Molnár, "Social force model for pedestrian dynamics," *Physical Review E*, vol. 51, pp. 4282–4286, May 1995.

[41] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Köpf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, *PyTorch: An Imperative Style, High-Performance Deep Learning Library*. Curran Associates Inc., 2019.

[42] S. Fujimoto, H. van Hoof, and D. Meger, "Addressing function approximation error in actor-critic methods," *CoRR*, vol. abs/1802.09477, 2018.