

# CoCo: Online Mixed-Integer Control via Supervised Learning

Abhishek Cauligi<sup>1</sup>, Preston Culbertson<sup>2</sup>, Edward Schmerling<sup>1</sup>  
 Mac Schwager<sup>1</sup>, Bartolomeo Stellato<sup>3</sup>, and Marco Pavone<sup>1</sup>

**Abstract**—Many robotics problems, from robot motion planning to object manipulation, can be modeled as mixed-integer convex programs (MICPs). However, state-of-the-art algorithms are still unable to solve MICPs for control problems quickly enough for online use and existing heuristics can typically only find suboptimal solutions that might degrade robot performance. In this work, we turn to data-driven methods and present the Combinatorial Offline, Convex Online (CoCo) algorithm for quickly finding high quality solutions for MICPs. CoCo consists of a two-stage approach. In the offline phase, we train a neural network classifier that maps the problem parameters to a *logical strategy*, which we define as the discrete arguments and relaxed big-M constraints associated with the optimal solution for that problem. Online, the classifier is applied to select a candidate logical strategy given new problem parameters; applying this logical strategy allows us to solve the original MICP as a convex optimization problem. We show through numerical experiments how CoCo finds near optimal solutions to MICPs arising in robot planning and control with 1 to 2 orders of magnitude solution speedup compared to other data-driven approaches and solvers.

## I. INTRODUCTION

Planning and control using MICP has been an extensive area of study within the robotics community. MICPs can be used as a modeling framework to capture the rich set of behaviors and logical constraints that arise in problems such as planning for systems with contact [1], [2], motion planning [3], [4], and dexterous manipulation [5]. Despite their popularity, MICPs have rarely been put into practice for real-world control tasks with demanding performance requirements of 10-100Hz operational rates due to computational constraints. Indeed, the tremendous strides made in accelerating MICP solution times by several orders of magnitude in the past few decades rely on multithreaded implementations that are inapplicable on embedded systems commonly found on robot hardware. Thus, although algorithms such as branch-and-bound [6] provide certificates of optimality for MICPs, finding the optimal solution can be challenging in practice due to the  $\mathcal{NP}$ -hard nature of solving MICPs.

Alternate techniques used to make MICPs amenable for real-time control include terminating branch-and-bound when a feasible solution is first found or by simply rounding fractional integer solutions. However, such methods can degrade robot performance if a poor quality feasible solution

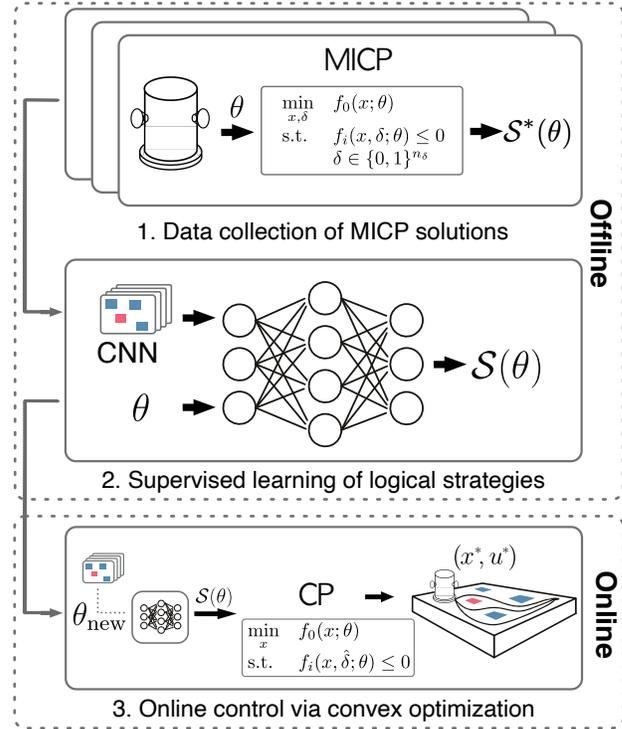


Fig. 1: Our algorithm CoCo is a data-driven approach that seeks to accelerate finding high-quality solutions to parametrized MICPs. The approach proposes a logical strategy  $\hat{S}(\theta)$ , which is a candidate discrete solution  $\hat{\delta}$  satisfying the logical constraints of the system. Given  $\hat{\delta}$ , the MICP can be approximately solved as a convex optimization problem. Here, we show CoCo applied to the free-flying spacecraft robot motion planning problem. (1) Offline, CoCo solves a set of MICPs for a representative set of planning problems and constructs the optimal logical strategy  $S^*(\theta)$  using the discrete optimizer  $\delta^*$  and set of relaxed constraints  $\mathcal{T}_M(\theta)$ . (2) Thereafter, a convolutional neural network classifier is trained to learn a mapping between problem parameters  $\theta$  and the logical strategy  $S^*(\theta)$ . (3) Online, this classifier predicts the logical strategy  $\hat{S}(\theta)$  associated with new parameters  $\theta$  and uses the candidate discrete solution  $\hat{\delta}$  to solve convex optimization problems until a feasible solution is found.

is returned. A promising approach that has emerged in recent years is to apply techniques from machine learning to accelerate finding solutions for numerical-optimization based robot controllers [7]. Although melding supervised learning techniques and MICP-control has been considered [8], [9], these approaches have yet to demonstrate the ability to scale to the large number of discrete decision variables or problem parameters typically found in robotics. Further, these approaches are agnostic to how the discrete decision variables are utilized (e.g., piecewise affine constraints, mixed logical dynamics) and do not take into consideration the logical structure of the problem in their solution approach.

In this work, we seek to develop a learning-based methodology for applying MICP-based control with the following desiderata in mind:

<sup>1</sup>Department of Aeronautics and Astronautics, Stanford University. {acauligi, schmerling, schwager, pavone}@stanford.edu

<sup>2</sup>Department of Mechanical Engineering, Stanford University. pculbertson@stanford.edu

<sup>3</sup>Department of Operations Research and Financial Engineering, Princeton University. bstellato@princeton.edu

This work was supported in part by NASA under the NASA Space Technology Research Fellowship Grants NNX16AM78H and 80NSSC18K1180.

- 1) *Performant control*: The controller should be able to find high quality solutions with respect to some performance metric for the control task.
- 2) *Speed*: The online solution procedure should be capable of providing real-time decision making.
- 3) *Generalizability*: Discrete decision variables are used to encode a rich set of behaviors in robotics and we seek an approach that can leverage the underlying structure present in many robot control tasks.
- 4) *Scalability*: The approach should be capable of solving MICPs with a high dimensional parameter space and 10s-100s of discrete decision variables.

*Related work*: The use of data-driven methods for quickly solving numerical optimization-based controllers is a nascent area of research. In [7], [10], supervised learning is used for learning warm starts for a quadratic program (QP)-based controller. Non-convex optimization-based controllers are considered as well, with [11], [12] learning warm starts for a sequential quadratic program (SQP)-based trajectory optimization library. The authors in [13] leverage differentiable convex optimization to develop a learning-based approach for tuning a QP-controller. However, the shortcoming of these preceding approaches is that they are limited to continuous optimization problems.

In comparison, there is a paucity of approaches that have investigated the use of data-driven techniques for accelerating integer program solutions in control [14]. For general discrete optimization problems, a popular approach has been to consider branch-and-bound as a sequential decision making problem and apply reinforcement learning approaches [15]. The shortcoming of such techniques for control is that they still rely on solving branch-and-bound online and often require multiple neural network forward passes at each node.

More recently, both [8], [5] propose a supervised learning approach using a neural network to warm start the binary variable assignments for an mixed-integer quadratic program (MIQP) controller and a k-nearest neighbors approach is proposed in [9]. In these approaches, the supervised framework maps problem parameters to a candidate discrete solution and, by fixing the discrete decision variables to the candidate solution values, the MIQP is solved as a QP. However, these approaches do not address scalability to a high dimensional parameter space or large number of discrete decision variables. Further, their solution strategy of directly proposing a candidate discrete solution without considering how the variables are employed in the MICP limits their applicability to higher dimensional parametric programs often found in robotics. Our proposed approach accommodates a broad set of systems for which MICPs are used as a modeling tool, such as mixed logical dynamical systems and piecewise affine (PWA) systems.

Our work draws inspiration from the field of explicit model predictive control (MPC) [16]. Unlike standard multiparametric optimization approaches that learn a solution map corresponding to polyhedral partitions of the parameter space, our approach learns strategies for regions of a nonlinear transformation (learned via a neural network) of the parameter space of the problem.

*Statement of Contributions*: Towards filling the gaps in existing works, we present the Combinatorial Offline, Convex Online (CoCo) framework. CoCo is a data-driven framework for solving MICPs and entails a two stage approach (as detailed in Figure 1) consisting of an offline and online phase. We introduce the concept of logical strategies and show how they enable for CoCo to find high quality solutions for a broad class of problems modeled as MICPs, including problems with 100s of binary decision variables and a large input parameter dimension.

To the best of our knowledge, CoCo uniquely satisfies the four aforementioned desiderata. This paper extends a conference publication [17] and provides the following additional contributions:

- 1) Demonstration of how the notion of task-specific logical strategies can be exploited to solve problems with a varying number of discrete variables.
- 2) Additional numerical results to compare CoCo against a commercial MICP solver and benchmark data-driven methods used to find feasible solutions for MICPs.
- 3) A thorough analysis of how CoCo can be deployed in practical situations and used across a variety of tasks.

## II. TECHNICAL BACKGROUND

This section introduces the parametrized MICPs studied in this work, the big-M constraint formulation approach, and the concept of well-posedness of MICP solutions.

### A. Parametrized MICPs

This work considers discrete optimization problems of the specific form known as parametrized mixed-integer convex programs. Given problem parameters  $\theta$ , we can define the following parametrized MICP with continuous decision variables  $x \in \mathbf{R}^{n_x}$ , and binary decision variables  $\delta \in \{0, 1\}^{n_\delta}$ :

$$\begin{aligned} & \underset{x, \delta}{\text{minimize}} && f_0(x; \theta) \\ & \text{subject to} && f_i(x, \delta; \theta) \leq 0, \quad i = 1, \dots, m_f \\ & && h_i(\delta; \theta) \leq 0, \quad i = 1, \dots, m_I \\ & && \delta \in \{0, 1\}^{n_\delta}. \end{aligned} \quad (1)$$

The objective function  $f_0$  and inequality constraints  $f_i$  are assumed convex with respect to  $x$ . The purely integer constraints  $h_i$  are assumed linear with respect to  $\delta$ .

We note here that the binary decision variables  $\delta$  are the “complicating” variables in the sense that (1) becomes much easier to solve if  $\delta$  are temporarily held fixed and the resulting convex program solved in terms of  $x$ . Indeed, if an optimal discrete solution  $\delta^*$  for (1) is provided, then the continuous optimizer  $x^*$  for the MICP can be easily found by solving a single convex optimization problem,

$$\begin{aligned} & \underset{x}{\text{minimize}} && f_0(x; \theta) \\ & \text{subject to} && f_i(x, \delta^*; \theta) \leq 0, \quad i = 1, \dots, m_f. \end{aligned}$$

Thus, we see that identifying an optimizer  $\delta^*$  for (1) allows a user to quickly find the optimal solution for an MICP and circumvent, e.g., exploring a full branch-and-bound tree. To this end, a host of supervised learning approaches from [8], [9], [18], [19], [20] utilize this insight to generate a candidate  $\hat{\delta}$  given problem parameters  $\theta$  for an MICP.

## B. Big-M Formulation

In MICPs, binary variables  $\delta$  are often introduced in conjunction with what is known as big-M formulation to capture high-level discrete or logical behavior of the system. These big-M constraints then enforce the desired high-level behavior on the continuous variables  $x$  for the robot task at hand, e.g., contact task assignment or hybrid control logic [21].

As an example of a desired logical behavior, consider enforcing the constraint,

$$f_1(x; \theta) \leq 0 \vee f_2(x; \theta) \leq 0, \quad (2)$$

where  $\vee$  indicates an ‘‘or’’ relationship. The big-M method can be used to encapsulate this constraint by introducing an auxiliary term,

$$M_i(\theta) := \sup_x f_i(x; \theta), \quad i = 1, 2$$

where  $M_1(\theta)$  is the least upper bound on the value attainable by  $f_1(x; \theta)$  and  $M_2(\theta)$  the least upper bound for  $f_2(x; \theta)$ . Then, the logical behavior can be enforced through the introduction of two binary variables  $\delta_1$  and  $\delta_2$ ,

$$f_1(x; \theta) \leq M_1(\theta)(1 - \delta_1) \quad (3)$$

$$f_2(x; \theta) \leq M_2(\theta)(1 - \delta_2) \quad (4)$$

$$\delta_1 + \delta_2 \geq 1 \quad (5)$$

The equivalence between (2) and (3)-(5) can be verified by noting that when  $\delta_1 = 1$ , then  $f_1(x; \theta) \leq 0$  is automatically recovered. Similarly,  $\delta_2 = 1$  leads to the constraint  $f_2(x; \theta) \leq 0$ . Enforcing the behavioral constraint that one of the two constraints must be satisfied is accomplished through the final constraint (5).

We note that the right hand sides of (3)-(4) are linear in terms of  $\delta_i$  and thus denote a ‘‘big-M’’ constraint as an affine expression of  $\delta_i$ :

$$g_i(x; \theta) \leq a_i(\theta)(1 - \delta_i), \quad (6)$$

where  $a_i(\theta)$  are positive-valued constants precomputed using  $M(\theta)$ , or any upper bound for the constraint, to impose the desired logical behavior.

In this work, we use big-M constraints exclusively to relate the continuous variables  $x$  and binary variables  $\delta$ . If  $m_M$  is the number of big-M constraints used, we can write more specifically the class of MICPs studied as:

$$\begin{aligned} & \underset{x, \delta}{\text{minimize}} && f_0(x; \theta) \\ & \text{subject to} && f_i(x; \theta) \leq 0, \quad i = 1, \dots, m_f \\ & && g_i(x; \theta) \leq a_i(\theta)(1 - \delta_i), \quad i = 1, \dots, m_M \\ & && h_i(\delta; \theta) \leq 0, \quad i = 1, \dots, m_I \\ & && \delta \in \{0, 1\}^{n_\delta}. \end{aligned} \quad (\mathcal{P}(\theta))$$

## C. Well-Posedness of MICP Solutions

Here, we distinguish between two common classes of MICPs and will later show how this difference is crucial in formulating an effective solution approach. Specifically, as the inclusion of binary variables  $\delta$  in (1) renders the problem non-convex, an MICP may admit multiple globally optimal solutions  $(x^*, \delta^*)$ . This leads to the distinction between *well-posed* and *completely well-posed* MICPs as introduced in [21]. Well-posed MICPs admit a unique continuous minimizer  $x^*$ ,

but may admit multiple discrete optimizers  $\delta^*$ . Completely well-posed problems assume that an MICP admits a single global minimizer  $(x^*, \delta^*)$ .

Completely well-posed MICPs can be used to model many systems, such as those with PWA constraints. In the case of PWA constraints, a continuous solution  $x^*$  can only be attained by a particular set of mode transitions that are uniquely encoded by a single discrete optimizer  $\delta^*$ . More broadly however, the assumption of a unique discrete optimizer is a limiting one and cannot accommodate many applications. For example, mixed logical dynamical systems [21] enforce logic rules on discrete-time dynamical systems with both continuous and discrete decision variables, but allow for multiple discrete optimizers  $\{\delta^*\}$  that satisfy the propositional logic constraints. Thus, we assume that the problems we treat are well-posed MICPs. Although this entails assuming that the MICP admits a single  $x^*$ , we note that this is a mild assumption in practice. For example, if problem parameters  $\theta$  yield multiple  $x^*$  for  $\mathcal{P}(\theta)$ , a small perturbation in the initial condition typically breaks ties and leads to a single minimizer  $x^*$ . Indeed, in many robotics problems, slight regularization terms can be added to the objective function to ensure a unique  $x^*$  while still accomplishing the control task.

Given a continuous optimizer  $x^*$ , the set of binary optimizers  $\{\delta^*\}$  for a well-posed MICP can be characterized by identifying which big-M constraints of the form (6) are *relaxed*, where a constraint  $g_i(x^*; \theta)$  is denoted relaxed if:

$$g_i(x^*; \theta) > 0. \quad (7)$$

That is, a big-M constraint is considered relaxed if, after plugging in the values from  $x^*$ , the constraint (6) is satisfied if and only if the binary variable  $\delta_i^*$  attains value 0. <sup>1</sup> If  $\delta$  is subject to purely integer constraints (e.g., a cardinality constraint), then the purely integer constraints are also included in identifying the list of relaxed constraints. In contrast, a big-M constraint is considered *enforced* if the values for  $x^*$  automatically lead to constraint satisfaction, i.e.,  $g_i(x^*; \theta) \leq 0$  which allows for both cases of  $\delta_i^* = 0$  or  $\delta_i^* = 1$  (subject to purely integer constraints). Thus, the set  $\{\delta^*\}$  is then comprised of  $\delta^*$  that attain value  $\delta_i^* = 0$  for the set of relaxed constraints and only differ in the values associated with the enforced constraints.

As an example, we consider a continuous optimizer  $x^*$  that satisfies the logical expression from (2) for which  $f_1(x^*; \theta) \leq 0$ , but  $f_2(x^*; \theta) > 0$ . In this case,  $f_2(x^*; \theta)$  is a relaxed constraint, while  $f_1(x^*; \theta)$  is an enforced constraint. This value of  $x^*$  then admits two binary optimizers  $(\delta_1^*, \delta_2^*) = (0, 1)$  and  $(\delta_1^*, \delta_2^*) = (1, 1)$ . We see here that as  $f_2(x^*; \theta)$  is the relaxed constraint,  $\delta_2^*$  attains value 1 for both possible binary optimizers.

## III. TECHNICAL APPROACH

This section defines logical strategies, presents the subsequent development of task-specific logical strategies, and demonstrates how they can be used for solving MICPs.

<sup>1</sup>We stress here that a relaxed constraint does not connote a relaxation of the binary variable  $\delta_i$  over its convex hull,  $\delta_i \in [0, 1]$ .

### A. Logical Strategies

MICPs are used to encode a rich set of logical constraints and behaviors for robotic systems and existing approaches seek to tackle this broad class of problems by identifying an optimal discrete solution  $\delta^*$  associated with problem parameters  $\theta$ . For well-posed MICPs, a solution approach that only proposes a single candidate binary optimizer  $\hat{\delta}$  given  $\theta$  leads to an ill-posed supervised learning problem, as there may exist a set of target values  $\{\delta^*\}$  given a  $\theta$ .

To address this shortcoming, we present the notion of logical strategies, named as such because they take into consideration how binary variables  $\delta$  are used to enforce high-level logical behavior in well-posed MICPs. A logical strategy  $\mathcal{S}(\theta)$  leverages the idea that, while there may exist multiple discrete optimizers given a continuous optimizer  $x^*$ , the problem  $\mathcal{P}(\theta)$  can be solved by identifying the set of relaxed big-M constraints.

Given the preceding definition of a relaxed big-M constraint from (7), we now define a logical strategy. Given problem parameters  $\theta$  and continuous optimizer  $x^*$  for the problem  $\mathcal{P}(\theta)$ , let  $\delta^*(\theta)$  be a particular binary optimizer from the set of discrete optimizers  $\{\delta^*(\theta)\}$  and  $\mathcal{T}_M(\theta)$  be the set of relaxed big-M constraints for the problem,

$$\mathcal{T}_M(\theta) = \{i \mid g_i(x^*; \theta) > 0\}. \quad (8)$$

We then define the logical strategy  $\mathcal{S}(\theta)$  as a tuple  $(\delta^*(\theta), \mathcal{T}_M(\theta))$ .

The utility of the logical strategy definition becomes apparent when revisiting the original MICP in (1). If the optimal logical strategy  $\mathcal{S}^*(\theta)$  is provided for an MICP, then the continuous optimizer  $x^*$  for the MICP can be found by solving a single convex optimization problem,

$$\begin{aligned} & \underset{x}{\text{minimize}} && f_0(x, \delta^*; \theta) \\ & \text{subject to} && f_i(x; \theta) \leq 0, \quad i = 1, \dots, m_f \\ & && g_i(x; \theta) \leq 0, \quad i \in \mathcal{T}_M^c(\theta) \end{aligned}$$

where  $\mathcal{T}_M^c(\theta)$  is the set of enforced constraints, i.e., the complement of  $\mathcal{T}_M(\theta)$ . We note here that if  $\mathcal{P}(\theta)$  admits a set of discrete solutions  $\{\delta^*\}$  given a particular  $x^*$ , then any one of the  $\delta^*$  can be included in  $\mathcal{S}(\theta)$ . This leads to the insight that, rather than proposing a candidate  $\hat{\delta}(\theta)$ , a solution approach should rather propose a candidate logical strategy  $\hat{\mathcal{S}}(\theta)$ . As we demonstrate later, considering logical strategies improves performance of a supervised learning-based solution approach.

### B. Task-Specific Logical Strategies for MICPs

For large problems, proposing a candidate logical strategy starts to become intractable as the number of candidate logical strategies becomes too large. However, underlying problem structure allows us to consider the logical strategy associated with each constraint individually. To that end, we introduce the idea of task-specific logical strategies that consider a common structure arising in many robotics problems known as separability.

To demonstrate an example of separability, we consider the robot motion planning problem shown in Figure 1. Four binary variables are used to enforce obstacle avoidance, where each binary variable is associated with lying on one side of an axis-aligned rectangular obstacle at a particular time. Thus, the

binary variables  $\delta$  are decoupled on the basis of which obstacle they are associated with.

Task-specific logical strategies seek to exploit this problem structure. Formally, the underlying mixed logical constraints can be written as a conjunction of Boolean formulas:

$$F_1 \wedge F_2 \wedge \dots \wedge F_\ell,$$

where  $F_i$  is a distinct sub-formula of literals involving continuous and binary variables, and each binary variable  $\delta_j$  is associated with only one sub-formula  $F_i$ . When the mixed logical constraint consists of  $\ell$  such Boolean formulas, then the logical strategy  $\mathcal{S}(\theta)$  can itself be split into  $\ell$  sub-formula strategies  $\mathcal{S}_1(\theta), \dots, \mathcal{S}_\ell(\theta)$ . We note that this decomposition is only possible because of the separability in the problem that allows for each  $\delta_j$  to be considered only in relation to a single sub-formula  $\mathcal{S}_i(\theta)$  and the value of the continuous solution  $x^*$ .

An example of such separability arising in a robotics problem is the MICP formulation of collision avoidance constraints [3]. Consider an axis-aligned, 2D rectangular obstacle  $m$  that is parametrized by the coordinates of its lower-left hand corner  $(x_{\min}^m, y_{\min}^m)$  and upper right-hand corner  $(x_{\max}^m, y_{\max}^m)$ . If the 2D position of the robot is  $p = (p^1, p^2) \in \mathbf{R}^2$ , then the collision avoidance constraints with respect to obstacle  $m$  are:

$$x_{\max}^m - M\delta^{m,1} \leq p^1 \leq x_{\min}^m + M\delta^{m,2} \quad (9)$$

$$y_{\max}^m - M\delta^{m,3} \leq p^2 \leq y_{\min}^m + M\delta^{m,4} \quad (10)$$

where  $M$  is chosen to be a sufficiently large number. As written in (15) and (16),  $\delta^{m,i} = 1$  indicates that robot is on one side of face  $i$  of the obstacle and in violation of that keep-out constraint. To ensure that the robot does not collide with obstacle  $m$ , a final constraint

$$\sum_{i=1}^4 \delta^{m,i} \leq 3, \quad (11)$$

is enforced. Note that each binary variable depends only on the three other variables associated with the same obstacle. Thus, each task-specific logical strategy considers only the binary variables and big-M constraints for obstacle  $m$ .

Task-specific logical strategies offer several advantages for a supervised learning approach to solving parametrized MICPs. First, as the number of possible logical strategies can grow exponentially in terms of the number of binary variables  $n_\delta$ , considering each sub-formula strategy separately  $\mathcal{S}_i(\theta)$  reduces the number of binary variables in each task-specific logical strategy, thereby resulting in a smaller number of values that each sub-formula can attain. Second, this reduced number of candidate sub-formulas leads to improved supervision in a learning-based approach as the number of target values is reduced. Finally, in the case when additional Boolean formulas are added (i.e., additional binary variables are added to the MICP), the sub-formula strategies can be queried at inference time for these new formulas.

Thus, task-specific logical strategies can lead to improved performance in problems with separability by reducing the number of class labels and thereby improving supervision for a data-driven approach. However, as we show, the separability of the constraints has performance limits, especially in applications where additional sub-formulas strategies are queried at test time. Practically, in the context of robotics,

---

### Algorithm 1: CoCo Offline

---

**Require:** Batch of training data  $\{\theta_i\}_{i=1,\dots,T}$ , problem  $\mathcal{P}(\theta)$

- 1: Initialize strategy dictionary  $\mathcal{S} \leftarrow \{\}$ , train set  $\mathcal{D} \leftarrow \{\}$
- 2:  $k \leftarrow 0$
- 3: **for** each  $\theta_i$  **do**
- 4:   Solve  $\mathcal{P}(\theta)$
- 5:   **if**  $\mathcal{P}(\theta)$  is optimal **then**
- 6:     Construct optimal strategy  $\mathcal{S}^*$
- 7:     **if**  $\mathcal{S}^*$  not in  $\mathcal{S}$  **then**
- 8:       Add  $\mathcal{S}^*$  to  $\mathcal{S}$
- 9:       Assign class label  $y_k$  to  $\mathcal{S}^*$
- 10:        $k \leftarrow k + 1$
- 11:     **end if**
- 12:     Identify class label  $y_i$  for strategy class  $\mathcal{S}^*$
- 13:     Add  $(\theta_i, y_i)$  to  $\mathcal{D}$
- 14:   **end if**
- 15: **end for**
- 16: Choose network weights  $\phi$  which minimize cross-entropy loss  $\mathcal{L}(h_\phi(\theta_i), y_i)_{i=1,\dots,T}$  via stochastic gradient descent
- 17: **return**  $h_\phi, \mathcal{S}$

---

the performance of the controller can be assessed beforehand and only deployed for tasks that are sufficiently similar to the test set (e.g., a maximum number of obstacles for which task-specific logical strategies are queried).

#### IV. COMBINATORIAL OFFLINE, CONVEX ONLINE

We now present our proposed approach CoCo, short for Combinatorial Offline, Convex Online. CoCo consists of a two-stage approach for training and deploying a neural network classifier that maps problem parameters  $\theta$  to a candidate logical strategy  $\hat{\mathcal{S}}(\theta)$ .

Algorithm 1 details the offline portion of CoCo. The algorithm takes as input a set of problem parameters  $\{\theta_i\}$ , where each  $\theta_i$  is sampled from a parameter distribution  $p(\Theta)$  representative of the problems encountered in practice. We refer to  $\mathcal{S}$  as the strategy dictionary and  $\mathcal{D}$  as the train set, both of which are initially empty (Line 1). The strategy dictionary  $\mathcal{S}$  stores the set of logical strategies  $\{\mathcal{S}^{(i)}\}$  constructed during the offline phase and the train set  $\mathcal{D}$  stores the set training tuples  $\{(\theta_i, y_i)\}$  used for training the neural network classifier, where  $y_i$  is the class label associated with logical strategy  $\mathcal{S}^{(i)}$ . For each  $\theta_i$ , the MICP  $\mathcal{P}(\theta)$  is solved (Line 4). If an optimal solution is found to the MICP, then the primal solution  $(x^*, \delta^*)$  is used to construct the logical strategy  $\mathcal{S}^*$  for this problem (Line 6) and added to the strategy dictionary  $\mathcal{S}$  if it is not already included (Lines 7-11). The class label  $y_i$  associated with  $\mathcal{S}^*$  is identified (Line 12) and the tuple  $(\theta_i, y_i)$  added to  $\mathcal{D}$  (Line 13). Finally, a neural network classifier  $h_\phi$  with output dimension  $|\mathcal{S}|$  is trained using the elements of  $\mathcal{D}$  and the weights  $\phi$  are chosen in order to approximately minimize a cross-entropy loss over the training samples (Line 16).

#### V. NUMERICAL EXPERIMENTS

In this section, we compare CoCo with commercial solvers and other data-driven approaches for solving MICPs. We present results on three benchmark problems in robotics that are modeled as MICPs: the control of an underactuated cart-pole with multiple contacts, dexterous manipulation for task-specific grasping, and the robot motion planning problem.

---

### Algorithm 2: CoCo Online

---

**Require:** Problem parameters  $\theta$ , strategy dictionary  $\mathcal{S}$ , trained neural network  $h_\phi, n_{\text{evals}}$

- 1: Compute class scores  $h_\phi(\theta)$
- 2: Identify top  $n_{\text{evals}}$ -scoring strategies in  $\mathcal{S}$
- 3: **for**  $j = 1, \dots, n_{\text{evals}}$  **do**
- 4:   **if**  $\mathcal{P}(\theta)$  is feasible for strategy  $\mathcal{S}^{(j)}$  **then**
- 5:     **return** Feasible solution  $(x^*, \delta^*)$
- 6:   **end if**
- 7: **end for**
- 8: **return** failure

---

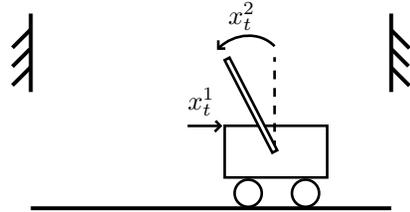


Fig. 2: 4D cart-pole with wall system.

#### A. Implementation Details

For each system, we first generate a dataset by sampling  $\theta$  from  $p(\Theta)$ , until a sufficient number of problems  $\mathcal{P}(\theta)$  are solved. For each system, we separate 90% of the problems for training and the remaining 10% for evaluation. For the cart-pole and dexterous manipulation problems, the neural network architecture consists of a standard ReLU feedforward network with three layers and 32 neurons per layer. For the free-flyer system, we used a CNN architecture with four convolutional layers followed by a feedforward network with three layers and 128 neurons per layer.

We implemented each example in Python and used the PyTorch machine learning library to implement our neural network models with the ADAM optimizer for training. The MICPs were written using the cvxpy modeling framework and solved using Mosek. We further benchmark CoCo against the commercial solver, the regression framework from [8], and the k-nearest neighbors framework from [9]. We disable presolve and multithreading to better approximate the computational resources of an embedded processor. The network architecture chosen for the regressor was identical to the CoCo classifier, updated with the appropriate number of integer outputs. The code for our algorithm is available at <https://github.com/StanfordASL/CoCo>.

#### B. Cart-Pole with Soft Walls

We first study the cart-pole with wall system shown in Figure 2, a well-known underactuated, multi-contact problem in robot control [1], [2]. The system consists of a cart and pole and the optimal control problem entails regulating the system from initial state  $x_0$  to a goal  $x_g$ . The non-convexity of the problem stems from four binary variables  $\delta_t \in \mathbf{R}^4$  introduced at each time step to enforce the logical constraint that the contact force from the wall only becomes active when the tip of the pole makes contact with either wall. The parameter space  $\theta \in \mathbf{R}^8$  for this problem is comprised of the initial state  $x_0 \in \mathbf{R}^4$  and goal state  $x_g \in \mathbf{R}^4$ . We refer the reader to the Appendix for a full derivation of system constraints.

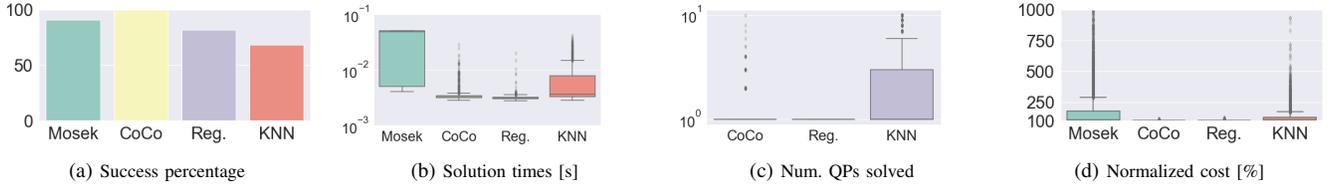


Fig. 3: For the cart-pole system, CoCo finds near-global solutions for a majority of problems.

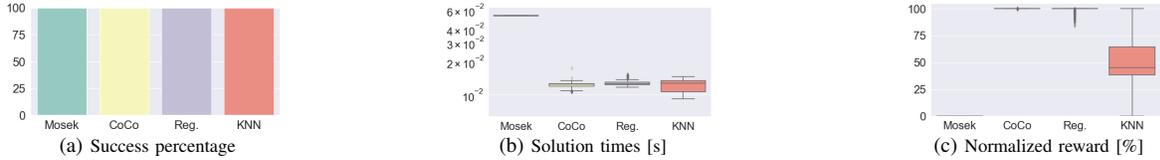


Fig. 4: Simulation results for manipulation example. CoCo reduces solution times for (10) between 2–3 orders of magnitude.

## 1) Results

Numerical results for the cart-pole system are given in Figure 3. We set the horizon  $N$  to the value 10, resulting in a total of 40 binary variables. The training set consists of 90 thousand problems generated using parameters sampled from the parameter distribution  $p(\Theta)$  and evaluation metrics presented for a test set of ten thousand problems. Figure 3a reports the percent of feasible solutions found over the test set. For the commercial solver, branch-and-bound is timed out after 50ms and, for CoCo, ten candidate logical strategies are evaluated before the algorithm terminates with failure. Computation times shown in Figure 3b include the inference step to generate a candidate binary solution (i.e. the forward pass of the network, nearest neighbor lookup, etc.) plus solution time for solving convex relaxations before a feasible solution was found. Figure 3c and Figure 3d report the number of convex relaxations solved per problem and the cost of the feasible solution relative to the globally optimal solution, respectively. Note that that Figure 3c does not include the number of convex relaxations solved by Mosek as the `cvxpy` interface does not provide this information.

We see that CoCo outperforms the commercial solver Mosek and the two other benchmarks. As shown in Figure 3a, CoCo finds feasible solutions for 99% of the problems, compared to 82% and 69% for the regressor and KNN, respectively. Mosek (timed out at 50ms) finds feasible solutions for 91% of the test set and only 44% of these solutions correspond to the globally optimal solution. As Figures 3c to 3d show, CoCo finds the globally optimal solution after one QP solve for 98% of its feasible solutions.

### C. Task-Oriented Optimization of Dexterous Grasps

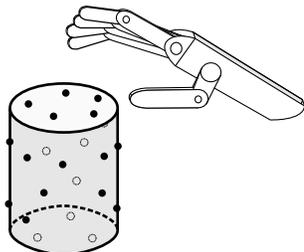


Fig. 5: Schematic of dexterous grasping problem. Here, a robotic hand with  $n = 5$  fingers chooses from  $M$  potential contact points to optimize a task-specific grasp metric.

The next problem considered is that of grasp optimization for task-specific dexterous grasping shown in Figure 5. Dexterous grasping with multi-fingered hands is a challenging problem due to both the number of contact modes that must be assigned for a stable grasp and because the resulting grasp must be able to execute the desired task under consideration. Task-agnostic grasp optimization problems generally entail solving challenging non-convex optimization problems, making them prohibitively expensive for applications replanning [22].

Thus, we are interested in enabling online computation of optimal dexterous grasps for fast replanning and regrasping. Specifically, we focus our attention on the problem of task-specific grasping where the grasp sequence is chosen with a particular task such as pushing or rotating a tool in mind. Task-specific grasp optimization can be posed as an MICP where the binary variables indicate which fingers are in contact and the objective function is the grasp metric chosen to capture the quality of a grasp for a particular task [23].

Specifically, we consider the problem of choosing  $n$  contact points for a multifingered robot hand from a set of points  $p_1, \dots, p_M \in \mathbb{R}^3$ , sampled from the object surface in order to optimize the task-oriented grasp metric from [23]. The “control” actions optimized for are the local contact forces  $f_i \in \mathbb{R}^3$  such that they also satisfy *friction cone* and grasp matrix constraints. We then use binary variables  $\delta$  to ensure that contact forces are not applied at all candidate points and enforce the constraint,

$$f_i^z \leq f_i^{z,\max} \delta_i,$$

where  $f_i^z$  is the z-axis component of the local contact force.

We use the task-based grasp metric introduced in [23] and consider *task wrenches*  $\hat{F}_t$ , which are specific directions in wrench space that characterize the applied wrenches necessary to complete the task. A task can then be expressed by a set of wrenches which must be generated and this set can be characterized as the positive span of  $T$  task vectors. Given a task wrench  $\hat{F}_t$ , let  $\alpha_t$  be its associated grasp quality and  $w_t \geq 0$  its task weighting. Then, the grasp quality metric from the  $T$  task vectors is given by,

$$\mu(\delta, \hat{F}_1, \dots, \hat{F}_T) = \sum_{t=1}^T w_t \alpha_t.$$

This grasp metric corresponds to the volume of the polyhedron defined by the vectors  $w_t \alpha_t \hat{F}_t$  and can be computed by solving  $T$  second-order cone programs (SOCPs). Thus, this problem is a mixed-integer second-order cone program (MISOCP) with

$M$  binary variables and the parameter vector  $\theta \in \mathbf{R}^{12}$  consists of the desired weights for a task vector  $\hat{F}_t$  which correspond to the basis vectors  $\pm e_i \in \mathbf{R}^6$  for  $i = 1, \dots, 6$ .

### 1) Results

The numerical experiments consist of planning grasps for a four finger manipulator  $n = 4$ . We consider a set of  $M$  candidate grasp points, with  $M$  equal to 30, for a single rigid body. The training set consists of 4,500 problems and the weights  $w_i$ , where  $w_i > 0$ , are generated by calculating the softmax of a vector sampled from a multivariate normal distribution with covariance matrix  $\Sigma = 10\mathbf{I}$ .

Figure 4 show the results for this system. As the primary point of comparison, we compare the optimality of the feasible solutions found and computation time. Indeed, we see in Figure 4c that CoCo finds the globally optimal grasp for 99% of the problems while maximizing the grasp metric is challenging for the benchmarks. As any grasp mode sequence with four contacts leads to a feasible solution for the problem, we see in Figures 4b to 4c that timing out Mosek at 50ms leads to highly suboptimal solutions for this particular problem. Moreover, finding a high quality solution after solving only one SOCP relaxation leads to solution times on the order of tens of milliseconds for CoCo. Thus, CoCo allows for high quality feasible solutions for MISOCPs, whereas a commercial solver would lead to highly suboptimal, low quality grasp solutions.

### D. Free-Flying Space Robots

A fundamental problem in robotics that is inherently combinatorial is that of motion planning in the presence of obstacles. Here, we study a free-flying spacecraft robot that must navigate around obstacles on a planar workspace with linear dynamics. We show how the use of task-specific logical strategies allows for (1) this problem to become tractable for application of CoCo and (2) the learned strategy classifier  $\hat{h}_\phi$  to be used at test time for MICPs with a different number of binary variables  $n_\delta$  than from the training set.

The system state  $x_t \in \mathbf{R}^4$  consists of the position  $p_t \in \mathbf{R}^2$  and velocity  $v_t \in \mathbf{R}^2$ . The planning problem is to regulate the robot from an initial state  $x_0$  to a goal state  $x_g$  while satisfying dynamics and actuator constraints. The crucial constraint that renders the problem non-convex is the safety constraint  $x_t \in \mathcal{X}_{\text{safe}}$ , as  $\mathcal{X}_{\text{safe}}$  is typically a highly non-convex region of the workspace and solving the planning problem requires a global combinatorial search. For the free-flying space robot, a popular approach to solving the motion planning problem has been to pose it as an MICP [3], [24]. In this formulation, given  $N_{\text{obs}}$  obstacles, the workspace is first decomposed into keep-in and keep-out zones and binary variables  $\delta$  used to enforce collision avoidance with the keep out regions. Due to the  $\ell_2$ -norm constraints imposed on the thruster forces, this problem is a mixed-integer quadratically constrained quadratic program (MIQCQP) with  $4N_{\text{obs}}N$  binary variables. The parameters  $\theta$  for this problem include the initial state  $x_0$ , goal state  $x_g$ , and position of obstacles  $\{(x_{\min}^m, y_{\min}^m, x_{\max}^m, y_{\max}^m)\}_{m=1}^{N_{\text{obs}}}$ .

### 1) Use of Task-Specific Logical Strategies

We show here how task-specific strategy decompositions are necessary to efficiently solve the free-flyer motion planning problem with CoCo and how the underlying structure of the problem can be leveraged in order to do so.

From inspection of (15) and (16), we note that a binary variable  $\delta_t^{m,i}$  appears in constraints only with the other three binary variables for obstacle  $m$  at time  $t$ . Thus, our key insight here is to decompose the logical strategy on a per obstacle basis. That is, a logical strategy  $\mathcal{S}(\theta; m)$  is associated with the constraint used to enforce collision avoidance with obstacle  $m$  and encodes information about the binary variable assignment for that obstacle  $\{\delta_t^m\}_{t=1}^N$ . Rather than training  $N_{\text{obs}}$  separate classifiers for each obstacle, we train a single classifier with  $\theta$  and append an encoding to specify which obstacle logical strategy is being queried.

Specifically, this is accomplished in this work using the architecture shown in Figure 1. First, a synthetic image of the obstacles is generated using obstacle coordinates  $(x_{\min}^m, y_{\min}^m)$  and  $(x_{\max}^m, y_{\max}^m)$ . The strategy  $\mathcal{S}(\theta; m)$  being queried is indicated by coloring in obstacle  $m$  with a different color in the image. A convolutional pass is then computed over the synthetic image and the output is flattened and appended with the remaining problem parameters  $\theta$  before being input to a fully connected feedforward network.

For inference, a batch of input images and problem parameters are constructed, a single forward pass computed, and the strategies for the corresponding sub-formula ranked. One issue with decomposing the strategy queries is that the full binary variable assignment  $\delta$  must be reconstructed from the individual  $\delta_m$ . In this work, we consider the  $n_{\text{evals}}$  highest scoring logical strategy candidates for each  $\mathcal{S}(\theta; m)$  and this leads to  $n_{\text{evals}}^{N_{\text{obs}}}$  candidates for  $\delta$ . As this can be a prohibitively large number of solution candidates, we instead randomly sample  $m_{\text{evals}}$  (where  $m_{\text{evals}} \ll n_{\text{evals}}^{N_{\text{obs}}}$ ) enumerations from the set of  $n_{\text{evals}}^{N_{\text{obs}}}$  candidate assignments for  $\delta$ . Additionally, we ensure that the  $\delta_m$  corresponding to the highest scoring candidate  $\mathcal{S}(\theta; m)$  for each obstacle is included in this set of  $m_{\text{evals}}$  candidate binary solutions.

The results for comparing task-specific strategies used in CoCo are shown in Figure 6, where we also include the Machine Learning Optimizer (MLOPT) from [18] as an additional benchmark. We note immediately in Figure 6a that using task-specific logical strategies with CoCo leads to performance gains compared to the approach used for MLOPT, with 92% feasible solutions found for CoCo compared to only 8% for MLOPT. This disparity is attributable to the fact that CoCo encodes 458 logical strategies versus approximately 67,000 for MLOPT over the 90,000 training problems, leading to a sparser set of training labels per class for MLOPT compared to CoCo. We also see in Figures 6b to 6d that, although Mosek finds a slightly higher number of feasible solutions, CoCo finds the globally optimal solution for 90% of the problems for which a feasible solution is found and at twice the speed as Mosek. Further, we see in Figure 6a that the regressor and KNN benchmarks fare poorly on the logical constraints and find feasible solutions for only 49% of the test set.

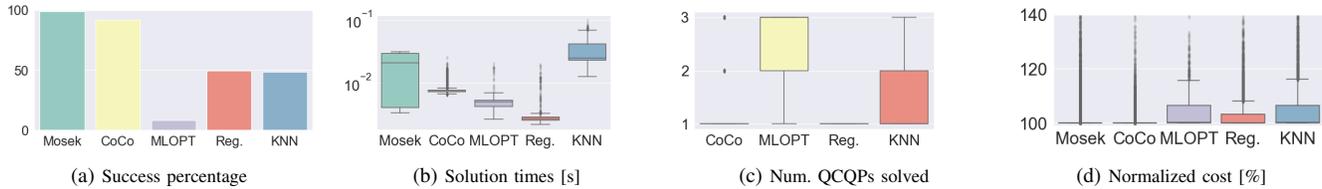


Fig. 6: Simulation results for the free-flyer show how task-specific strategies are necessary for using CoCo.

## 2) Generalization

One important consequence of using a convolutional pass to query the strategy sub-formula  $\mathcal{S}(\theta; m)$  is that it can be used for inference in problems with a different number of  $N_{\text{obs}}$  than from the training set. Here, we evaluate the ability of CoCo to generalize to a distribution of problems with a varying number of binary obstacles. We train multiple strategy classifiers corresponding to horizons of  $N = \{5, 7, 9, 11\}$  with environments of eight obstacles. Figure 7 shows the performance of applying these networks in solving problems with an increasing number of obstacles  $N_{\text{obs}} = \{6, \dots, 12\}$ . As shown, we see that the efficacy of the strategy classifier diminishes with an increasing horizon length  $N$  due to a corresponding increase in the number of strategies. Intuitively, we also see that performance decreases with an increase in obstacles simply due to the increased difficulty of the planning problem. However, we note that the performance dropoff remains roughly linear rather than an exponential decrease of performance stemming from including additional binary variables. Thus, given knowledge of the robot operating environment, the use of a trained classifier can be limited to scenarios in which the number of obstacles does not lead to a dramatic dropoff in performance.

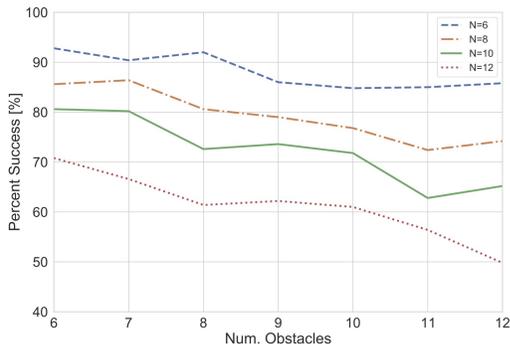


Fig. 7: Generalization results for a CoCo strategy classifier network demonstrates roughly linear dropoff in performance as obstacles are added. The strategy classifier network was trained on an environment with  $N_{\text{obs}} = 8$  for various horizons  $N$ .

## 3) Timeout Performance

Here, we compare CoCo with a commercial solver that is timed out with a prespecified termination time (i.e., the incumbent solution from branch-and-bound is returned). For CoCo, we terminate CoCo either when a feasible solution is found or after the termination time has been exceeded. Figure 8 compares the percent of feasible solutions found between CoCo and Mosek. We see that CoCo finds a feasible solution for the majority of problems within about five milliseconds, which is approximately the time required to compute a forward pass

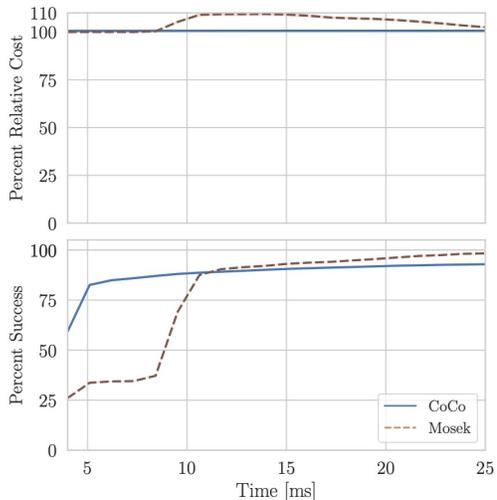


Fig. 8: We compare the performance of CoCo (blue) and Mosek (brown) after timing out both on a test set of MIQCQPs. We see in the bottom plot that timing out Mosek leads to a reduced number of feasible solutions found compared with CoCo. As the cutoff time increases, we see in the top plot that if Mosek finds a feasible solution, then these feasible solutions are suboptimal compared to CoCo, which generally finds the globally optimal solution even when timed out early.

of the CNN and solve a single convex relaxation. However, Mosek requires twice the computation time before it finds a feasible solution for the majority of problems.

Moreover, we see that the solutions found by CoCo are effectively the globally optimal solution for that problem, whereas the Mosek's incumbent solution from branch-and-bound is often suboptimal until branch-and-bound terminates. Thus, in tasks where Mosek is allowed to run its full course, Mosek will indeed find the globally optimal solution, but a designer can weigh the tradeoffs between quickly finding a high quality feasible solution using CoCo or allowing Mosek to terminate. Finally, we further note that in applications requiring a certificate of optimality, a feasible solution found by CoCo can be used as the incumbent and provide a tighter upper bound for branch-and-bound.

## VI. CONCLUSION

In this work, we presented CoCo, a data-driven framework to find high quality feasible solutions for MICPs used in robot planning and control problems. We demonstrated how problem structure arising in robot tasks can be utilized effectively in a supervised learning framework. Specifically, we introduced the notion of task-relevant logical strategies to exploit such problem structure and showed how they improve the performance of the trained strategy classifier. We showed through numerical experiments that CoCo improves solution speeds by 1-2 orders of magnitude with only a slight loss of optimality, compared to low-quality feasible solutions found by the commercial solver and benchmark data-driven approaches. Finally, we showed

how CoCo can uniquely be used to solve problems with a varying number of discrete decision variables and how this allows for solving a new set of tasks online. To this end, we believe that a promising direction of work is to improve the classifier performance given new tasks online. One future approach could be to explore meta-learning to allow for CoCo to adapt network parameters online for improved performance using information gained from solving problems online.

## REFERENCES

- [1] R. Deits, T. Koolen, and R. Tedrake, "LVIS: Learning from value function intervals for contact-aware robot controllers," in *Proc. IEEE Conf. on Robotics and Automation*, 2019.
- [2] T. Marcucci and R. Tedrake, "Warm start of mixed-integer programs for model predictive control of hybrid systems," *IEEE Transactions on Automatic Control*, 2020.
- [3] B. Landry, R. Deits, P. R. Florence, and R. Tedrake, "Aggressive quadrotor flight through cluttered environments using mixed integer programming," in *Proc. IEEE Conf. on Robotics and Automation*, 2016.
- [4] P. Culbertson, S. Bandyopadhyay, and M. Schwager, "Multi-robot assembly sequencing via discrete optimization," in *IEEE/RSJ Int. Conf. on Intelligent Robots & Systems*, 2019.
- [5] F. R. Hogan, E. R. Grau, and A. Rodriguez, "Reactive planar manipulation with convex hybrid MPC," in *Proc. IEEE Conf. on Robotics and Automation*, 2018.
- [6] J. Lee and S. Leyffer, Eds., *Mixed Integer Nonlinear Programming*. Springer-Verlag, 2012.
- [7] X. Zhang, M. Bujarbaruah, and F. Borrelli, "Safe and near-optimal policy learning for model predictive control using primal-dual neural networks," in *American Control Conference*, 2019.
- [8] D. Masti and A. Bemporad, "Learning binary warm starts for multi-parametric mixed-integer quadratic programming," in *European Control Conference*, 2019.
- [9] J.-J. Zhu and G. Martius. (2019) Fast non-parametric learning to accelerate mixed-integer programming for online hybrid model predictive control. Available at <https://arxiv.org/pdf/1911.09214.pdf>.
- [10] S. W. Chen, T. Wang, N. Atanasov, V. Kumar, and M. Morari. (2019) Large scale model predictive control with neural networks and primal active sets. Available at <https://arxiv.org/pdf/1910.10835.pdf>.
- [11] Z. Wang, T. Taubner, and M. Schwager, "Multi-agent sensitivity enhanced iterative best response: A real-time game theoretic planner for drone racing in 3D environments," *Robotics and Autonomous Systems*, vol. 125, 2020.
- [12] G. Tang, W. Sun, and K. Hauser, "Learning trajectories for real-time optimal control of quadrotors," in *IEEE/RSJ Int. Conf. on Intelligent Robots & Systems*, 2018.
- [13] A. Agrawal, S. Barratt, S. Boyd, and B. Stellato, "Learning convex optimization control policies," in *Learning for Dynamics & Control*, 2019.
- [14] Y. Bengio, A. Lodi, and A. Prouvost. (2018) Machine learning for combinatorial optimization: a methodological tour d'Horizon.
- [15] H. Dai, E. B. Khalil, Y. Zhang, B. Dilkina, and L. Song, "Learning combinatorial optimization algorithms over graphs," in *Conf. on Neural Information Processing Systems*, 2017.
- [16] D. Malyuta and B. Açikmeşe, "Approximate multiparametric mixed-integer convex programming," *IEEE Control Systems Letters*, vol. 4, no. 1, pp. 157–162, 2019.
- [17] A. Cauligi, P. Culbertson, B. Stellato, D. Bertsimas, M. Schwager, and M. Pavone, "Learning mixed-integer convex optimization strategies for robot planning and control," in *IEEE CDC*, 2020, in Press.
- [18] D. Bertsimas and B. Stellato. (2019) Online mixed-integer optimization in milliseconds. Available at <https://arxiv.org/abs/1907.02206>.
- [19] B. Karg and S. Lucia, "Deep learning-based embedded mixed-integer model predictive control," in *European Control Conference*, 2018.
- [20] Y. Löhr, M. Klaučo, M. Fikar, and M. Mönnigmann, "Machine learning assisted solutions of mixed integer MPC on embedded platforms," in *IFAC World Congress*, 2020.
- [21] A. Bemporad and M. Morari, "Control of systems integrating logic, dynamics, and constraints," *Automatica*, 1999.
- [22] C. Ferrari and J. Canny, "Planning optimal grasps," in *Proc. IEEE Conf. on Robotics and Automation*, 1992.
- [23] R. Haschke, J. J. Steil, I. Steuwer, and H. Ritter, "Task-oriented quality measures for dextrous grasping," in *Proc. IEEE Int. Symp. on Computational Intelligence in Robotics and Automation*, 2005.
- [24] T. Schouwenaars, B. De Moor, E. Feron, and J. How, "Mixed integer programming for multi-vehicle path planning," in *European Control Conference*, 2001.

## APPENDIX

We review the system dynamics and constraints for each MICP studied in this work.

### A. Cart-Pole with Soft Walls

As depicted in Figure 2, the system consists of a cart and pole and the optimal control problem entails regulating the system to a goal  $x_g$ :

$$\begin{aligned}
 & \underset{x_0:N, u_0:N-1, \delta}{\text{minimize}} && \|x_N - x_g\|_2 + \sum_{\tau=0}^{N-1} \|x_\tau - x_g\|_2 + \|u_\tau\|_2 \\
 & \text{subject to} && x_{t+1} = Ax_t + Bu_t + Gs_t, \quad t = 0, \dots, N-1 \\
 & && u_{\min} \leq u_t \leq u_{\max}, \quad t = 0, \dots, N-1 \\
 & && s_t = \begin{cases} \kappa\lambda_t + \nu\gamma_t & \text{if } \lambda_t \geq 0 \text{ and } \kappa\lambda_t + \nu\gamma_t \geq 0 \\ 0, & \text{otherwise} \end{cases} \\
 & && t = 0, \dots, N-1 \\
 & && x_{\min} \leq x_t \leq x_{\max}, \quad t = 0, \dots, N \\
 & && x_0 = x_{\text{init}} \\
 & && \delta \in \{0, 1\}^{4 \times (N-1)},
 \end{aligned} \tag{12}$$

where the state  $x_t \in \mathbf{R}^4$  consists of the position of the cart  $x_t^1$ , angle of the pole  $x_t^2$ , and their derivatives  $x_t^3$  and  $x_t^4$ , respectively. The force applied to the cart is  $u_t \in \mathbf{R}$  and  $s_t \in \mathbf{R}^2$  are the contact forces imparted by the two walls. The relative distance of the tip of the pole with respect to the left and right walls is  $\lambda_t \in \mathbf{R}^2$  and the time derivative of this relative distance  $\gamma_t \in \mathbf{R}^2$ . Finally,  $\kappa$  and  $\nu$  are parameters associated with the soft contact model used.

As the contact force  $s_t$  becomes active only when the tip of the pole makes contact with either wall, we must introduce binary variables to enforce the logical constraints given in (12). We denote the relative distance of the tip of the pole with respect to the left and right walls as  $\lambda_t^1$  and  $\lambda_t^2$ , respectively:

$$\begin{aligned}
 \lambda_t^1 &= -x_t^1 + \ell x_t^2 - d \\
 \lambda_t^2 &= x_t^1 - \ell x_t^2 - d,
 \end{aligned}$$

where  $\ell$  is the length of the pole and  $d$  half the distance between the walls. The time derivatives of  $\lambda_t^1$  and  $\lambda_t^2$  are

$$\begin{aligned}
 \gamma_t^1 &= -x_t^3 + \ell x_t^4 \\
 \gamma_t^2 &= x_t^3 - \ell x_t^4.
 \end{aligned}$$

To constrain contact forces  $s_t$  to become active only when the pole tip strikes a wall, we introduce four binary variables  $\delta_t^i$ ,  $i = 1, \dots, 4$ . Using the formulation from [21], we enforce the following constraints for  $k = 1, 2$ :

$$\begin{aligned}
 \lambda_{\min}^k (1 - \delta_t^{(2k-1)}) &\leq \lambda_t^k \leq \lambda_{\max}^k \delta_t^{(2k-1)} \\
 s_{\min}^k (1 - \delta_t^{(2k)}) &\leq \kappa\lambda_t^k + \nu\gamma_t^k \leq s_{\max}^k \delta_t^{(2k)}
 \end{aligned}$$

Finally, we impose constraints on  $s_t$ , for  $k = 1, 2$ :

$$\nu\gamma_{\max}^k (\delta_t^{(2k-1)} - 1) \leq s_t^k - \kappa\lambda_t^k - \nu\gamma_t^k \leq s_{\min}^k (\delta_t^{(2k)} - 1)$$

There are then a total of  $4N$  integer variables in this MIQP.

### B. Task-Oriented Optimization of Dexterous Grasps

The task-specific dexterous grasping problem entails choosing  $n$  contact points out of  $M$  in order to optimize the task-oriented grasp metric from [23]. Each contact point  $p_i \in \mathbf{R}^3$  is sampled from the object surface and contacts between the finger and the object are modeled as point contacts with friction. The force that can be applied by each finger in the

local contact force is  $f_i = (f_i^x, f_i^y, f_i^z) \in \mathbf{R}^3$ , where the local coordinate frame has the  $x$ - and  $y$ -axes tangent to the surface, and the  $z$ -axis along the inward surface normal. Intuitively,  $f_i^z$  is the component of the contact force which is normal to the object surface, and  $f_i^x, f_i^y$  are its tangential components.

Under this contact model, we constrain the contact force  $f_i$  to lie within the *friction cone*  $\mathcal{K}^{(i)}$ . Let us further define the contact force vector  $f = (f_1, \dots, f_M) \in \mathbf{R}^{3M}$ , which is the vector of all contact forces.

Using the definition of a grasp matrix from [22], we can express the wrench applied to the object (from all contact forces) as  $Gf$ , where  $G = [G_1, \dots, G_M]$ .

However, contact forces may not be applied at all candidate points. To this end, we introduce the logical variables  $\delta_i \in \{0, 1\}$ , with  $\delta_i = 1$  iff point  $p_i$  is selected for the grasp. Thus, we enforce the constraint

$$f_i^z \leq \delta_i,$$

which constrains the normal forces of all unused grasps to be zero, and to be bounded by unity otherwise. Thus, for a choice of grasps  $\delta = (\delta_1, \dots, \delta_M)$ , the set of possible object wrenches is defined as  $\mathcal{W}(\delta) = \{Gf \mid f \in \mathcal{K}, f_i^z \leq \delta_i\}$ .

In [23], the authors propose a task-based grasp metric using *task wrenches*  $\hat{F}_t$ , which are specific directions in wrench space that characterize the applied wrenches necessary to complete the task. For instance, if the desired task is to push the object along the  $+x$ -axis, then this task could be described using  $\hat{F} = (1, 0, 0, 0, 0, 0)$ , and so on. For a task described by a single wrench, the grasp quality can be defined as

$$\mu_1(\delta, \hat{F}_t) = \sup \left\{ \alpha \geq 0 \mid \alpha \hat{F}_t \in \partial \mathcal{W}(\delta) \right\},$$

where  $\partial \mathcal{W}$  denotes the boundary of  $\mathcal{W}$ .

However, most tasks are best described by a set of wrenches which must be generated, rather than a single direction in wrench space. Thus, the authors propose describing this set as the positive span of  $T$  task vectors; in turn, the grasp metric is defined as

$$\mu(\delta, \hat{F}_1, \dots, \hat{F}_T) = \sum_{t=1}^T w_t \alpha_t,$$

where  $w_i \geq 0$  are the relative weightings of the task vectors, and  $\alpha_t = \mu_1(\delta, \hat{F}_t)$ . This can, in turn, be computed by solving  $T$  SOCPs.

We seek  $\delta^*$  which maximizes this grasp metric, which yields a MISOCP:

$$\begin{aligned} & \text{maximize}_{\alpha, f_{1:M}, \delta} && \sum_{t=1}^T w_t \alpha_t \\ & \text{subject to} && Gf^t = \alpha_t \hat{F}_t, \quad t = 1, \dots, T \\ & && f_i^t \in \mathcal{K}^{(i)}, \quad i = 1, \dots, M, \quad t = 1, \dots, T \\ & && f_i^{z,t} \leq \delta_i, \quad i = 1, \dots, M, \quad t = 1, \dots, T \\ & && \sum_{i=1}^M \delta_i \leq n \\ & && \delta \in \{0, 1\}^M \end{aligned} \quad (13)$$

### C. Free-Flying Space Robots

We let  $p_t \in \mathbf{R}^2$  be the robot position and  $v_t \in \mathbf{R}^2$  the velocity in the 2-dimensional plane. The robot state is  $x_t = (p_t, v_t)$  and the input  $u_t \in \mathbf{R}^2$  consists of the forces produced by the thruster. Letting  $\mathcal{X}_{\text{safe}}$  be the free space which the robot

must navigate through, the optimal control problem is to plan a collision free trajectory towards a goal state  $x_g$ :

$$\begin{aligned} & \text{minimize}_{x_{0:N}, u_{0:N-1}, \delta} && \|x_N - x_g\|_2 + \sum_{\tau=0}^{N-1} \|x_\tau - x_g\|_2 + \|u_\tau\|_2 \\ & \text{subject to} && x_{t+1} = Ax_t + Bu_t, \quad t = 0, \dots, N-1 \\ & && \|u_t\|_2 \leq u_{\max}, \quad t = 0, \dots, N-1 \\ & && x_{\min} \leq x_t \leq x_{\max}, \quad t = 0, \dots, N \\ & && x_0 = x_{\text{init}} \\ & && x_t \in \mathcal{X}_{\text{safe}}, \quad t = 0, \dots, N \\ & && \delta \in \{0, 1\}^{4N_{\text{obs}} \times N}. \end{aligned} \quad (14)$$

The constraint  $x_t \in \mathcal{X}_{\text{safe}}$  is the primary constraint of interest as it renders the problem non-convex. In this work, we consider axis-aligned rectangular obstacles. An obstacle  $m$  is parametrized by the coordinates of its lower-left hand corner  $(x_{\min}^m, y_{\min}^m)$  and upper right-hand corner  $(x_{\max}^m, y_{\max}^m)$ . Given the state  $x_t$ , the collision avoidance constraints with respect to obstacle  $m$  at time  $t$  are:

$$x_{\max}^m - M\delta_t^{m,1} \leq x_t^1 \leq x_{\min}^m + M\delta_t^{m,2} \quad (15)$$

$$y_{\max}^m - M\delta_t^{m,3} \leq x_t^2 \leq y_{\min}^m + M\delta_t^{m,4} \quad (16)$$