

# SOMTP: Self-Supervised Learning-Based Optimizer for MPC-Based Safe Trajectory Planning Problems in Robotics\*

Yifan Liu<sup>1</sup>, You Wang<sup>1,†</sup> and Guang Li<sup>1</sup>

**Abstract**—Model Predictive Control (MPC)-based trajectory planning has been widely used in robotics, and incorporating Control Barrier Function (CBF) constraints into MPC can greatly improve its obstacle avoidance efficiency. Unfortunately, traditional optimizers are resource-consuming and slow to solve such non-convex constrained optimization problems (COPs) while learning-based methods struggle to satisfy the non-convex constraints. In this paper, we propose SOMTP algorithm, a self-supervised learning-based optimizer for CBF-MPC trajectory planning. Specifically, first, SOMTP employs problem transcription to satisfy most of the constraints. Then the differentiable SLPG correction is proposed to move the solution closer to the safe set and is then converted as the guide policy in the following training process. After that, inspired by the Augmented Lagrangian Method (ALM), our training algorithm integrated with guide policy constraints is proposed to enable the optimizer network to converge to a feasible solution. Finally, experiments show that the proposed algorithm has better feasibility than other learning-based methods and can provide solutions much faster than traditional optimizers with similar optimality.

## I. INTRODUCTION

Constrained optimization problems (COPs) have been widely used in many areas, such as trajectory planning in robotics, power systems, scheduling, and logistics. In some scenarios, COP may impose high demands on solving time. MPC is a special COP that generates a sequence of optimal control inputs and states. In recent years, MPC has gained extensive application in collision avoidance trajectory planning for autonomous driving and robots [1]–[3]. To achieve efficient collision avoidance trajectory planning for MPC, [4] proposes an efficient CBF-MPC planner by incorporating discrete-time non-convex CBFs as constraints into MPC. CBF-MPC is then widely used in the robot [3], [5], [6]. CBF could enforce the system to avoid obstacles even when the reachable set is far away from the obstacles, as well as ensure forward invariance of the safe set [4]. However, the addition of CBF also increases the complexity of the COP and thus lengthens the time required to solve it.

MPC-based trajectory planning problem needs to be solved within a single time step. And if it is resolved sooner, the robot can respond faster to external changes and uncertainty. Traditional optimization methods for such non-convex COPs usually require large computational resources and a

long solving time. Learning-based methods can significantly reduce the solving time and have been applied to many optimization problems with or without simple constraints. However, in terms of MPC-based trajectory planning, fulfilling complex non-convex constraints becomes a major challenge for these methods [4], [7]. What’s more, such sequential optimization problems may be more challenging since many robots’ kinematic models, which act as equality constraints, are nonlinear. All in all, there is a demand for optimization methods to speed up optimization while strictly satisfying constraints.

In this paper, we propose SOMTP, a self-supervised learning-based optimizer, to solve the CBF-MPC-based safe trajectory planning problem with CBF constraints to avoid obstacles. SOMTP is highly concerned with constraint satisfaction and is capable of rapidly and effectively solving the CBF-MPC while ensuring feasibility. Specifically, first, we transcribe the original sequential optimization problem into nonlinear programming with only CBF constraints by employing the single-shooting method [8]. Then, the following two algorithms are proposed to satisfy the CBF constraints: (1) the SLPG correction algorithm; and (2) the ALM-based training method with guide policy constraints. The SLPG correction procedure is inspired by DC3 in [7] which can be incorporated into the training process. Unlike DC3, SLPG correction is to optimize a non-convex projection-type problem by performing sequential linearization, quadratic penalty method, and gradient descent. Compared with traditional projection-based method, the SLPG correction can deal with non-convex constraints and is more efficient in time and resource consuming. Instead of obtaining the true projection on the safe set, SLPG correction is intended to move the pre-solution from the optimizer network several steps closer to the safe set by generating a differentiable approximate solution within iteration limits. After that, we propose the ALM-based training method to achieve both optimality and feasibility by training the optimizer network and updating Lagrangian multipliers. Furthermore, in the ALM-based training method, the corrected solution from SLPG is again utilized as the guide policy to guide the network training. As a result, the optimizer network is expected to gradually converge to a solution of the CBF-MPC while constraints are satisfied.

In summary, our key contributions are as follows:

- SOMTP is a self-supervised learning-based optimization algorithm for the CBF-MPC-based trajectory planning problem and holds referential significance for other COPs and optimal control problems with non-convex

\*This work was supported by the Fundamental Research Funds for the Central Universities 226-2022-00086.

<sup>1</sup>Yifan Liu, You Wang and Guang Li are with State Key Laboratory of Industrial Control Technology, Institute of Cyber Systems and Control, Zhejiang University, Hangzhou, 310027, China. yifanliu@zju.edu.cn, king\_wy@zju.edu.cn, guangli@zju.edu.cn

<sup>†</sup> Corresponding author is You Wang.

constraints. Meanwhile, our strategy for CBF constraints is also meaningful for CBF-based safe RL.

- SLPG performs a differentiable correction, which can move the solution closer to the non-convex safe set and speed up the reduction of the violation. It is then treated as the guide policy in ALM-based training.
- ALM-based training algorithm is proposed to enable the optimizer network to converge to the feasible solution of the MPC.
- Guide policy constraints are integrated into the ALM-based training algorithm to guide the learning process and accelerate convergence.
- Experiments demonstrate that our SOMTP has better feasibility than other learning-based approaches. Compared with traditional optimization methods, it can provide high-quality, feasible solutions much faster while still maintaining a similar level of optimality.

The paper is organized as follows: Section II presents the related works. Section III presents the background and preliminaries. In section IV, SOMTP algorithm is proposed. Comparisons and experiments are given in V. Finally, conclusions are given in Section VI.

## II. RELATED WORKS

CBF-MPC trajectory planning belongs to the non-convex sequential COP and has been widely used in robotics. We intend to solve the above non-convex COP with learning-based methods. And our method's related works fall into the following three categories.

**Traditional optimization methods.** Before optimizing MPC, traditional methods need to transcribe such an optimal control problem into a nonlinear programming problem (which is also a COP) [9]. The main optimization methods for such nonlinear COP are the interior point method (IPM) and sequential quadratic programming (SQP). [10]. To obtain the COP's solution, these methods often require iterating multiple times and calculating the Hessian matrix, both of which can be time-consuming and resource-consuming. Moreover, the ALM algorithm mentioned above is also a traditional optimizer for solving COP, which can obtain solution and its associated Lagrangian multipliers through iterative methods while avoiding numerical instabilities [11].

**Learning-based methods.** Learning-based optimization methods are mainly divided into two types, i.e., supervised learning methods (SLMs) and self-supervised learning methods (SSLMs). SLMs concern training models that can map a problem instance's representation to the target solutions from a traditional optimizer [12]. These methods require the use of a traditional optimizer to generate the target optimal solutions and pay much attention to the errors (e.g., MSE or MAE) between the network's outputs and the target solutions [13]–[16]. SSLMs don't need additional traditional optimizers. Instead, they can directly train the optimizer network using the objective function and constraint violations [17]. One of the SSLMs is to add constraint violations to the loss function for gradient descent. On top of this, the DC3 algorithm proposed in [7] incorporates a correction process

for inequality constraints during training, the details of which can be seen in III. And [18] directly employs the DC3 to solve the MPC problem. DC3 is the most similar work to ours, with the following three major differences between our SOMTP and DC3: (1) We propose a differentiable SLPG correction for non-convex CBF constraints that can provide better direction for correction and push the pre-solution closer to the safe set; (2) The ALM-based training method can update associate Lagrangian multipliers and quadratic penalty terms during the training procedure; (3) To enhance the efficiency of ALM-based training, guide policy constraints are incorporated with SLGP correction as the guide policy. Besides, [17] proposes PDL that uses primal-dual learning to alternately train two networks, one the optimizer network and one the Lagrangian multipliers' network.

**Reinforcement learning (RL).** Safe RL expects to handle constraints during training, which has given us some inspiration. Some directly enforce constraints by projecting pre-solutions onto a safe set using convex optimization layers (e.g., [19]) in the case of general convex constraints [20], [21]. However, these projection-based approaches are inefficient since they often require a large amount of computing resources during training. What's more, the projection-type problem may be infeasible when the initial value is very distant from the safe set. Different from projection-based methods, our SLPG correction focuses on non-convex constraints and is intended to move the pre-solution several steps closer to the safe set within iteration limits. What's more, in recent years, CBF has gained popularity in MPC and safe RL methods due to the fact that it ensures the forward invariance of the safe set [22], [23].

## III. BACKGROUND

Trajectory planning aims to plan a trajectory for the robot from its initial position to its target state while avoiding obstacles and satisfying kinematic constraints. In recent years, optimization-based trajectory planning for mobile robotics like autonomous vehicles has been studied in many works. And the proposed algorithm in this study is intended to solve the aforementioned optimization-based trajectory planning problem with a learning-based method. The necessary preliminaries are as follows:

### A. CBF-MPC Trajectory Planning

In this study, we consider a 2-D mobile robot with the following nonlinear kinematic model:

$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k) \quad (1)$$

where the states and controls of the robotic system at time-step  $k$  are  $\mathbf{x}_k \in \mathbb{R}^3$  and  $\mathbf{u}_k \in \mathbb{R}^{n_u}$ , respectively.  $\mathbf{x}_k = [X \ Y \ \phi]$ , where  $X$  and  $Y$  represent the coordinates in the local frame O'X'Y' in Fig. 1.  $\phi$  represents the yaw angle of the robot in the same frame.

Assume the robot's initial state and goal state are  $\mathbf{x}_{in}$  and  $\mathbf{x}_{go}$ , respectively. And there are  $n_{obs} \geq 0$  obstacles in the local cost map. To avoid the obstacles, we use  $\mathcal{S}$  as a safe set,

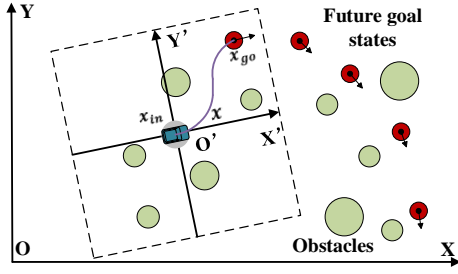


Fig. 1. Frames of the robotic system.

which is the superlevel set of a continuously differentiable function  $H : \mathcal{X} \subset \mathbb{R}^3$ :

$$\mathcal{S} = \{x \in \mathcal{X}, H(x) > 0\} \quad (2)$$

where  $H$  is the control barrier function (CBF) on  $\mathcal{S}$  if there exists  $\gamma \in (0, 1]$  such that for all  $x \in \mathcal{S}$ ,  $H$  satisfies

$$\sup_{u_k} \Delta H(x_k, u_k) = H(x_{k+1}) - H(x_k) \geq -\gamma H(x_k) \quad (3)$$

Given a CBF  $H$ , we can define the following set to guarantees the forward invariance of  $\mathcal{S}$  and the robots' safety.

$$G_{cbf} = \{u_k \in \mathcal{U}, -\Delta H(x_k, u_k) - \gamma H(x_k) \leq 0\} \quad (4)$$

After combining the above discrete-time CBF constraints, the loss function of the CBF-MPC trajectory planner is obtained which can be seen as follows:

$$\min_{u(\cdot), x(\cdot)} J(\cdot) = \sum_{k=0}^N \|x_k - x_{go}\|_Q^2 + \sum_{k=0}^{N-1} \|u_k\|_R^2 \quad (5)$$

$$\text{s.t. } x_{k+1} = f(x_k, u_k), \quad (5a)$$

$$x_0 = x_{in}, \quad (5b)$$

$$u_k \in [u_{min}, u_{max}], \quad (5c)$$

$$-\Delta H(x_k, u_k | O_j) - \gamma H(x_k | O_j) \leq 0 \quad (5d)$$

where  $N$  is the prediction horizon,  $Q$  and  $R$  are definite weighting coefficient matrix. The number of (5d) is  $N \cdot n_{obs}$ . To help represent these obstacles, we use  $O_j = [X_{o,j} \ Y_{o,j} \ R_{o,j}]$ ,  $j \in [0, n_{obs})$ , to represent the  $j$ -th obstacles, where  $X_{o,j}$  and  $Y_{o,j}$  represent the coordinates of the obstacle in the local frame  $O'X'Y'$  in Fig. 1.  $R_{o,j}$  represents the radius of the obstacle's smallest circumscribed circle. For the CBF function  $H(\cdot)$  in (5d), we employ the following function to simplify the obstacle avoidance problem:

$$H(x | O_j) = (x[0] - X_{o,j})^2 + (x[1] - Y_{o,j})^2 - (R_{o,j} + R + l_{ex})^2 \quad (6)$$

where  $R$  represents the radius of the robot and  $l_{ex}$  is the expansion length of the obstacle. And the discrete-time CBF in (5d) is generally non-convex according to [4].

The above CBF-MPC belongs to the non-convex optimal control problem and needs to be solved using some nonlinear optimization algorithms such as SQP and IPM. However, these traditional optimization methods often require large computing resources and a long time to work.

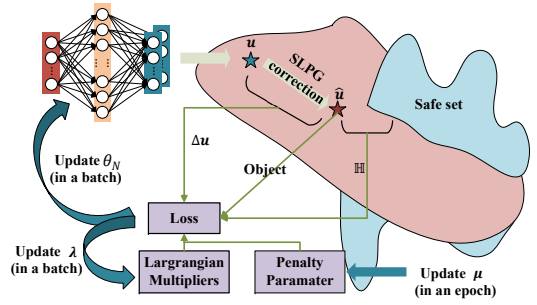


Fig. 2. The structure of SOMTP algorithm.

## B. DC3

DC3 is a learning-based method that intends to solve COPs by directly integrating two processes, i.e., equality completion and inequality correction, into the training procedure [7]. Consider the following optimization problem as an instance:

$$\min_y f_x(y), \text{ s.t. } h_x(y) = 0, \quad g_x(y) \leq 0 \quad (7)$$

where  $x$  is the problem data and  $y \in \mathbb{R}^n$  is the corresponding solution. Loss function, equality, and inequality constraints are denoted by  $f_x$ ,  $h_x$  and  $g_x$ , respectively. Overall, DC3 will employ a neural network to output a partial set of solution  $p \in \mathbb{R}^m$ ,  $m \leq n$ .  $p$  will then be completed to a full set of solution  $y_f = [p \ \phi_x(p)]^T \in \mathbb{R}^n$ , where they assume access to the function  $\phi_x$  to satisfy the equality constraints. Next, the  $y_f$  will be corrected to  $\hat{y}$  by performing gradient descent on the violations of the inequality constraints during the correction process. Specifically, the correction is to modify the  $y_f$  with the equation  $\hat{y} = y_f + \Delta y = [p - \gamma \Delta p \ \phi_x(p) - \gamma \partial \phi_x(p) / \partial p \Delta p]^T$  where  $\Delta p = \nabla_p \|\text{ReLU}(g_x(y_f))\|^2$  and  $\gamma$  is the constant learning rate.

After the above procedures, the overall network will be trained using backpropagation on the following loss function:

$$l_{soft} = f_x(\hat{y}) + \lambda_h \|h_x(\hat{y})\|^2 + \lambda_g \|\text{ReLU}(g_x(\hat{y}))\|^2 \quad (8)$$

where  $\lambda_h, \lambda_g > 0$  are hand-tuned hyper-parameters.

## IV. ALGORITHM

To achieve safe trajectory planning, we employ the MPC planner with CBF to keep the robot away from obstacles. As it can be seen, the optimization problem in (5) is a complex, non-convex, nonlinear, sequential optimal control problem with nonlinear equality and inequality constraints. To solve the aforementioned problem, inspired by optimization algorithms, safe RL algorithms, and the learning-with-correction framework in DC3, we finally propose the learning-based optimization algorithm SOMTP. The structure of SOMTP can be seen in Fig. 2.

### A. Problem Transcription

Before solving the problem in (5), states and obstacles need to be transformed to the local coordinate system, with the current state as the zero point, and thus the initial state

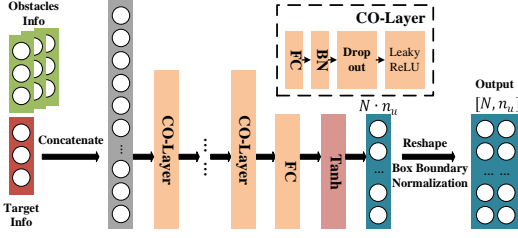


Fig. 3. The structure of the network. Each CO-Layer in the network has the structure in the upper right corner of the figure.

$\mathbf{x}_0 = \mathbf{x}_{in} = [0 \ 0 \ 0]$  can be seen as known. According to (5a) and (5b), the sequential trajectory with respect to the control inputs  $\mathbf{u}$  can be simplified as:

$$\mathbf{x}_k = \mathbf{F}(\mathbf{u}_0, \dots, \mathbf{u}_{k-1}), k \geq 1 \quad (9)$$

After employing the above single-shooting method, the trajectory optimization problem in (5) can be transcribed as a non-linear programming problem in (10) [8].

$$\min_{\mathbf{u}(\cdot)} \widehat{J}(\mathbf{u}, \mathbf{x}_{go}) \quad (10)$$

$$\text{s.t. } \mathbf{u}_k \in [\mathbf{u}_{min}, \mathbf{u}_{max}], k \in [0, N) \quad (10a)$$

$$H_{k,j}^{cbf}(\mathbf{u}_0, \dots, \mathbf{u}_k, \mathbf{O}_j) \leq 0, j \in [0, n_{obs}) \quad (10b)$$

where  $\widehat{J}(\cdot) = \sum_{k=1}^N \|\mathbf{F}(\mathbf{u}_{0:k-1}) - \mathbf{x}_{go}\|_{\mathbf{Q}}^2 + \sum_{k=0}^{N-1} \|\mathbf{u}_k\|_{\mathbf{R}}^2$ ,  $H_{k,j}^{cbf}(\cdot) = -\Delta H(\mathbf{F}(\mathbf{u}_{0:k-1}), \mathbf{u}_k | \mathbf{O}_j) - \gamma H(\mathbf{F}(\mathbf{u}_{0:k-1}) | \mathbf{O}_j)$ .

Take a neural network denoted as  $\theta_N$  to be the solver of the optimization problem (10), the structure of which is shown in Fig. 3. The output of the final activation layer tanh will be bounded within the range of [-1,1]. This output will then be forwarded to the final layer, which applies box boundary normalization to the result to rescale it to the range in (10a). And in this way, the problem in (10) will be transcribed to the following optimization problem (11)

$$\min_{\mathbf{u}=\pi_{\theta_N}} \widehat{J}(\mathbf{u}, \mathbf{x}_{go}) \quad (11)$$

$$\text{s.t. } H_{k,j}^{cbf}(\mathbf{u}_0, \dots, \mathbf{u}_k, \mathbf{O}_j) \leq 0, k \in [0, N), j \in [0, n_{obs}) \quad (11a)$$

where  $\pi_{\theta_N}$  represents the result of the  $\theta_N$ .

It is obvious that, following the transcription, the only constraints that require consideration in (11) are the CBF constraints denoted as (11a). To deal with the CBF constraints, we implement the following two policies: the SLPG correction and the policy-guided ALM-based training method.

### B. SLPG Correction

$\pi_{\theta_N}$  cannot directly guarantee the CBF constraints. In DC3, they insert an inequality correction procedure that can hopefully pull the results closer to or within the safe set by performing gradient descent on the violations of the inequality constraints. In terms of the CBF constraints in (11a), the correction procedure in DC3 is to repeat the following:

$$\hat{\mathbf{u}}_{dc3} = \mathbf{u} - \gamma_d \nabla_{\mathbf{u}} \sum \|\text{ReLU}(H_{k,j}^{cbf})\|^2 \quad (12)$$

where  $\hat{\mathbf{u}}_{dc3}$  is the corrected results with DC3-correction and  $\gamma_d$  is the associated step-size, which is man-tuned [7].

However, the (11) is a complex non-convex problem with nonlinear implicit function (9), and the CBF constraints show highly nonlinear and non-convex features. Therefore, the gradient descent-based correction method in DC3 may fail to find a suitable direction for correction. What's more, the corrected results  $\hat{\mathbf{u}}_{dc3}$  may conflict with the constraints in (10a).

Several safe RL methods employ projection algorithms to project the results to the safe set constraints during training procedures [20], [21]. Here, we also expect to design a projection-type correction procedure to deal with the CBF constraints. Define  $\Delta \mathbf{u}$  as the input's correction value, and it can be obtained from the following optimization problem:

$$\min_{\Delta \mathbf{u}(\cdot)} \sum_{k=0}^{N-1} \|\Delta \mathbf{u}_k\|_{\mathbf{R}}^2 \quad (13)$$

$$\text{s.t. } \mathbf{x}_k^c = \mathbf{F}(\mathbf{u}_{0:k-1} + \Delta \mathbf{u}_{0:k-1}) \quad (13a)$$

$$\mathbf{u}_k + \Delta \mathbf{u}_k \in [\mathbf{u}_{min}, \mathbf{u}_{max}], \quad (13b)$$

$$H_{k,j}^{cbf}(\mathbf{u}_{0:k} + \Delta \mathbf{u}_{0:k}, \mathbf{O}_j) \leq 0, j \in [0, n_{obs}) \quad (13c)$$

Moreover, the projection problem in the current safe RL algorithm is equivalent to a quadratic programming (QP) problem with linear or convex constraints. However, the problem (13) is not a QP problem at all with nonlinear and non-convex constraints, and solving this optimization problem will consume much time.

Inspired by SQP and SLQP algorithms that solve complex non-convex problems by solving a sequential of QP problem, we perform a first-order Taylor expansion on the constraints in (13) and optimize several QP sub-problems. The first-order Taylor expansion is performed at the points  $(\mathbf{u}, \mathbf{x})$  which are the results from  $\pi_{\theta_N}$ . Thus, we can obtain a QP-type sub-problem as follows:

$$\min_{\Delta \mathbf{u}(\cdot)} \sum_{k=0}^{N-1} \|\Delta \mathbf{u}_k\|_{\mathbf{R}}^2 \quad (14)$$

$$\text{s.t. } \mathbf{x}_k^c = \mathbf{x}_k + [\nabla_{\mathbf{u}_{0:k-1}} \mathbf{F}] \Delta \mathbf{u}_{0:k-1} \quad (14a)$$

$$\mathbf{u}_k + \Delta \mathbf{u}_k \in [\mathbf{u}_{min}, \mathbf{u}_{max}], \quad (14b)$$

$$H_{k,j}^{cbf}(\mathbf{u}_{0:k}, \mathbf{O}_j) + [\nabla_{\mathbf{u}_{0:k}} H_{k,j}^{cbf}] \Delta \mathbf{u}_{0:k} \leq 0, j \in [0, n_{obs}) \quad (14c)$$

where (14a) can indeed be ignored since it has already been inserted into (14c) during the transcription.

However, the QP problem in (14) may be infeasible since the original solution from  $\theta_N$  is insufficient during training, sometimes very distant from the safe set, and fails to meet the CBF constraints. This kind of situation happens frequently during the beginning stages of training. Since the purpose of the correction is to find a proper direction to reduce the violations of the CBF constraints, we adopt the Quadratic Penalty Method [24] to transfer the CBF constraints to the loss function. Instead of solving (14), we turn to solving the QP problem in (15) to pull the solution closer to the safe

---

**Algorithm 1:** SLPG Correction

---

**Input** : Solution from  $\pi_{\theta_N}$ :  $\mathbf{u}$ . Obstacles:  $\mathbf{O}$ . Max correction steps:  $n_m$ . Max QP steps:  $i_m$ . Box constraints:  $[\mathbf{u}_{min}, \mathbf{u}_{max}]$ .  $\lambda_c$ .

**Output:** Corrected solution:  $\hat{\mathbf{u}}$ .  $\Delta \mathbf{u}$

**Init** :  $\Delta \mathbf{u}$

```
1  $\mathbf{u}^0 \leftarrow \mathbf{u}$ ;  
2 for  $n \in [0, n_m]$  do  
3   SequentialLinear( $\mathbf{u}^n$ ), Init( $\Delta \mathbf{u}^0$ );  
4   for  $i \in [0, i_m]$  do  
5      $\mathbf{d}^i \leftarrow \nabla_{\Delta \mathbf{u}^i} J_{corr}$ ;  
6      $\gamma_c \leftarrow$  Armijo-BoxLinearSearch with  
       condition (18);  
7      $\Delta \hat{\mathbf{u}}^{i+1} \leftarrow \Delta \mathbf{u}^i - \gamma_c \mathbf{d}^i$ ; in (16);  
8      $\Delta \mathbf{u}^{i+1} = \text{Clamp}(\mathbf{u}^n + \Delta \hat{\mathbf{u}}^{i+1}, \mathbf{u}_{min}, \mathbf{u}_{max}) - \mathbf{u}$   
9   end  
10   $\mathbf{u}^{n+1} \leftarrow \mathbf{u}^n + \Delta \mathbf{u}^{i_m}$  |  $\Delta \mathbf{u}, \hat{\mathbf{u}} \leftarrow \Delta \mathbf{u} + \Delta \mathbf{u}^{i_m}, \mathbf{u}^{n+1}$ ;  
11  Break if Max(ReLU( $H_{k,j}^{cbf}(\hat{\mathbf{u}}_{0:k}, \mathbf{O}_j)$ )) is tiny  
12 end
```

---

set.

$$\min_{\Delta \mathbf{u}(\cdot)} \sum_{k=0}^{N-1} \|\Delta \mathbf{u}_k\|_R^2 + \lambda_c \sum_{k=0}^{N-1} \sum_{j=0}^{n_{obs}-1} \|\text{ReLU}(\mathcal{H}_{k,j}^{cbf}(\cdot))\|^2 \quad (15)$$

$$\text{s.t. } \mathbf{u}_k + \Delta \mathbf{u}_k \in [\mathbf{u}_{min}, \mathbf{u}_{max}], \quad (15a)$$

where  $\mathcal{H}_{k,j}^{cbf}$  is equivalent to (14c),  $\lambda_c > 0$  denotes the penalty parameter. Employing  $J_{corr}$  to represent the loss function in (15).

In terms of the box constraints QP problem in (15), we perform gradient descent on the  $J_{corr}$  to obtain the suitable direction. The clamp function is also executed to ensure compliance with the box constraints. Details are provided below:

$$\Delta \hat{\mathbf{u}}^{i+1} = \Delta \mathbf{u}^i - \gamma_c \nabla_{\Delta \mathbf{u}^i} J_{corr} \quad (16)$$

where the appropriate step size  $\gamma_c$  is determined by the linear search method that will be proposed later.

$$\Delta \mathbf{u}^{i+1} = \text{Clamp}(\mathbf{u} + \Delta \hat{\mathbf{u}}^{i+1}, \mathbf{u}_{min}, \mathbf{u}_{max}) - \mathbf{u} \quad (17)$$

Furthermore, it was suggested by [25] that naive clamping for box constraints could damage the direction of descent and convergence. To address this, we modify the Armijo condition in [26] to propose the Armijo-box condition in (18) for the box constraints in the linear search method for  $\gamma_c$ , ensuring that  $\mathbf{u} + \Delta \hat{\mathbf{u}}^{i+1}$  stays as much inside the box constraints as possible.

$$\text{(Armijo}(\cdot) \text{ and } \mathbf{u} + \Delta \hat{\mathbf{u}}^{i+1}(\gamma_c) \in [\mathbf{u}_{min}, \mathbf{u}_{max}]) \\ \text{or iteration counter} \geq \text{maximum iteration} \quad (18)$$

Following the aforementioned procedures, we finally present the SLPG correction approach, which includes the key steps, i.e., the Sequential Linearization, the Quadratic Penalty method and the Gradient Decent with modified

linear search condition. Pseudocode of SLPG can be seen in Algorithm 1.

### C. Augmented Lagrangian-based Training Algorithm with Guide Policy Constraints

After the SLPG correction, we obtained the corrected input  $\hat{\mathbf{u}}$ . By incorporating SLPG into the training process, the overall objective of the training can be defined as addressing a following optimization problem (19) which is equivalent to (11). (19b) is also employed as the **guide policy constraints** that can convert SLPG into a sub-guide policy during the training.

$$\min_{\hat{\mathbf{u}}=\text{SLPG}(\pi_{\theta_N})} \widehat{J}(\hat{\mathbf{u}}, \mathbf{x}_{go}) \quad (19)$$

$$\text{s.t. } H_{k,j}^{cbf}(\hat{\mathbf{u}}_{0:k}, \mathbf{O}_j) \leq 0, k \in [0, N], j \in [0, n_{obs}] \quad (19a)$$

$$\Delta \mathbf{u} = \hat{\mathbf{u}} - \pi_{\theta_N} = 0 \quad (19b)$$

Thus, we can build the following loss function to train the policy network, which can be seen as the ALM of the (19),

$$\begin{aligned} \mathcal{L}(\hat{\mathbf{u}}, \lambda, \mu) = & \widehat{J}(\hat{\mathbf{u}}, \mathbf{x}_{go}) + \sum_{i=0}^N \sum_{j=0}^{n_{obs}} \lambda_{i,j}^c \mathbb{H}_{i,j} + \frac{\mu^c}{2} \sum_{i=0}^N \sum_{j=0}^{n_{obs}} \mathbb{H}_{i,j}^2 \\ & + \lambda^{du} \cdot \text{abs}(\Delta \mathbf{u}) + \frac{\mu^{du}}{2} \|\Delta \mathbf{u}\|^2 \end{aligned} \quad (20)$$

where  $\mathbb{H}_{i,j} = \text{ReLU}(H_{i,j}^{cbf}(\hat{\mathbf{u}}_{0:i}, \mathbf{O}_j))$ ,  $\lambda = [\lambda^c \ \lambda^{du}]$  is the Lagrange multiplier and belongs to a positive set.  $\mu = [\mu^c \ \mu^{du}]$  represents the corresponding penalty parameter, with all values being positive. And  $\text{abs}(\cdot)$  computes the absolute value of each element. It is important to note that the presence of the penalty on  $\Delta \mathbf{u}$  is due to an implicit expectation in the original problem. This expectation is to guide the original policy network  $\theta_N$  to approach the corrected policy from SLPG in order to enable  $\theta_N$  to directly generate a feasible solution, thus making (19) be equivalent to (11). What's more, such policy-guided training can optimize the training direction and speed up convergence.

A gradient search is then performed on  $\mathcal{L}(\hat{\mathbf{u}}, \lambda, \mu)$  to obtain the optimal tuple, which is shown as follows:

$$\theta_N^{n+1} = \theta_N^n - \eta_\theta \nabla_{\theta_N^n} \mathcal{L}(\hat{\mathbf{u}}^n, \lambda^n, \mu) \quad (21)$$

$$\lambda_{i,j}^{c,n+1} = \lambda_{i,j}^{c,n} + \mu^{c,m} \mathbb{E}_{batch}[\mathbb{H}_{i,j}] \quad (22)$$

$$\lambda^{du,n+1} = \lambda^{du,n} + \mu^{du,m} \mathbb{E}_{batch}[\text{abs}(\Delta \mathbf{u})] \quad (23)$$

where  $batch$  denotes the training batch size and  $\eta_\theta$  represents the learning rate.  $n$  and  $m$  denote the  $n$ -th training step and  $m$ -th epoch, respectively. What's more, the penalty parameters will be updated at the end of each epoch when the following condition occurs:

$$\text{if } \mathbb{E}_{epoch}[\|\mathbb{H}\|^2] < \beta_c / \epsilon_c : \beta_c = \mathbb{E}_{epoch}[\|\mathbb{H}\|^2] \\ \mu^{c,m+1} = \min(\epsilon_c \cdot \mu^{c,m}, \mu_{max}^c) \quad (24)$$

$$\text{if } \mathbb{E}_{epoch}[\|\Delta \mathbf{u}\|^2] < \beta_{du} / \epsilon_{du} : \beta_{du} = \mathbb{E}_{epoch}[\|\Delta \mathbf{u}\|^2] \\ \mu^{du,m+1} = \min(\epsilon_{du} \cdot \mu^{du,m}, \mu_{max}^{du}) \quad (25)$$

---

**Algorithm 2:** SOMTP (During Training)

---

**Input** : Dataset:  $\mathbb{D}$ ; Penalty:  $\mu$ ; Step size:  $\eta_\theta$ ;  
Penalty updating values:  $\epsilon_c, \epsilon_{du}$ ; Max  
penalty:  $\mu_{max}^c, \mu_{max}^{du}$   
**Output:** Policy Network:  $\theta_N$   
**Init** :  $\theta_N, \lambda, \mu, \beta_c, \beta_{du}$

- 1 **for**  $m \leftarrow 0$  **to**  $M$  **do**
- 2     Shuffle dataset  $\mathbb{D}$ ;
- 3     **for**  $n \leftarrow 0$  **to**  $N$  **do**
- 4         From  $\mathbb{D}$  choose Data  $d$ ;  $\mathbf{u} \leftarrow \theta_N(d)$ ;
- 5          $\hat{\mathbf{u}}, \Delta \mathbf{u} \leftarrow \text{Alg. SLPG}(\mathbf{u}, d)$ ;
- 6         Calculate  $\mathcal{L}(\hat{\mathbf{u}}, \lambda, \mu)$  using (20);
- 7         Update  $\theta_N$  using (21);
- 8         Update  $\lambda$  using (22) and (23);
- 9     **end**
- 10     If the condition occurs, update  $\mu$  and  $[\beta_c, \beta_{du}]$   
       using (24) and (25)
- 11 **end**
- 12 **Return**  $\theta_N$

---

where  $\epsilon_c, \epsilon_{du} > 1$  are constants and *epoch* denotes the training epoch. As for now, the overall training process has been presented, which can also be seen in Algorithm 2.

## V. EXPERIMENTAL RESULTS

In this section, experiments are carried out to compare the effectiveness of SOMTP to baselines on CBF-MPC-based safe trajectory planning problem. The chosen experimental agent is the autonomous vehicle with the following kinematic model:

$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k) = \begin{bmatrix} X + v \cos \phi dt \\ Y + v \sin \phi dt \\ \phi + v \tan q / L dt \end{bmatrix}$$

where  $\mathbf{u}_k = [v \ q]^T$ . Specifically, we compare SOMTP against the following baselines:

- **IPOPT, SQPmethod** (Traditional Optimizer): We employ the IPM-based optimizer IPOPT [27] and the SQP-based optimizer SQPmethod (with OSQP [28] to solve the sub-QP problem) in CasADi [29]. In addition, CasADi is also utilized as an algorithmic differentiation tool in both optimizers;
- **MSE, MAE**: trained to minimize the L2 norm error or the L1 norm error between the optimizer network's outputs and the optimal solutions from the traditional optimizer (IPOPT);
- **Penalty**: trained to minimize a soft loss in (8);
- **DC3**: proposed in [7];
- **PDL**: proposed in [17];
- **ALM** (ablation study): SOMTP without SLPG correction at both train and test time;
- **SOMTP-w/o $\Delta\mathbf{u}$**  (ablation study): SOMTP without guide policy constraints in (19b) so that the ALM loss in (20) will not contain the penalty on  $\Delta\mathbf{u}$ ;

The same neural network is employed across all experiments, with the structure in Fig. 3 and five CO-Layers. Each CO-Layer has a hidden layer fully connected with 2000 nodes, which is then followed by a dropout layer (rate 0.3). The optimizer network is trained using PyTorch [30], and the training process is executed on a system equipped with a GeForce RTX 3080 GPU and an Intel Xeon 2.9GHz CPU. To train the neural network, we generate the dataset with a total of 1 million examples (with train/test/validation ratio 18/1/1). For each example, the goal state  $\mathbf{x}_{go}$  and obstacles  $\mathbf{O}_j$  are randomly generated within the local cost map (6 m  $\times$  6 m). Each example has three obstacles ( $n_{obs} = 3$ ), each with a circular shape and a radius  $R_o \in [0, 0.5]$  m. The prediction horizon  $N$  is 20.  $\mathbf{Q}$  and  $\mathbf{R}$  are  $\text{diag}([2.0, 2.0, 1.0])$  and  $\text{diag}([1.0, 1.5])$ , respectively. In addition,  $(\gamma, R, l_{ex}) = (0.5, 0.3, 0.1)$ , which is used in (5) and (6).

### A. Results on the Test Dataset

Table. I compare the performance of the SOMTP algorithm with traditional optimizers and other learning-based methods on the test dataset. For our SLPG correction procedure, we use  $(n_m^{train}, i_m^{train}) = (2, 2)$  and  $(n_m^{test}, i_m^{test}) = (10, 2)$ . The following five indicators are selected to compare the results:

- **Obj.**: mean object value from the loss function in (5);
- **Mean CBF**: mean violations of the cumulative CBF non-convex constraints  $\mathbb{E}_{test} [\sum_{i=0}^N \sum_{j=0}^{n_{obs}} \mathbb{H}_{i,j}]$ ;
- **Max CBF**: maximum violations on CBF constraints;
- **Infe. (%)**: infeasibility rate over the test dataset = (number of infeasible instances) / (total number of instances)  $\times 100\%$ ;
- **Time (ms)**: mean time cost to solve.

According to Table. I, SOMTP has a remarkably low infeasibility rate that is only surpassed by IPOPT. In addition, SOMTP is capable of providing the result 79 $\times$  faster than IPOPT with similar optimality. The results also indicate that SSLMs are more feasible than SLMs. According to the results, the SLPG correction employed at test time can effectively reduce infeasibility. Nevertheless, the results of ALM and SOMTP-w/o $\Delta\mathbf{u}$  indicate that only applying SLPG correction to the training procedure without a guide policy constraint does not significantly contribute to its feasibility. By incorporating the guide policy constraints with  $\Delta\mathbf{u}$  into the training procedure, the feasibility of SOMTP (or SOMTP $\not\leftarrow$ ) is significantly enhanced. This suggests that the guide policy constraints in (19b) and (20) can guide the learning process to a better point and accelerate convergence. What's more, though some baselines have very low Obj., they violate the obstacle constraints, and their strategies are likely to find the shortest trajectory between the current position and the target point in the obstacle-free task.

Overall, SOMTP achieves state-of-the-art (SOTA) in learning-based optimization algorithms for the CBF-MPC-based trajectory planning problem, with a much lower infeasibility rate, fewer violations of CBF constraints, and similar optimality.

TABLE I  
RESULTS ON TEST DATASET (WITH 50000 INSTANCES).<sup>1</sup>

Algorithm	Obj.	Mean CBF	Max CBF	Infe. (%)	Time (ms)
IPOPT	165.31	0.0000	0.0000	0.00	<b>83.676</b>
SQPmethod	170.72	0.0000	0.0010	0.35	115.18
<b>SOMTP *</b>	167.93	<b>0.0000</b>	<b>0.0252</b>	<b>0.07</b>	<b>1.057</b>
SOMTP, $\not\leq$	167.92	0.0000	0.0555	0.11	0.727
SOMTP-w/o $\Delta u$	165.57	0.0002	0.1537	0.64	1.841
SOMTP-w/o $\Delta u$ , $\not\leq$	165.55	0.0005	0.1537	0.85	0.717
ALM	164.70	0.0003	0.1746	0.87	<b>0.724</b>
DC3	163.23	0.0006	0.2091	1.54	1.037
DC3, $\not\leq$	163.23	0.0007	0.2129	1.54	0.741
Penalty	162.88	0.0008	0.1792	1.74	0.728
PDL	<b>160.53</b>	0.0027	0.1807	3.71	0.744
MSE	169.74	0.1007	0.2703	16.56	0.760
MAE	164.89	0.1015	0.2794	17.48	0.739

$\not\leq$  denotes the absence of correction phase (like SLPG correction) at test time.

<sup>1</sup> To address the speed discrepancy in algorithmic differentiation like the calculation of gradients, traditional optimizers and learning-based algorithms that necessitate correction procedures utilize CasADi to provide gradients during testing.

\* SOMTP can provide the output 79 $\times$  faster than IPOPT with similar optimality and a very low infeasibility rate. Furthermore, it is far more feasible than other learning-based algorithms and shows a lower infeasibility rate compared to the traditional optimizer SQPmethod.

TABLE II  
RESULT ON ROBOT IN TEN  
DIFFERENT PLANNING TASKS.

Alg.	Suc.	Dist.
IPOPT	90%	0.1112
SOMTP	90%	0.2008
PDL	80%	0.2253
DC3	80%	0.2714
MAE	40%	0.0734
MSE	40%	0.0734

## B. Results on the Robot's Navigation

We also apply our algorithm to the robot in 10 different tasks to compare the performance in continuous trajectory planning. The robot will repeatedly solve the CBF-MPC based trajectory planning problem in (5) in each time-step until it reaches the target area or encounters an obstacle. Results can be seen in Table. II and Fig. 4. The task is considered successfully achieved when the robot reaches the target area. The success rate (the number of successful instances / the total number of instances) is denoted as [Suc.], whereas the average final weighted distance to the target in the success instances is represented as [Dist.]. Experiments show that during robot navigation, continuous trajectory planning and control will amplify the impact of infeasibility, causing frequent collisions with MAE and MSE. Over all the tasks, SOMTP achieves a high success rate and can sometimes complete tasks that IPOPT fails. Nevertheless, the SSLM-based optimizers ultimately achieve larger weighted distances than IPOPT. This could be attributed to the fact

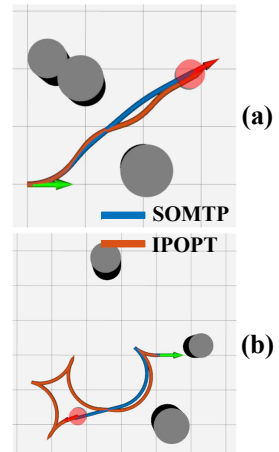


Fig. 4. Continuous trajectory planning tasks on robot. Each grid is 1 m in width. The target area is denoted by a red circle. Green arrays represent the initial states, while red arrays represent the target states.

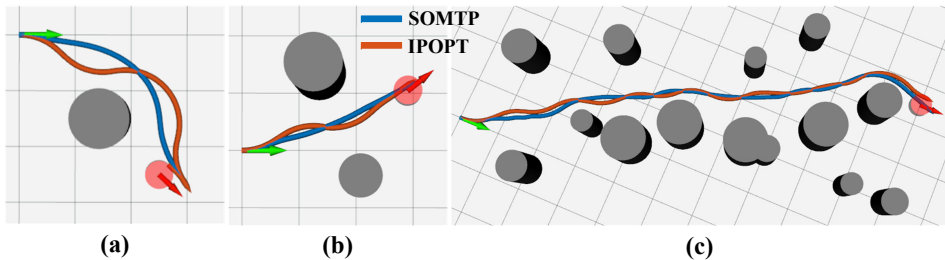


Fig. 5. Apply optimizers to robot's trajectory planning tasks with variable quantities of obstacles. Each grid is 1 m in width. The target area is denoted by a red circle. Green arrays represent the initial states, while red arrays represent the target states.

that SSLM-based methods have a higher tendency to train for tasks where the target is more away from the current position, which leads to larger losses.

Ultimately, to confirm the practicality and robustness of the SOMTP-based optimizer, we tested its performance on tasks with variable quantities of obstacles and long-distance planning. The figures in Fig. 5 illustrate the practicality and robustness of the algorithm.

## VI. CONCLUSIONS

We propose SOMTP, a self-supervised learning-based optimization algorithm for CBF-MPC-based trajectory planning problem, which belongs to a complex non-convex COP. The core components of SOMTP consist of three parts: (1) the problem transcription can transcribe the problem into a neural network and satisfy most of the constraints; (2) the SLPG correction can pull the initial solution closer to the non-convex safe set and provide a guide policy for the following training process; and (3) the ALM-based training

process with guide policy constraints integrated ensures that the network reaches a feasible point with regard to both feasibility and optimality. Experiments demonstrate that our SOMTP has significantly greater feasibility compared to previous learning-based algorithms, as well as being considerably faster than traditional optimizers while maintaining a similar level of optimality.

What's more, we believe that the SOMTP algorithm holds referential significance for other COPs and optimal control problems with non-convex constraints. And our strategy for CBF constraints is also meaningful for CBF-based safe RL. Therefore, in future work, we will try to extend SOMTP to other COPs as well as explore safe RL strategies based on the SOMTP algorithm.

## REFERENCES

- [1] J. Ji, A. Khajepour, W. W. Melek, and Y. Huang, "Path planning and tracking for vehicle collision avoidance based on model predictive control with multiconstraints," *IEEE Transactions on Vehicular Technology*, vol. 66, no. 2, pp. 952–964, 2017.
- [2] M. Ammour, R. Orjuela, and M. Basset, "A mpc combined decision making and trajectory planning for autonomous vehicle collision avoidance," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 12, pp. 24 805–24 817, 2022.
- [3] Z. Jian, Z. Yan, X. Lei, Z. Lu, B. Lan, X. Wang, and B. Liang, "Dynamic control barrier function-based model predictive control to safety-critical obstacle-avoidance of mobile robot," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2023, pp. 3679–3685.
- [4] J. Zeng, B. Zhang, and K. Sreenath, "Safety-critical model predictive control with discrete-time control barrier function," in *American Control Conference (ACC)*, 2021, pp. 3882–3889.
- [5] W. Xiao, C. G. Cassandras, and C. A. Belta, "Bridging the gap between optimal trajectory planning and safety-critical control with applications to autonomous vehicles," *Automatica*, vol. 129, p. 109592, 2021.
- [6] Y. Yu, D. Shan, O. Benderius, C. Berger, and Y. Kang, "Formally robust and safe trajectory planning and tracking for autonomous vehicles," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 12, pp. 22 971–22 987, 2022.
- [7] P. L. Donti, D. Rolnick, and J. Z. Kolter, "Dc3: A learning method for optimization with hard constraints," in *International Conference on Learning Representations (ICLR)*, 2021.
- [8] J. T. Betts, *Practical Methods for Optimal Control and Estimation Using Nonlinear Programming*, 2nd ed. Society for Industrial and Applied Mathematics, 2010. [Online]. Available: <https://epubs.siam.org/doi/abs/10.1137/1.9780898718577>
- [9] C. Kirches, L. Wirsching, H. G. Bock, and J. P. Schlöder, "Efficient direct multiple shooting for nonlinear model predictive control on long horizons," *Journal of Process Control*, vol. 22, no. 3, pp. 540–550, 2012.
- [10] J. Nocedal and S. J. Wright, *Numerical Optimization (2nd ed.)*. Springer, 2006.
- [11] R. Andreani, E. G. Birgin, J. M. Martínez, and M. L. Schuverdt, "On augmented lagrangian methods with general lower-level constraints," *SIAM Journal on Optimization*, vol. 18, no. 4, pp. 1286–1309, 2008.
- [12] J. Kotary, F. Fioretto, P. Van Hentenryck, and B. Wilder, "End-to-end constrained optimization learning: A survey," in *International Joint Conference on Artificial Intelligence (IJCAI)*, 2021, pp. 4475–4482.
- [13] S. Koziel and L. Leifsson, *Surrogate-based modeling and optimization*. Springer, 2013.
- [14] J. Kotary, F. Fioretto, and P. Van Hentenryck, "Learning hard optimization problems: A data generation perspective," *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 34, pp. 24 981–24 992, 2021.
- [15] A. S. Zamzam and K. Baker, "Learning optimal solutions for extremely fast ac optimal power flow," in *IEEE International Conference on Communications, Control, and Computing Technologies for Smart Grids (SmartGridComm)*, 2020, pp. 1–6.
- [16] F. Fioretto, T. W. Mak, and P. Van Hentenryck, "Predicting ac optimal power flows: Combining deep learning and lagrangian dual methods," in *Proceedings of the AAAI conference on artificial intelligence (AAAI)*, vol. 34, no. 01, 2020, pp. 630–637.
- [17] S. Park and P. Van Hentenryck, "Self-supervised primal-dual learning for constrained optimization," in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, vol. 37, no. 4, 2023, pp. 4052–4060.
- [18] C. Diehl, J. Adamek, M. Krüger, F. Hoffmann, and T. Bertram, "Differentiable constrained imitation learning for robot motion planning and control," *arXiv preprint arXiv:2210.11796*, 2022.
- [19] A. Agrawal, B. Amos, S. Barratt, S. Boyd, S. Diamond, and J. Z. Kolter, "Differentiable convex optimization layers," *Advances in neural information processing systems (NeurIPS)*, vol. 32, 2019.
- [20] T.-Y. Yang, J. Rosca, K. Narasimhan, and P. J. Ramadge, "Projection-based constrained policy optimization," in *International Conference on Learning Representations (ICLR)*, 2020.
- [21] Y. Emam, G. Notomista, P. Glofelfel, Z. Kira, and M. Egerstedt, "Safe reinforcement learning using robust control barrier functions," *IEEE Robotics and Automation Letters*, pp. 1–8, 2022.
- [22] R. Cheng, G. Orosz, R. M. Murray, and J. W. Burdick, "End-to-end safe reinforcement learning through barrier functions for safety-critical continuous control tasks," in *Proceedings of the AAAI conference on artificial intelligence (AAAI)*, vol. 33, no. 01, 2019, pp. 3387–3395.
- [23] X. Zhang, Y. Peng, W. Pan, X. Xu, and H. Xie, "Barrier function-based safe reinforcement learning for formation control of mobile robots," in *2022 International Conference on Robotics and Automation (ICRA)*, 2022, pp. 5532–5538.
- [24] C. Broyden and N. Attia, "Penalty functions, newton's method, and quadratic programming," *Journal of optimization theory and applications*, vol. 58, pp. 377–385, 1988.
- [25] Y. Tassa, N. Mansard, and E. Todorov, "Control-limited differential dynamic programming," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2014, pp. 1168–1175.
- [26] L. Armijo, "Minimization of functions having lipschitz continuous first partial derivatives," *Pacific Journal of mathematics*, vol. 16, no. 1, pp. 1–3, 1966.
- [27] A. Wächter and L. T. Biegler, "On the implementation of an interior-point filter line-search algorithm for large-scale nonlinear programming," *Mathematical programming*, vol. 106, pp. 25–57, 2006.
- [28] B. Stellato, G. Banjac, P. Goulart, A. Bemporad, and S. Boyd, "OSQP: an operator splitting solver for quadratic programs," *Mathematical Programming Computation*, vol. 12, no. 4, pp. 637–672, 2020. [Online]. Available: <https://doi.org/10.1007/s12532-020-00179-2>
- [29] J. A. E. Andersson, J. Gillis, G. Horn, J. B. Rawlings, and M. Diehl, "CasADi – A software framework for nonlinear optimization and optimal control," *Mathematical Programming Computation*, vol. 11, no. 1, pp. 1–36, 2019.
- [30] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, "Pytorch: An imperative style, high-performance deep learning library," *Advances in neural information processing systems (NeurIPS)*, vol. 32, 2019.