

A Hybrid Learning-to-Optimize Framework for Mixed-Integer Quadratic Programming

Viet-Anh Le

Mu Xie

Rahul Mangharam

Department of Electrical & Systems Engineering, University of Pennsylvania, Philadelphia, PA 19104, USA

VIETANH@SEAS.UPENN.EDU

MUX2001@SEAS.UPENN.EDU

RAHULM@SEAS.UPENN.EDU

Editors: G. Sukhatme, L. Lindemann, S. Tu, A. Wierman, N. Atanasov

Abstract

In this paper, we propose a learning-to-optimize (L2O) framework to accelerate solving parametric mixed-integer quadratic programming (MIQP) problems, with a particular focus on mixed-integer model predictive control (MI-MPC) applications. The framework learns to predict integer solutions with enhanced optimality and feasibility by integrating supervised learning (for optimality), self-supervised learning (for feasibility), and a *differentiable quadratic programming (QP) layer*, resulting in a hybrid L2O framework. Specifically, a neural network (NN) is used to learn the mapping from problem parameters to optimal integer solutions, while a differentiable QP layer is integrated to compute the corresponding continuous variables given the predicted integers and problem parameters. Moreover, a *hybrid loss function* is proposed, which combines a supervised loss with respect to the global optimal solution, and a self-supervised loss derived from the problem's objective and constraints. The effectiveness of the proposed framework is demonstrated on two benchmark MI-MPC problems, with comparative results against purely supervised and self-supervised learning models.

Keywords: Learning to optimize, mixed-integer quadratic programming, mixed-integer model predictive control.

1. Introduction

Mixed-integer optimization is fundamental to many control applications involving discrete decision-making, such as autonomous driving (Quirynen et al., 2024), traffic signal coordination with connected automated vehicles (Le and Malikopoulos, 2024), multi-robot pickup and delivery (Camisa et al., 2022), motion planning and task assignment for robot fleets (Salvado et al., 2018), and signal temporal logic specifications (Belta and Sadraddini, 2019). However, solving a mixed-integer program (MIP) is NP-hard because it requires combinatorial search over discrete decision variables. Consequently, computation time can grow exponentially with problem size or constraint complexity, making MIPs generally unsuitable for real-time control or decision-making.

Advancements in machine learning and differentiable programming provide a promising opportunity to accelerate MIP solvers through learning-to-optimize (L2O) frameworks. The literature on L2O for MIP problems is still limited but has gained increasing attention in recent years. The current state of the art can be categorized into two main approaches: (i) supervised learning (SL) and (ii) self-supervised learning (SSL). Classical SL approaches, e.g., (Cauligi et al., 2021, 2022; Le et al., 2025), train neural networks (NNs) to minimize the supervised loss between predictions and the optimal integer solutions generated by an optimization solver such as Gurobi (Gurobi Optimization, LLC, 2021). A major drawback of SL approaches is that they do not guarantee feasibility,

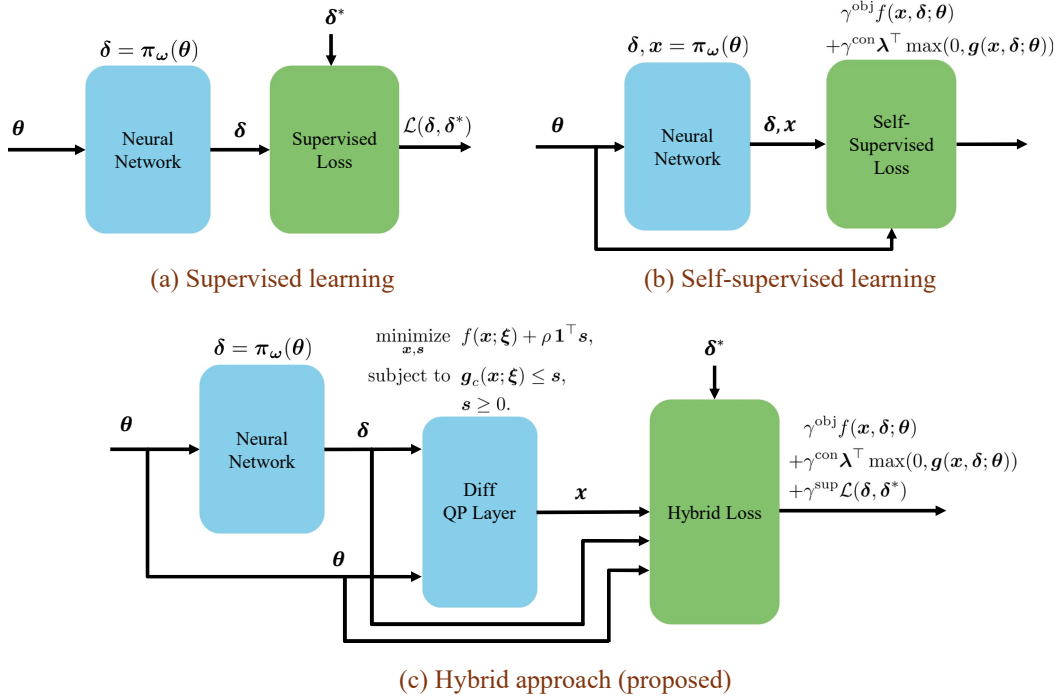


Figure 1: Architecture of the proposed hybrid framework (c) compared with supervised learning (a) and self-supervised learning (b). In our framework, the NN takes the problem parameters θ to predict the integer solution δ , while the QP layer computes the continuous solution x based on θ and δ . In conventional SL and SSL, the NN is trained to predict the integer solution without considering the continuous solution or to predict both the integer and continuous solutions, respectively.

i.e., the resulting convex continuous problems obtained by fixing the learned integer variables can be infeasible. SSL approaches (Tang et al., 2024; Boldocky et al., 2025), on the other hand, do not rely on labeled data and can improve feasibility by using a loss function that combines the objective function with a penalty for constraint violation. For example, (Tang et al., 2024) proposed a framework for mixed-integer nonlinear programming that solves the integer-relaxed problem, combined with integer correction layers to ensure integrality and a projection step to improve feasibility. (Boldocky et al., 2025) considered parametric MIQP problems within a differentiable predictive control framework, which constrains the integer solutions and optimal control inputs to a neural state-feedback law. However, trained SSL models may produce feasible but suboptimal solutions, since differentiable programming techniques such as gradient descent may converge to locally optimal solutions.

Addressing the limitations of both SL (infeasibility) and SSL (suboptimality), we propose a novel hybrid L2O framework that strategically combines both training paradigms with an integrated differentiable QP layer. The hybrid approach proposed in this paper shares conceptual similarities with physics-informed machine learning (PINN) by embedding the known optimization structure of the problem (the QP layer) as an inductive bias, much as PINNs embed physical laws (the PDEs). First, we propose a novel architecture in which a differentiable QP layer is integrated into the network to better incorporate the optimization structure. During training, the QP layer takes the NN predictions of the integer decision variables as input and outputs the corresponding optimal solu-

tions for the continuous variables, but it may be infeasible. To address this issue and ensure that gradients can always be computed, we propose a simple yet effective approach that introduces a differentiable layer for the relaxed QP problem. We prove that, if the penalty weight is chosen sufficiently large, the relaxed problem yields either the exact optimal solution when the original QP is feasible or the minimally infeasible solution otherwise. Second, we propose a new loss function design, called a *hybrid loss function*, defined as a weighted sum of SL and SSL losses. This approach allows the framework to balance the feasibility-optimality trade-off in L2O. The overall architecture of our framework in comparison with SL and SSL can be illustrated in Fig. 1.

Our framework differs from existing work in the relevant literature in the following aspects. First, a major difference lies in how we encode the dependency between continuous variables, integer variables, and problem parameters during training. In SL approaches (Cauligi et al., 2021, 2022), the goal is to learn the mapping from problem parameters to the optimal integer variables, while disregarding the continuous variables during training. In online prediction, the continuous variables are then obtained by solving a QP given the NN predictions of the integer variables. In contrast, SSL approaches (Tang et al., 2024; Boldocky et al., 2025) train NNs to predict both the continuous and integer variables from the problem parameters. Thus, the prediction obtained directly from the NN does not explicitly account for the dependency between continuous and integer variables. In our approach, we incorporate the continuous variables into training by integrating a QP layer, based on the fact that the optimal continuous variables are the solutions of a parametric QP given the integer variables and the MIQP problem parameters. Therefore, our approach better incorporates the underlying optimization structure into both training and prediction than supervised and self-supervised learning. Second, we combine supervised and self-supervised learning objectives to define a hybrid loss function, rather than relying on either one alone. This hybrid loss allows the framework to balance the optimality of supervised learning with training labels and the feasibility improvement of self-supervision during training. Thus, our proposed framework can be viewed as a compromise between SL and SSL. We show through numerical examples that the hybrid loss function achieves near-global optimality and minimal constraint violation for most problem instances.

2. Preliminaries

This section provides a brief discussion on the parametric MIQPs and L2O for accelerating solving MIQP problems.

2.1. Parametric MIQPs

We consider a parametric MIQP problem that takes the following form:

$$\underset{\mathbf{x} \in \mathcal{X}, \delta \in \mathcal{I}}{\text{minimize}} \quad f(\mathbf{x}, \delta; \boldsymbol{\theta}), \tag{1a}$$

$$\text{subject to} \quad \mathbf{g}(\mathbf{x}, \delta; \boldsymbol{\theta}) \leq 0, \tag{1b}$$

where \mathbf{x} is the vector of continuous optimization variables, δ is the vector of integer optimization variables, and $\boldsymbol{\theta}$ is the vector of problem parameters. We let \mathcal{X} and \mathcal{I} be the domains for continuous and integer optimization variables, and $\mathbf{g}(\cdot) = [g_1(\cdot), \dots, g_r(\cdot)]$ be the vector of r linear constraints. We assume that \mathcal{X} and \mathcal{I} are non-empty. In this work, we consider a convex quadratic objective

function and linear constraints, i.e.,

$$\begin{aligned} f(\mathbf{x}, \boldsymbol{\delta}; \boldsymbol{\theta}) &= \frac{1}{2} \begin{bmatrix} \mathbf{x} \\ \boldsymbol{\delta} \end{bmatrix}^\top \mathbf{Q}(\boldsymbol{\theta}) \begin{bmatrix} \mathbf{x} \\ \boldsymbol{\delta} \end{bmatrix} + \mathbf{p}(\boldsymbol{\theta})^\top \begin{bmatrix} \mathbf{x} \\ \boldsymbol{\delta} \end{bmatrix}, \\ \mathbf{g}(\mathbf{x}, \boldsymbol{\delta}; \boldsymbol{\theta}) &= \mathbf{G}(\boldsymbol{\theta}) \begin{bmatrix} \mathbf{x} \\ \boldsymbol{\delta} \end{bmatrix} - \mathbf{h}(\boldsymbol{\theta}). \end{aligned} \quad (2)$$

where $\mathbf{Q}(\boldsymbol{\theta}) \succeq 0$. Note that given known parameters $\boldsymbol{\theta}$ and integer variables $\boldsymbol{\delta}$, the optimal solution of the continuous variables can be obtained by solving a QP problem, if it is feasible, as follows:

$$\underset{\mathbf{x} \in \mathcal{X}}{\text{minimize}} \quad f(\mathbf{x}; \boldsymbol{\xi}), \quad (3a)$$

$$\text{subject to} \quad \mathbf{g}_c(\mathbf{x}; \boldsymbol{\xi}) \leq 0, \quad (3b)$$

where $\mathbf{g}_c(\cdot)$ denotes the components of $\mathbf{g}(\cdot)$ that involve at least one continuous decision variable, and $\boldsymbol{\xi} = [\boldsymbol{\delta}^\top, \boldsymbol{\theta}^\top]^\top$. The objective function and constraint function in (3) can be expressed as:

$$f(\mathbf{x}; \boldsymbol{\xi}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q}_x(\boldsymbol{\xi}) \mathbf{x} + \mathbf{p}_x(\boldsymbol{\xi})^\top \mathbf{x}, \quad (4)$$

$$\mathbf{g}_c(\mathbf{x}; \boldsymbol{\xi}) = \mathbf{G}_x(\boldsymbol{\xi}) \mathbf{x} - \mathbf{h}_x(\boldsymbol{\xi}). \quad (5)$$

Mixed-integer MPC: A typical application that can greatly benefit from L2O approaches is model predictive control (MPC). In MPC, we solve a parametric optimization problem at every time step, in which the problem parameters may include, for instance, a given initial state, target state, to name a few. Using machine learning to solve or assist in solving parametric MPC problems enables real-time implementation of complex MPC, such as nonlinear MPC or mixed-integer MPC. We consider the state $\mathbf{x}_t \in \mathcal{X} \subset \mathbb{R}^{n_x}$ and control inputs $\mathbf{u}_t \in \mathcal{U} \subset \mathbb{R}^{n_u}$ as the continuous decision variables and let $\boldsymbol{\delta}_t \in \mathcal{I} \subset \mathbb{Z}^{n_\delta}$ be the integer decision variables. Let H be the control horizon length, and $\mathbf{x}_{0:H}, \mathbf{u}_{0:H-1}, \boldsymbol{\delta}_{0:H-1}$ be the concatenated vectors over the control horizon. Given a vector of problem parameters $\boldsymbol{\theta} \in \mathbb{R}^{n_p}$, a parametric MI-MPC can be written as:

$$\begin{aligned} \underset{\mathbf{x}_{0:H}, \mathbf{u}_{0:H-1}, \boldsymbol{\delta}_{0:H-1}}{\text{minimize}} \quad & c_H(\mathbf{x}_H; \boldsymbol{\theta}) + \sum_{t=0}^{H-1} c_t(\mathbf{x}_t, \mathbf{u}_t, \boldsymbol{\delta}_t; \boldsymbol{\theta}), \\ \text{subject to} \quad & \mathbf{x}_0 = \mathbf{x}_{\text{init}}(\boldsymbol{\theta}), \\ & \mathbf{x}_{t+1} = \mathbf{f}(\mathbf{x}_t, \mathbf{u}_t, \boldsymbol{\delta}_t; \boldsymbol{\theta}), \quad t = 0, \dots, H-1, \\ & \mathbf{g}_t(\mathbf{x}_t, \mathbf{u}_t, \boldsymbol{\delta}_t; \boldsymbol{\theta}) \leq 0, \quad t = 0, \dots, H-1, \\ & \mathbf{g}_H(\mathbf{x}_H; \boldsymbol{\theta}), \end{aligned} \quad (6)$$

where the stage cost $c_t(\cdot)$ and terminal cost $c_H(\cdot)$ are convex quadratic, while the dynamics $\mathbf{f}(\cdot)$, the inequality constraints $\mathbf{g}_t(\cdot)$ and $\mathbf{g}_H(\cdot)$ are assumed to be linear functions. The objective function and constraints are functions of the parameter vector $\boldsymbol{\theta} \in \Theta$, where $\Theta \subseteq \mathbb{R}^{n_p}$ is the admissible set of parameters.

2.2. Learning to Optimize for MIQPs

Due to the combinatorial nature, the problem (1) is computationally demanding to solve, where finding the optimal solution may scale exponentially with the problem size. Meanwhile, solving (3)

while the integers are fixed is significantly cheaper to solve than the MIQP. An interesting approach to accelerating solving problems of the form (1) is to learn a map between the vector of problem parameters θ and the discrete optimizer δ^* by an NN, $\delta^* = \pi_\omega(\theta)$, where ω is the vector of network weights. Overall, there are two main approaches for learning π_ω , including (i) supervised learning and (ii) self-supervised learning.

Supervised Learning: The NN π_ω can be trained by a classical SL approach. In SL, we collect the optimal solutions $\delta^{i,*}$ corresponding to each θ^i , for $i = 1, \dots, M$, obtained from a solver. We then use the dataset $\{\theta^i, \delta^{i,*}\}$ to train a NN that minimizes the following loss function:

$$\underset{\omega}{\text{minimize}} \quad \frac{1}{M} \sum_{i=1}^M \mathcal{L}(\pi_\omega(\theta^i), \delta^{i,*}), \quad (7)$$

where \mathcal{L} denotes a supervised loss (e.g., Huber or cross-entropy) between the predicted outputs and labels. However, the prediction from an SL model may not ensure that the resulting QP is feasible, although the original MIQP is feasible.

Self-Supervised Learning: Self-supervised learning, contrary to SL, does not rely on labeled data for training the model. It instead trains models by minimizing the objective function and constraint violation directly from model predictions. In generic SSL, an NN is trained to predict both the integer and continuous variables, i.e., $(\delta^i, \mathbf{x}^i) = \pi_\omega(\theta^i)$, and to train the NNs given a dataset with M training instances, the following self-supervised loss function is used:

$$\underset{\omega}{\text{minimize}} \quad \frac{1}{M} \sum_{i=1}^M f(\mathbf{x}^i, \delta^i; \theta^i) + \boldsymbol{\lambda}^\top \max(0, \mathbf{g}(\mathbf{x}^i, \delta^i; \theta^i)), \quad (8a)$$

$$\text{subject to} \quad (\delta^i, \mathbf{x}^i) = \pi_\omega(\theta^i), \quad (8b)$$

In (8), we include a penalty for constraint violation with max penalty (implemented via a ReLU) function. $\boldsymbol{\lambda} \in \mathbb{R}_{>0}^r$ is a vector of penalty parameters that balances the trade-off between minimizing the objective function and satisfying the constraints. Although the penalty methods lack formal guarantees, they often outperform their hard constraint counterparts in practice. Training the NN with a self-supervised loss function can improve, though it does not guarantee, feasibility. Nevertheless, a major drawback of this approach is that since the self-supervised loss (8) is non-convex with respect to ω , gradient-based methods may not converge and may converge to a sub-optimal solution. Moreover, in MPC applications, designing an NN to predict the optimal continuous decision variables from the problem parameters is generally challenging, as it is difficult to enforce the system dynamics on the network outputs (Cauligi et al., 2021), unless the optimal integer and continuous solutions at each time step are constrained to follow a state-feedback law, as in differentiable predictive control (Boldocky et al., 2025).

Remark 1 (Differentiating through discrete operations) *In many approaches, the NNs are designed to directly output discrete values, using discrete operations such as rounding to produce those outputs. These discrete operations lead to non-differentiability and hinder the use of standard differentiable programming for network training. A common approach to address this issue is the straight-through estimator (STE) (Bengio et al., 2013), which enables backpropagation through discrete operations. During the forward pass, STE applies a non-differentiable operation to obtain discrete values. During the backward pass, it bypasses the non-existent gradients of these operations by replacing them with those of smooth surrogate functions. This approach was used in*

self-supervised learning frameworks for mixed-integer programming (Tang et al., 2024; Boldocky et al., 2025). We also use this technique in our framework.

We observe that the strength of SL in finding global solutions corresponds to the weakness of SSL, and vice versa, the strength of SSL in improving feasibility corresponds to the weakness of SL. Therefore, an interesting idea is to combine SL and SSL to exploit the benefits of both approaches.

3. Proposed Framework

In this section, we present an L2O framework for MIQPs in which a differentiable QP layer is incorporated, and a hybrid loss function combining SL and SSL is proposed.

3.1. Differentiable QP Layers for Feasible and Infeasible Problems

Given known δ^* and θ , the optimal solution of the continuous variables can be obtained by solving the QP problem (3), if it is feasible. Thus, we can consider the QP problem as a differentiable layer within deep learning architectures, denoted by $\mathbf{x} = \text{QP}(\delta, \theta)$. Therefore, it leads to the following problem in which we approximate the integer solutions by NNs, and the continuous solutions by a QP layer:

$$\underset{\omega}{\text{minimize}} \quad f(\mathbf{x}, \delta; \theta), \quad (9a)$$

$$\text{subject to} \quad \delta = \pi_{\omega}(\theta), \quad (9b)$$

$$\mathbf{x} = \text{QP}(\delta, \theta), \quad (9c)$$

$$\mathbf{g}(\mathbf{x}, \delta; \theta) \leq 0. \quad (9d)$$

To ensure the validity of the L2O framework using a QP layer and differentiable programming, we need the following assumption.

Assumption 1 *The QP problem (3) is strictly convex.*

As stated in Amos and Kolter (2017, Theorem 1), Assumption 1 is needed to ensure that the QP layer is subdifferentiable everywhere, and differentiable at all but a measure-zero set of points, because the solution of a strictly convex QP is continuous. The question is how to compute the gradient of the optimal solution \mathbf{x}^* with respect to the argument, i.e., $\frac{\partial \mathbf{x}^*}{\partial \xi}$. If the problem is feasible, these derivatives can be obtained by differentiating the KKT conditions (sufficient and necessary conditions for optimality) of (3). However, since the methods in (Amos and Kolter, 2017; Agrawal et al., 2019) rely on KKT conditions, it assumes the QP problem is feasible. On the other hand, in our framework, the optimization problems might be infeasible during training, given different values of the integer variables from the NN. Thus, we cannot directly incorporate the differentiable QP layer in (Amos and Kolter, 2017; Agrawal et al., 2019) into our framework. In this section, we present a simple yet efficient way to handle infeasibility during training, as described below.

To this end, we introduce slack variables \mathbf{s} for the constraints, leading to the following QP:

$$\underset{\mathbf{x} \in \mathcal{X}, \mathbf{s}}{\text{minimize}} \quad f(\mathbf{x}; \xi) + \rho \mathbf{1}^\top \mathbf{s}, \quad (10a)$$

$$\text{subject to} \quad \mathbf{g}_c(\mathbf{x}; \xi) \leq \mathbf{s}, \quad (10b)$$

$$\mathbf{s} \geq 0. \quad (10c)$$

For ease of notations, in the rest of this section, we omit the argument ξ while mentioning the terms involving it. Since (10) is feasible given any realization of ξ as long as the domain for x is non-empty, we can apply the KKT conditions. First, we formulate the Lagrangian of (10) as follows:

$$L(x, s, \mu, \kappa) = \frac{1}{2}x^\top Q_x x + p_x^\top x + \rho \mathbf{1}^\top s + \mu^\top (G_x x - h_x - s) - \kappa^\top s, \quad (11)$$

where $\mu \geq 0$ and $\kappa \geq 0$ are the dual variables on the constraints, and $\rho > 0$ is a penalty weight for the slack variables. The KKT conditions for stationarity, primal feasibility, and complementary slackness are

$$Q_x x^* + p_x + G_x^\top \mu^* = 0, \quad (12a)$$

$$\rho \mathbf{1} - \mu^* - \kappa^* = 0, \quad (12b)$$

$$D(\mu^*)(G_x x^* - h_x - s^*) = 0, \quad (12c)$$

$$D(\kappa^*)s^* = 0, \quad (12d)$$

where $D(\cdot)$ is the operation creating a diagonal matrix from a vector. Taking the differentials of the KKT conditions (12), we obtain

$$dQ_x x^* + Q_x dx + dp_x + dG_x^\top \mu^* + G_x^\top d\mu = 0, \quad (13a)$$

$$d\mu + d\kappa = 0, \quad (13b)$$

$$D(G_x x^* - h_x - s^*)d\mu + D(\mu^*)(dG_x x^* + G_x dx - dh_x - ds) = 0, \quad (13c)$$

$$D(s^*)d\kappa + D(\kappa^*)ds = 0, \quad (13d)$$

or in the matrix form as follows:

$$\begin{bmatrix} Q_x & \mathbf{0} & G_x & \mathbf{0} \\ D(\mu^*)G_x & -D(\mu^*) & D(G_x x^* - h_x - s^*) & \mathbf{0} \\ \mathbf{0} & D(\kappa^*) & \mathbf{0} & D(s^*) \\ \mathbf{0} & \mathbf{0} & I & I \end{bmatrix} \begin{bmatrix} dx \\ ds \\ d\mu \\ d\kappa \end{bmatrix} = \begin{bmatrix} -dQ_x x^* - dp_x - dG_x^\top \mu^* \\ -D(\mu^*)(dG_x x^* - dh_x) \\ \mathbf{0} \\ \mathbf{0} \end{bmatrix}. \quad (14)$$

Using these equations, we can form the Jacobians of x^* and s^* with respect to any of the problem parameters. The details can be found in (Amos and Kolter, 2017), (Agrawal et al., 2019).

Note that given any x , $s^* = \max(0, g_c(x))$. In other words, the use of slack variables in (10) is equivalent to using the max penalty function. The following theorem shows that if the original QP (3) is feasible, then by choosing a sufficiently large value for ρ , solving (10) yields the same solution as (3).

Theorem 2 *If the original QP (3) is feasible and let x'^* and μ'^* be the optimal solutions and multipliers of the problem, respectively. If the penalty weight ρ is chosen such that $\rho > \|\mu'^*\|_\infty$, then the optimal solutions of (10) and the original QP (3) are the same.*

The proof of this theorem follows directly from Proposition 5.25 in (Bertsekas, 2014) and is therefore omitted. Therefore, if the original QP is feasible, we get the exact derivative of the optimal solution by using the KKT conditions of the relaxed problem. In the infeasible case, we show in the following theorem that, under some mild conditions, the obtained solution from (10) is a point that minimizes the constraint violation, and among all the solutions with minimal constraint violation, the obtained solution also minimizes the original objective function.

Theorem 3 *If the original QP (3) is infeasible and assume that either one of the following properties hold:*

- $\mathbf{Q}_x \succ 0$, which means $f(\mathbf{x})$ is coercive.
- \mathcal{X} is compact (closed and bounded).

Let us denote $v(\mathbf{x}) := \mathbf{1}^\top \max(0, \mathbf{g}_c(\mathbf{x}))$. If the penalty weight ρ is chosen sufficiently large, then (i) a minimizer $(\mathbf{x}_\rho^*, \mathbf{s}_\rho^*)$ of (10) achieves minimal total violation, i.e., $v(\mathbf{x}_\rho^*) = v^*$ and (ii)

$$\mathbf{x}_\rho^* \in \arg \min_{\mathbf{x} \in \mathcal{X}} \{f(\mathbf{x}) : v(\mathbf{x}) = v^*\}. \quad (15)$$

The proof for Theorem 3 is given in Appendix A.

3.2. Hybrid Training Loss

Our training framework relies on a *hybrid loss function* that combines supervised and self-supervised losses. The main advantage of this hybrid loss is that it leverages the strengths of supervised learning (SL) in achieving global solution optimality and self-supervised learning (SSL) in improving constraint satisfaction. The hybrid loss used to train the neural network can be defined as follows:

$$\underset{\omega}{\text{minimize}} \quad \frac{1}{M} \sum_{i=1}^M \gamma^{\text{obj}} f(\mathbf{x}^i, \delta^i; \theta^i) + \gamma^{\text{con}} \boldsymbol{\lambda}^\top \max(0, \mathbf{g}(\mathbf{x}^i, \delta^i; \theta^i)) + \gamma^{\text{sup}} \mathcal{L}(\delta^i, \delta^{i,*}), \quad (16)$$

where γ^{obj} , γ^{con} , and $\gamma^{\text{sup}} \in \mathbb{R}_{\geq 0}$ are the weights for the objective value, constraint violation, and supervised loss, respectively. Note that in our framework, the constraints involving at least one continuous variable can be directly handled using the differentiable QP layer, while the constraints involving only the integer variables must be incorporated into the loss function.

4. Results and Discussions

We validate our hybrid L2O framework on two benchmark MI-MPC problems: (i) collision avoidance for robot navigation and (ii) simplified thermal energy tanks (Boldocky et al., 2025). In the first example, binary variables are used to formulate the collision-avoidance constraints, thus, there are several coupling constraints between the integer and continuous variables. Meanwhile, the second example involves integer decision variables and demonstrates the case where the integers appear in the objective function. The details of the two examples are provided in Appendix B. For each example, a multilayer perceptron network is constructed with four hidden layers, 128 neurons per layer, and ReLU activation functions. Our implementation and examples are available at <https://github.com/mlab-upenn/L2O-MIQP>.

We compare the proposed hybrid L2O framework with an SL model and an SSL model. Note that the SSL model considered here follows the architecture of our framework with the differentiable QP layer, rather than the conventional design in which the NN is trained to predict both integer and continuous solutions. However, the loss function used to train the SSL model does not include the supervised term. We evaluate the trained models using two metrics: constraint violation and optimality gap. We separate the violations into those associated with constraints involving continuous variables and those involving only integer variables. The optimality gap is expressed as a

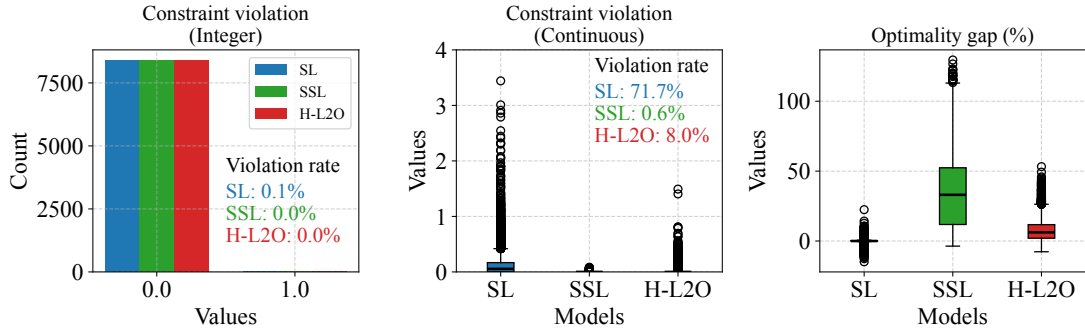


Figure 2: Statistical comparison of the three models: hybrid L2O (H-L2O), supervised learning (SL), and self-supervised learning (SSL), for the robot navigation example.

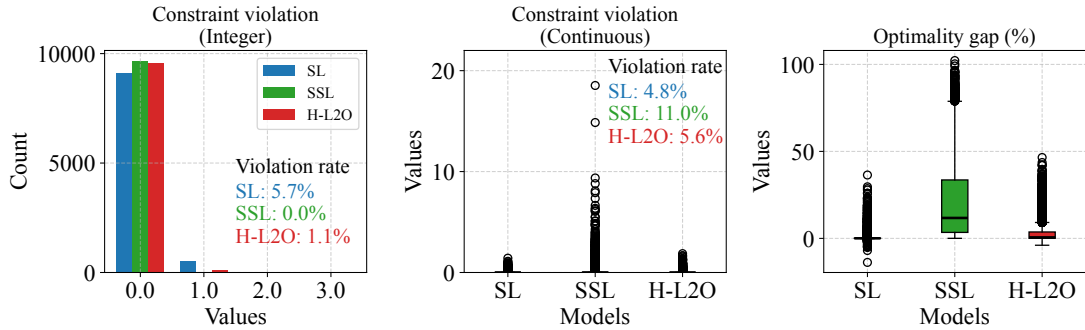


Figure 3: Statistical comparison of the three models: hybrid L2O (H-L2O), supervised learning (SL), and self-supervised learning (SSL), for thermal energy tank example.

percentage, computed as the ratio of the objective gap to the optimal objective value. We show the statistical comparison for the two examples in Figures 2 and 3, respectively. In each figure, the left panel shows the violations for integer-only constraints in the form of barplots, while the boxplots in the middle and right panels show the continuous-constraint violations and the optimality gap, respectively. For the figures showing constraint violations, we also report the violation rate, i.e., the percentage of validation problems in which the obtained solutions violate the constraints. For the robot navigation example, the results indicate that all three models satisfy the constraints involving only integer variables. However, for continuous-constraint violations, the SL model fails to ensure constraint satisfaction in 71.7% of the test cases (a consequence of SL training only for the integer solution, which does not guarantee a feasible continuous solution from the subsequent QP), whereas the hybrid L2O and SSL models exhibit much smaller violation rates of 8% and 0.6%, respectively, demonstrating improved constraint satisfaction. The plot of the optimality gap shows that the hybrid L2O model achieves better optimality than the SSL model, although it exhibits a slightly larger gap than the SL model. A similar trend in optimality is observed in the second example. Regarding constraint satisfaction, the hybrid L2O model outperforms the SL model for constraints involving only integer variables, achieving a lower violation rate of 1.1% compared to 5.7%. Meanwhile, the two models achieve comparable levels of satisfaction for continuous constraints (5.6% vs. 4.8% violation rate). In contrast, the SSL model perfectly satisfies the integer-only constraints but performs poorly on the continuous ones. Overall, the results show that the proposed hybrid L2O frame-

Table 1: Average solving time and standard deviation for GUROBI solver and the L2O solver.

Problem	Solver	Avg. time (ms)	Std. dev. (ms)
Robot navigation example	GUROBI	69.49	19.06
	L2O	7.55	1.31
Energy tank example	GUROBI	15.40	8.61
	L2O	1.31	0.39

work effectively balances feasibility and optimality, achieving optimality comparable to supervised learning while improving constraint satisfaction.

Finally, we report the computation times in Table 1, comparing the L2O approach (either SL, SSL, or the hybrid L2O) with the GUROBI solver (Gurobi Optimization, LLC, 2021). Note that the computation time for L2O approach includes both the NN prediction time and the time required to solve the relaxed QP problem. The results confirm that the L2O approach significantly reduces the overall solving time compared to a state-of-the-art MIQP solver. In particular, the L2O approach achieves approximately $9\times$ and $12\times$ faster computation for the robot navigation and energy tank examples, respectively.

Although our proposed hybrid L2O framework demonstrates some benefits over SL and SSL approaches, it still has certain limitations that can be addressed in future research. First, the choice of weights in the hybrid loss function significantly affects the optimality and feasibility performance and currently requires manual tuning. Second, integrating the differentiable QP layer considerably increases the training time compared to purely supervised learning. Finally, the current framework is limited to MIQPs, while integrating differentiable optimization layers into general mixed-integer convex or nonlinear programming problems remains an open challenge.

5. Conclusions

In this work, we developed a hybrid L2O framework for MIQPs that integrates a differentiable QP layer and a hybrid loss function combining supervised learning and self-supervised learning. We designed a model architecture in which a neural network learns the mapping from problem parameters to optimal integer variables, while the differentiable QP layer computes the corresponding continuous variables. This architecture enables the optimization structure to be incorporated into the learning process. To balance solution optimality and constraint feasibility during training, we defined a hybrid loss that linearly combines supervised and self-supervised terms. We validated the framework on two benchmark MPC examples, showing that the hybrid L2O approach effectively balances optimality and feasibility: it achieves better constraint satisfaction than purely supervised learning and better optimality than purely self-supervised learning. Future work will focus on extending the framework to mixed-integer convex and nonlinear programming problems.

Acknowledgments

This work was partially supported by US DoT Safety21 National University Transportation Center and NSF grants CISE-2431569.

References

- Akshay Agrawal, Brandon Amos, Shane Barratt, Stephen Boyd, Steven Diamond, and J Zico Kolter. Differentiable convex optimization layers. *Advances in neural information processing systems*, 32, 2019.
- Brandon Amos and J Zico Kolter. Optnet: Differentiable optimization as a layer in neural networks. In *International conference on machine learning*, pages 136–145. PMLR, 2017.
- Calin Belta and Sadra Sadraddini. Formal methods for control synthesis: An optimization perspective. *Annual Review of Control, Robotics, and Autonomous Systems*, 2(1):115–140, 2019.
- Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.
- Dimitri P Bertsekas. *Constrained optimization and Lagrange multiplier methods*. Academic press, 2014.
- Ján Boldocký, Shahriar Dadras Javan, Martin Gulán, Martin Mönnigmann, and Ján Drgoňa. Learning to solve parametric mixed-integer optimal control problems via differentiable predictive control. *arXiv preprint arXiv:2506.19646*, 2025.
- Andrea Camisa, Andrea Testa, and Giuseppe Notarstefano. Multi-robot pickup and delivery via distributed resource allocation. *IEEE Transactions on Robotics*, 39(2):1106–1118, 2022.
- Abhishek Cauligi, Preston Culbertson, Edward Schmerling, Mac Schwager, Bartolomeo Stellato, and Marco Pavone. Coco: Online mixed-integer control via supervised learning. *IEEE Robotics and Automation Letters*, 7(2):1447–1454, 2021.
- Abhishek Cauligi, Ankush Chakrabarty, Stefano Di Cairano, and Rien Quirynen. Prism: Recurrent neural networks and presolve methods for fast mixed-integer optimal control. In *Learning for Dynamics and Control Conference*, pages 34–46. PMLR, 2022.
- Gurobi Optimization, LLC. Gurobi optimizer reference manual, 2021. URL <http://www.gurobi.com>.
- Viet-Anh Le and Andreas A Malikopoulos. Distributed Optimization for Traffic Light Control and Connected Automated Vehicle Coordination in Mixed-Traffic Intersections. *IEEE Control Systems Letters*, 8:2721–2726, 2024.
- Viet-Anh Le, Panagiotis Kounatidis, and Andreas A. Malikopoulos. Combining Graph Attention Networks and Distributed Optimization for Multi-Robot Mixed-Integer Convex Programming. In *2025 64th IEEE Conference on Decision and Control*, 2025.
- Rien Quirynen, Sleiman Safaoui, and Stefano Di Cairano. Real-time mixed-integer quadratic programming for vehicle decision-making and motion planning. *IEEE Transactions on Control Systems Technology*, 2024.
- João Salvado, Robert Krug, Masoumeh Mansouri, and Federico Pecora. Motion planning and goal assignment for robot fleets using trajectory optimization. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 7939–7946. IEEE, 2018.

Bo Tang, Elias B Khalil, and Ján Drgoňa. Learning to optimize for mixed-integer non-linear programming. *arXiv preprint arXiv:2410.11061*, 2024.

Appendix A. Proof of Theorem 3

Proof We first prove (ii) given assuming that $\mathbf{x}_\rho^* \in S_v$. Let $S_v = \{\mathbf{x} \mid v(\mathbf{x}) = v^*\}$ be the set of all points that achieve this minimum violation. S_v is closed since it is a level set of a continuous function. Combining with the condition that either $f(\mathbf{x})$ is coercive or \mathcal{X} is compact, there exist a minimizer of $f(\mathbf{x})$ on S_v . From the definition of \mathbf{x}_ρ^* , we have

$$\mathbf{x}_\rho^* \in \arg \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) + \rho v(\mathbf{x}), \quad (17)$$

which means \mathbf{x}_ρ^* is also a minimizer of $f(\mathbf{x}) + \rho v(\mathbf{x})$ on S_v , i.e.,

$$\mathbf{x}_\rho^* \in \arg \min_{\mathbf{x} \in S_v} f(\mathbf{x}) + \rho v(\mathbf{x}) = \arg \min_{\mathbf{x} \in S_v} f(\mathbf{x}) + \rho v^*. \quad (18)$$

Since ρv^* is a constant, then $\mathbf{x}_\rho^* \in \arg \min_{\mathbf{x} \in S_v} f(\mathbf{x})$, or \mathbf{x}_ρ^* is a minimizer of $f(\mathbf{x})$ on S_v .

For (i), we prove that for a sufficiently large but finite ρ , \mathbf{x}_ρ^* must be in S_v . Let us assume, for the sake of contradiction, $\mathbf{x}_\rho^* \notin S_v$. From (17), we have

$$f(\mathbf{x}_\rho^*) + \rho v(\mathbf{x}_\rho^*) \leq f(\mathbf{x}^{**}) + \rho v(\mathbf{x}^{**}) = f(\mathbf{x}^{**}) + \rho v^*, \quad (19)$$

for any $\mathbf{x}^{**} \in S_v$, which leads to

$$\rho \leq \frac{f(\mathbf{x}^{**}) - f(\mathbf{x}_\rho^*)}{v(\mathbf{x}_\rho^*) - v^*}. \quad (20)$$

Next, since $v(\cdot)$ is a convex piecewise linear function, the Hoffman error bound property holds, i.e.,

$$v(\mathbf{x}_\rho^*) - v^* \geq c \operatorname{dist}(\mathbf{x}_\rho^*, S_v), \quad (21)$$

with $c > 0$. Moreover, $f(\mathbf{x})$ is a quadratic function, so it is Lipschitz continuous, and we have the following inequality

$$|f(\mathbf{x}^{**}) - f(\mathbf{x}_\rho^*)| \leq L_f \cdot \|\mathbf{x}^{**} - \mathbf{x}_\rho^*\|. \quad (22)$$

Therefore, if we choose $\mathbf{x}^{**} \in S_v$ such that $\operatorname{dist}(\mathbf{x}_\rho^*, S_v) = \|\mathbf{x}^{**} - \mathbf{x}_\rho^*\|$, i.e., \mathbf{x}^{**} is the closest point in S_v to \mathbf{x} , then from (19), we obtain that $\rho \leq \frac{L_f}{c}$ must be satisfied. Thus, we can select ρ such that $\rho > \frac{L_f}{c}$, which leads to a contradiction. The proof is thus complete. \blacksquare

Appendix B. Details of Numerical Examples

B.1. Collision Avoidance for Robot Navigation

In this example, we consider a navigation problem for a single robot operating in an environment with stationary obstacles. The robot is required to move from its initial position to a designated

goal while avoiding collisions with obstacles. This problem is formulated as an MIQP, where binary variables are used to formulate the collision avoidance constraints between the robot and the obstacles. The corresponding MI-MPC formulation follows the setup described in (Le et al., 2025).

We consider an MPC problem with a single robot and n_o obstacles, and let \mathcal{O} be the set of obstacles. We formulate the MPC problem with a control horizon of length H . Let $t \in \mathbb{Z}^+$ be the current time step. At every time step $k \in \mathbb{Z}^+$, let $\mathbf{p}(k) = [p^x(k), p^y(k)]^\top \in \mathbb{R}^2$, $\mathbf{v}(k) = [v^x(k), v^y(k)]^\top \in \mathbb{R}^2$, and $\mathbf{u}(k) = [u^x(k), u^y(k)]^\top \in \mathbb{R}^2$ be the vectors of positions, velocities, and accelerations for robot, respectively. Let $\mathbf{x}(k) = [\mathbf{p}(k), \mathbf{v}(k)]^\top$ be the state vector of robot. The dynamics of each robot are governed by a discrete-time double-integrator model as follows,

$$\begin{aligned} \mathbf{p}(k+1) &= \mathbf{p}(k) + \tau \mathbf{v}(k) + \frac{1}{2} \tau^2 \mathbf{u}(k), \\ \mathbf{v}(k+1) &= \mathbf{v}(k) + \tau \mathbf{u}(k), \end{aligned} \quad (23)$$

where $\tau \in \mathbb{R}_{>0}$ is the sampling time period, and compactly expressed as $\mathbf{x}(k+1) = \mathbf{f}(\mathbf{x}(k), \mathbf{u}(k))$. We assume that the states and control inputs of robots are subjected to the following bound constraints:

$$\begin{aligned} p_{\min}^x \leq p^x(k) \leq p_{\max}^x, \quad p_{\min}^y \leq p^y(k) \leq p_{\max}^y, \\ -v_{\max} \leq v^x(k), v^y(k) \leq v_{\max}, \quad -a_{\max} \leq u^x(k), u^y(k) \leq a_{\max}, \end{aligned} \quad (24)$$

where $[p_{\min}^x, p_{\max}^x, p_{\min}^y, p_{\max}^y]^\top \in \mathbb{R}^4$ is the boundary of the environment, $v_{\max} \in \mathbb{R}_{>0}$ and $a_{\max} \in \mathbb{R}_{>0}$ are the maximum speed and acceleration of the robots, respectively. More compactly, (24) is expressed as $\mathbf{x}(k) \in \mathcal{X}$ and $\mathbf{u}(k) \in \mathcal{U}$.

The mixed-integer constraints for collision avoidance between the robot and obstacle $o \in \mathcal{O}$ at time-step k are formulated by big-M formulation as follows,

$$\begin{aligned} \cos \alpha_o(p^x(k+1) - p_o^x) + \sin \alpha_o(p^y(k+1) - p_o^y) &\geq L_o + d_{\min} - M b_{1,o}(k), \\ -\sin \alpha_o(p^x(k+1) - p_o^x) + \cos \alpha_o(p^y(k+1) - p_o^y) &\geq W_o + d_{\min} - M b_{2,o}(k), \\ -\cos \alpha_o(p^x(k+1) - p_o^x) - \sin \alpha_o(p^y(k+1) - p_o^y) &\geq L_o + d_{\min} - M b_{3,o}(k), \\ \sin \alpha_o(p^x(k+1) - p_o^x) - \cos \alpha_o(p^y(k+1) - p_o^y) &\geq W_o + d_{\min} - M b_{4,o}(k), \end{aligned} \quad (25)$$

where $b_{1,o}(k)$, $b_{2,o}(k)$, $b_{3,o}(k)$ and $b_{4,o}(k)$ are binary decision variables satisfying

$$b_{1,o}(k) + b_{2,o}(k) + b_{3,o}(k) + b_{4,o}(k) \leq 3, \quad (26)$$

$[p_o^x, p_o^y]^\top$ is the center location, α_o is the rotation angle, and $2L_o$ and $2W_o$ are the length and width of obstacle $o \in \mathcal{O}$, respectively, while d_{\min} is the minimal distance between robot and obstacle to be considered as no collision. We define the binary decision variables $\delta(k) \in \{0, 1\}^{4n_o}$ at each time step k as the concatenated vector of $b_{1,o}(k)$, $b_{2,o}(k)$, $b_{3,o}(k)$, $b_{4,o}(k)$, for all $o \in \mathcal{O}$, and rewrite all the collision avoidance constraints as $\mathbf{g}_o(\mathbf{x}_{k+1}, \delta_k) \leq 0$.

The objective for the robot is to reach the goal, i.e., minimize the distance to the goal, while maintaining the minimum effort. Thus, we consider the following MPC cost given by a weighted sum of terminal cost \bar{c} and running cost c over the horizon i.e.,

$$\begin{aligned} \underset{\substack{\mathbf{x}(k+1) \in \mathcal{X}, \\ \mathbf{u}(k) \in \mathcal{U}}}{\text{minimize}} \quad & \bar{c}(\mathbf{x}(t+H)) + \sum_{k=t}^{t+H-1} c(\mathbf{u}(k), \mathbf{x}(k)), \end{aligned} \quad (27)$$

where

$$\begin{aligned}\bar{c}(\mathbf{x}(t+H)) &= \omega_{\text{pt}} \|\mathbf{p}(t+H) - \mathbf{p}_g\|_2^2, \\ c(\mathbf{u}(k), \mathbf{x}(k)) &= \omega_p \|\mathbf{p}(k) - \mathbf{p}_g\|_2^2 + \omega_u \|\mathbf{u}(k)\|_2^2\end{aligned}\quad (28)$$

with \mathbf{p}_g being the vector of goal positions, while ω_{pt} , ω_p , and ω_u are positive weights. Consequently, the cost function is quadratic in the continuous decision variables.

Therefore, the parametric MI-MPC problem can be given by

$$\begin{aligned}\text{minimize} \quad & \bar{c}(\mathbf{x}_H; \boldsymbol{\theta}) + \sum_{k=0}^{H-1} c(\mathbf{x}_k, \mathbf{u}_k; \boldsymbol{\theta}) \\ \text{subject to} \quad & \mathbf{x}_0 = \mathbf{x}_{\text{init}}(\boldsymbol{\theta}), \\ & \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k), \\ & \mathbf{g}_o(\mathbf{x}_{k+1}, \boldsymbol{\delta}_k) \leq 0, \\ & \mathbf{x} \in \mathcal{X}^{H+1}, \mathbf{u} \in \mathcal{U}^H, \boldsymbol{\delta}_k \in \{0, 1\}^{4n_o}.\end{aligned}\quad (29)$$

where the parameter vector $\boldsymbol{\theta} = [\mathbf{x}_0^\top, \mathbf{p}_g^\top]^\top \in \mathbb{R}^6$ contains the initial state and goal position. We consider an MPC problem with three obstacles and a control horizon of length 20, leading to $3 \times 20 \times 4 = 240$ binary decision variables. An NN $\pi_\omega : \mathbb{R}^6 \rightarrow \{0, 1\}^{240}$ is employed to predict the binary decision variables. The parameters for the simulation are set as follows: $\tau = 0.25$ s, $[p_{\min}^x, p_{\max}^x, p_{\min}^y, p_{\max}^y]^\top = [-0.5 \text{ m}, 3 \text{ m}, -3 \text{ m}, 0.5 \text{ m}]^\top$, $v_{\max} = 0.5$ m/s, $a_{\max} = 0.5$ m/s², $d_{\min} = 0.25$ m, $M = 10^3$, $\omega_{\text{pt}} = 10$, $\omega_p = 1$, $\omega_u = 1$, $w_s = 10^4$. The information of the three obstacles is:

1. Obstacle 1: $p_1^x = 1.0$ m, $p_1^y = 0.0$ m, $L_1 = 0.8$ m, $W_1 = 1.0$ m, $\alpha_1 = 0.0$ rad.
2. Obstacle 2: $p_2^x = 0.7$ m, $p_2^y = -1.1$ m, $L_2 = 1.0$ m, $W_2 = 0.8$ m, $\alpha_2 = 0.0$ rad.
3. Obstacle 3: $p_3^x = 0.4$ m, $p_3^y = -2.5$ m, $L_3 = 0.8$ m, $W_3 = 1.0$ m, $\alpha_3 = 0.0$ rad.

B.2. Simplified Thermal Energy Tank System

In the second example, we examine a simplified thermal energy tank system, adapted from (Boldock et al., 2025). The system dynamics are represented by a discrete-time linear time-invariant (LTI) model with both continuous and discrete control inputs, described as follows:

$$\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}_u\mathbf{u}_k + \mathbf{B}_d\delta_k + \mathbf{E}\mathbf{d}_k, \quad (30)$$

where $\mathbf{x}_k \in \mathbb{R}^2$ is the state vector, $\mathbf{u}_k \in \mathbb{R}^2$ is the continuous control input, $\delta_k \in \{0, 1, 2, 3\}$ is the discrete control input, and $\mathbf{d}_k \in \mathbb{R}^2$ represents known disturbances at time step k . The dynamics matrices \mathbf{A} , \mathbf{B}_u , \mathbf{B}_d , and \mathbf{E} are:

$$\mathbf{A} = \begin{bmatrix} 0.9983 & 0.001 \\ 0 & 0.9966 \end{bmatrix}, \quad \mathbf{B}_u = 0.075 \mathbb{I}_2, \quad \mathbf{B}_d = \begin{bmatrix} 0 \\ 0.0825 \end{bmatrix}, \quad \mathbf{E} = -0.0833 \mathbb{I}_2. \quad (31)$$

The system is subject to the following state and input constraints:

$$0 \leq x_{1,k} \leq 8.4, \quad 0 \leq x_{2,k} \leq 3.6, \quad 0 \leq u_{1,k}, u_{2,k} \leq 8. \quad (32)$$

In addition to the constraints on the continuous control inputs, we also impose constraints on the changes between consecutive time steps of the discrete control input as follows,

$$-1 \leq \delta_k - \delta_{k-1} \leq 1, \quad k = 1, \dots, H-1 \quad (33)$$

The control objective is to minimize the expected cumulative cost over the prediction horizon, which includes both state tracking and control effort penalties. The stage cost function and terminal cost function are defined as

$$\begin{aligned} \bar{c}(\mathbf{x}_k, \mathbf{u}_k, \delta_k; \mathbf{r}_k) &= \|\mathbf{x}_k - \mathbf{r}_k\|_{\mathbf{Q}}^2 + \|\mathbf{u}_k\|_{\mathbf{R}}^2 + \rho \|\delta_k\|_2^2, \\ c(\mathbf{x}_H; \mathbf{r}_H) &= \|\mathbf{x}_H - \mathbf{r}_H\|_{\mathbf{Q}_t}^2, \end{aligned} \quad (34)$$

where \mathbf{r}_k is the reference state at time step k , and \mathbf{Q} , \mathbf{R} , \mathbf{Q}_t and ρ are weighting matrices and vectors, respectively, given by $\mathbf{Q} = \mathbf{Q}_t = \mathbb{I}_2$, $\mathbf{R} = 0.5 \mathbb{I}_2$, and $\rho = 0.1$. For simplicity, we assume that the reference state \mathbf{r}_k remains constant over the prediction horizon, i.e., $\mathbf{r}_k = \mathbf{r}$ for all k , where $\mathbf{r} = [4.2, 1.8]^\top$. Thus, the optimization problem is parameterized by the current state \mathbf{x}_t and the sequence of known future disturbances over the prediction horizon, $[\mathbf{d}_k, \mathbf{d}_{k+1}, \dots, \mathbf{d}_{k+H-1}]$. Accordingly, the parameter vector is defined as $\boldsymbol{\theta} = [\mathbf{x}_k^\top, \mathbf{d}_k^\top, \mathbf{d}_{k+1}^\top, \dots, \mathbf{d}_{k+H-1}^\top] \in \mathbb{R}^{2H+2}$.

The parametric MI-MPC problem is given by

$$\begin{aligned} \text{minimize} \quad & \bar{c}(\mathbf{x}_H; \mathbf{r}_H) + \sum_{k=0}^{H-1} c(\mathbf{x}_k, \mathbf{u}_k, \delta_k; \mathbf{r}_k) \\ \text{subject to} \quad & \mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}_u\mathbf{u}_k + \mathbf{B}_d\delta_k + \mathbf{E}\mathbf{d}_k, \\ & -1 \leq \delta_k - \delta_{k-1} \leq 1, \quad k = 1, \dots, H-1 \\ & \mathbf{x}_k \in \mathcal{X}, \quad \mathbf{u}_k \in \mathcal{U}, \quad \delta_k \in \{0, 1, 2, 3\}. \end{aligned} \quad (35)$$

An NN is trained to predict the integer control inputs, i.e., $\boldsymbol{\pi}_\omega : \mathbb{R}^{2H+2} \rightarrow \{0, 1, 2, 3\}^H$. For this example, we set the control horizon to $H = 20$, resulting in a parameter vector $\boldsymbol{\theta} \in \mathbb{R}^{42}$ and an NN output dimension of 20.