Conformal Mixed-Integer Constraint Learning with Feasibility Guarantees

Daniel Ovalle

Department of Chemical Engineering Carnegie Mellon University dovallev@andrew.cmu.edu

Ignacio E. Grossmann

Department of Chemical Engineering Carnegie Mellon University igOc@andrew.cmu.edu

Lorenz T. Biegler

Department of Chemical Engineering Carnegie Mellon University 1b01@andrew.cmu.edu

Carl D. Laird

Department of Chemical Engineering Carnegie Mellon University claird@andrew.cmu.edu

Mateo Dulce Rubio

Center for Data Science New York University mateo.d@nyu.edu

Abstract

We propose Conformal Mixed-Integer Constraint Learning (C-MICL), a novel framework that provides probabilistic feasibility guarantees for data-driven constraints in optimization problems. While standard Mixed-Integer Constraint Learning methods often violate the true constraints due to model error or data limitations, our C-MICL approach leverages conformal prediction to ensure feasible solutions are ground-truth feasible with probability at least $1-\alpha$, under a conditional independence assumption. The proposed framework supports both regression and classification tasks without requiring access to the true constraint function, while avoiding the scalability issues associated with ensemble-based heuristics. Experiments on real-world applications demonstrate that C-MICL consistently achieves target feasibility rates, maintains competitive objective performance, and significantly reduces computational cost compared to existing methods. Our work bridges mathematical optimization and machine learning, offering a principled approach to incorporate uncertainty-aware constraints into decision-making with rigorous statistical guarantees.

1 Introduction

Constraint Learning (CL) deals with inferring the functional form of constraints or objectives from observed data to embed them into mathematical optimization problems, particularly when parts of the system behavior are difficult to model explicitly [1]. This is especially relevant in settings where the relationship between decision and output variables is unknown but can be learned from data, enabling a supervised learning approach [2, 3]. It is also valuable in surrogate modeling scenarios where the relationship is known but difficult to optimize, either due to its nonlinearity, computational cost, or complexity [4–6]. That is, CL simply leverages modern machine learning systems to learn the structure of complex constraints directly from data, enabling prediction of constraint boundaries at unobserved points to efficiently navigate feasible spaces and identify optimal solutions.

The core principle that enables CL is that the predictions of many machine learning models can be formulated as closed-form algebraic expressions that are optimization-compatible and commonly admit exact mixed-integer programming (MIP) representations [1, 7–17]. This integration is commonly referred to as *Mixed-Integer Constraint Learning* (MICL). However, a critical challenge arises from the very nature of optimization algorithms: they systematically leverage the mathematical structure of the problem to identify extreme points. As a result, small approximation errors in the learned constraint model (e.g., arising from noisy data, model bias, or insufficient training) can be systematically exploited during the optimization. This behavior often leads to solutions that are infeasible with respect to the true system, meaning that the optimal decision variables found by optimization programs under CL frequently violate the original constraints when evaluated using the true underlying model [18]. This lack of *true* feasibility results in solutions that cannot be implemented in practice, or that are overly conservative due to the need for additional safety margins.

To improve the robustness of MICL solutions, recent research has proposed ensemble-based approaches where multiple predictive models are trained on different bootstrap samples of the data [19]. In particular, Maragno et al. [7] introduce a framework in which P models are independently trained on bootstrapped datasets, and the optimization is constrained so that a $1-\alpha$ fraction of these models must satisfy the imposed constraints. There is some empirical evidence that suggests that this heuristic yields solutions that are typically feasible with respect to the true underlying system, and its robustness improves as P increases. However, this approach lacks formal guarantees as there is no assurance that the final solution will satisfy the true constraints with high probability, particularly in the presence of model misspecification or low-quality data used in the MICL step. Moreover, the method suffers from key scalability issues: as P grows, the size of the optimization problem increases accordingly due to the replication of model constraints and auxiliary variables, rapidly making the problem computationally intractable for large ensembles.

In this paper, we introduce *Conformal Mixed-Integer Constraint Learning* (C-MICL), a novel model-agnostic framework for incorporating data-driven constraints into optimization problems with formal probabilistic feasibility guarantees. Our method supports both regression and classification settings and is computationally efficient without sacrificing significant solution quality. Specifically:

- **Probabilistic Guarantees:** We provide a new optimization formulation that guarantees ground-truth feasibility with probability at least $1-\alpha$ when the true constraint function is unknown or inaccessible, under a mild conditional independence assumption.
- **Model-Agnostic and MIP-Compatible:** C-MICL is compatible with any predictive model that admits a MIP encoding, enabling integration into standard optimization solvers.
- Scalability and Efficiency: Our method avoids query access or model refinement during
 optimization, does not scale with dataset size, and requires training at most two models,
 achieving orders-of-magnitude computational speedup over ensemble-based heuristics.
- Empirical Validation: We demonstrate the effectiveness of C-MICL on real-world case studies in both regression and classification settings. Extensive experiments show that our method consistently achieves target feasibility levels, while delivering competitive objective values with significantly reduced computation times compared to state-of-the-art baselines.

2 Literature Review

Mixed-Integer Constraint Learning. The mixed-integer constraint learning framework has enabled a wide range of developments in machine learning, including neural network verification [20–27], adversarial example generation [28–30], and applications in reinforcement learning [31–33]. It has also been successfully applied to domain-specific tasks in power systems [34, 35], healthcare [36, 37], supply chain management [38, 39], and chemical process design [40]. For a comprehensive overview of recent advances in constraint learning methodologies, we refer readers to Fajemisin et al. [1], and for a review focused on engineering applications, see Misener and Biegler [41]. However, existing CL frameworks generally lack guarantees of true system feasibility, particularly in settings where predictive uncertainty from noisy data or model misspecification is prevalent but typically ignored. The approach proposed in this work directly addresses this gap by providing formal probabilistic feasibility guarantees for uncertain learned constraints embedded into optimization problems.

Trust-Region-Filter for Optimization. Trust region filter methods have been proposed as a strategy to optimize over learned constraints while ensuring that the solution aligns with that of the true, underlying system [18]. These approaches operate by iteratively refining the learned constraint through local sampling and evaluation of the truth function, progressively narrowing the optimization domain within trust regions where the learned model is accurate [42, 43]. However, such methods rely on strong assumptions: they typically assume noise-free data, differentiability of the learned constraint, and direct query access to the truth function. In contrast, our framework is designed for settings where these assumptions do not hold. We account for noisy data, allow for non-differentiable learned models, and, most importantly, do not require access to or knowledge of the underlying function, enabling optimization with high-probability feasibility guarantees even under limited and imperfect information.

Conformal-Enhanced Optimization. Building on advances in conformal prediction [44–46], recent work has explored its integration with optimization. Lin et al. [47] proposed conformalized inverse optimization, targeting the estimation of unknown parameters embedded in the objective or constraints of an optimization problem, and providing uncertainty quantification over these learned parameters. Our approach addresses a different problem, namely learning and conformalizing the functional forms of constraints themselves, without relying on a fixed parametric structure, which allows for more expressive modeling of complex or unknown system behaviors. Separately, Zhao et al. [48] introduced conformalized chance constrained optimization, in which constraints involving random variables are required to hold with high probability. Their method encodes quantile thresholds from conformal prediction directly into the optimization formulation, which scales with the number of data points. Unlike this approach, ours is focused on uncertainty of model predictions rather than model stochasticity, and supports learning general constraint functions.

Other related efforts include Johnstone and Cox [49], who construct conformal ellipsoids for robust optimization, and Yeh et al. [50], who learn convex uncertainty sets for robust objectives using differentiable conformal layers. In contrast, C-MICL conformalizes feasibility regions rather than robustifying objective values. Similarly, Kiyani et al. [51] and Patel et al. [52] employ conformal sets for utility-aware or risk-sensitive decision-making, optimizing value-at-risk or ensuring robust output quality, whereas C-MICL focuses on constraint satisfaction under uncertainty from model estimation. Our framework supports general, nonconvex feasible regions and provides distribution-free probabilistic feasibility guarantees from fixed-size optimization formulations regardless of dataset size, allowing practitioners to leverage large datasets for improved model accuracy without inflating computational cost. A key distinction of C-MICL is that it centers on feasibility-aware decision-making rather than on robustifying objectives or actions for downstream optimization problems.

3 Preliminaries

3.1 Problem Set-Up

Consider a constrained optimization problem over a set of decision variables (x,y,z) where certain constraints depend on unknown or difficult-to-model functions relating decision variables x to outcomes y. We assume no access to the true underlying mapping denoted by $h(x) = \mathbb{E}[Y \mid X = x]$. Instead, we have access to a dataset $\mathcal{D}_{\text{train}} = \{(x_i,y_i)\}_{i=1}^{N_{\text{train}}}$, where each observation (x_i,y_i) is drawn independently from a common (unknown) underlying distribution \mathcal{P}_{XY} . Note that z are additional decision variables that may constrain x, but do not influence y through h(x).

The goal of Constraint Learning is to approximate h(x) with a predictive model $\hat{h}(x)$ trained on the observed data $\mathcal{D}_{\text{train}}$. Once the model $\hat{h}(x)$ is trained, we embed it in a mixed-integer optimization program, defining the following Mixed-Integer Constraint Learning (MICL) problem [1]:

$$\begin{aligned} & \min_{x,y,z} & & f(x,z) \\ & \text{s.t.} & & g(x,z) \leq 0 \\ & & \hat{h}(x) = y \\ & & (x,z) \in \mathcal{X} \\ & & y \in \mathcal{Y} \end{aligned} \tag{MICL}$$

where f(x,z) is a known real-valued objective function and g(x,z) encodes a set of known constraints. Moreover, $\mathcal{X} \subseteq \mathbb{R}^{m_1} \times \mathbb{Z}^{m_2}$ defines the feasibility set for the decision variables (x,z) where both x and z can take both real and integer values. Finally, \mathcal{Y} defines the feasible set for the outcome variable y, where feasibility must be satisfied with respect to the true (but unknown) function h(x).

The formulation above captures a general form of a mixed-integer constraint learning problem. Specific problem instances correspond to different definitions of the feasible prediction set \mathcal{Y} . In the regression setting, \mathcal{Y} is any subset of \mathbb{R} that is representable within a MIP framework. For instance, in our experiments, we define $\mathcal{Y} = [\underline{y}, \overline{y}]$, which enforces lower and upper bounds on the continuous output variable y. In the multi-label classification case, $\mathcal{Y} = \mathcal{K}^{\text{des}}$ is a subset of the set of possible labels $\mathcal{K} = \mathcal{K}^{\text{des}} \cup \mathcal{K}^{\text{und}}$, which we assume partitions into disjoint sets of *desired* classes \mathcal{K}^{des} and *undesired* classes \mathcal{K}^{und} .

By construction, all feasible solutions (x',z') from the optimization problem MICL are feasible with respect to the learned model $\hat{h}(x)$, meaning that $\hat{h}(x') \in \mathcal{Y}$. However, since h(x) is unknown and only approximated via $\hat{h}(x)$, ground-truth feasibility with respect to the true function h(x') is not guaranteed. We formalize this conceptual difference in the following definition.

Definition 3.1. (Ground-Truth Feasibility) Given a feasible solution (x', z') to the MICL optimization problem, the solution is said to achieve ground-truth feasibility if $y' = h(x') \in \mathcal{Y}$, where h(x) is the true constraint function.

Recent research by Maragno, Wiberg, Bertsimas, Birbil, den Hertog, and Fajemisin [7] proposes a heuristic approach called Wrapped model MICL (W-MICL) to address the challenge of ensuring ground-truth feasibility of MICL solutions. W-MICL trains an ensemble of P models using independent bootstrapped datasets rather than relying on a single predictive model. Each model $\hat{h}_p(x)$ provides a distinct prediction, collectively capturing model uncertainty. Within the MICL optimization formulation, the constraint $\hat{h}(x) \in \mathcal{Y}$ is replicated for each trained model, and feasibility $\hat{h}_p(x) \in \mathcal{Y}$ is enforced for a $(1-\alpha)$ fraction of them, implemented via standard big-M formulations in MIP [7]. While empirical results suggest that W-MICL improves the ground-truth feasibility of optimal solutions, the method remains a heuristic and does not offer formal feasibility guarantees.

Remark 3.1. We assume that the underlying predictive model $\hat{h}(x)$ admits a reformulation as a mixed-integer program (MIP), such as Neural Networks with ReLU activations (ReLU NNs) [8–10], Gradient-Boosted Trees (GBTs) [11, 12], Random Forests (RFs) [13–16], or Linear-Model Decision Trees (LMDTs) [17]. Then, the MICL problem is itself a MIP that can be solved using standard branch-and-bound or global optimization algorithms [53, 54]. While details on embedding these models into MIPs are provided in Appendix A, we emphasize that our formulation remains model-agnostic, requiring only that the predictive model $\hat{h}(x)$ is MIP-representable.

3.2 Conformal Prediction

Conformal prediction is a statistical framework for model-free uncertainty quantification that provides probabilistic guarantees for prediction sets constructed from any given (fixed) predictive model. This framework has gained significant popularity as an assumption-lean, model-agnostic, computationally efficient method offering rigorous statistical coverage guarantees without requiring model retraining or modifications to pre-trained machine learning systems [55–57].

Formally, suppose we have access to a *calibration dataset* of N data points $(x_1, y_1), \ldots, (x_N, y_N)$ drawn exchangeably from the underlying distribution \mathcal{P}_{XY} . Given a new data point (x_{N+1}, y_{N+1}) drawn exchangeably from \mathcal{P}_{XY} , conformal prediction provides a principled method for constructing a *conformal set* $\mathcal{C}_{\alpha}(x_{N+1})$ containing the *true* label y_{N+1} with probability at least $1-\alpha$, without making additional assumptions about the underlying distribution beyond exchangeability.

The core idea behind conformal prediction is to compute *conformal scores* $\{s(x_i, y_i)\}_{i=1}^N$, that capture model uncertainty over the calibration data. For instance, in regression settings, a common *conformal score function* corresponds to model absolute residuals $s(x, y) = |y - \hat{h}(x)|$. The

¹If y is a vector-valued outcome, $[y, \overline{y}]$ imposes elementwise constraints on $y \in \mathcal{Y}$.

²A random vector (Z_1, \ldots, Z_N) is exchangeable if $(Z_{\sigma(1)}, \ldots, Z_{\sigma(N)})$ follows the same distribution for all permutations σ of the indices $\{1, \ldots, N\}$ [57]. For instance, *i.i.d.* random variables are exchangeable.

conformal prediction set is then defined as:

$$C_{\alpha}(x_{N+1}) = \{ y : s(x_{N+1}, y) \le \widehat{q}_{1-\alpha} \}$$

$$\tag{1}$$

where $\widehat{q}_{1-\alpha}$ is the $(1-\alpha)(1+1/N)$ -quantile of the empirical distribution of conformal scores $\{s(x_1,y_1),\ldots,s(x_N,y_N)\}$. This construction yields the following key marginal guarantee.

Theorem 3.1 (Theorem 3.2 in [57]). Assume that $(X_1, Y_1), \ldots, (X_{N+1}, Y_{N+1})$ are exchangeable random variables. Then, for any user-specified coverage level $\alpha \in (0, 1)$, we have that the conformal set from eq. (1) satisfies

$$\mathbb{P}\left(Y_{N+1} \in \mathcal{C}_{\alpha}(X_{N+1})\right) \ge 1 - \alpha. \tag{2}$$

Recent advancements have extended the framework to construct *adaptive* prediction sets with conditional coverage guarantees, whose size varies according to the model's confidence at specific inputs, producing narrower intervals in regions of low uncertainty [58, 59]. This adaptive property is particularly relevant to our work, as it informs optimization algorithms about local prediction confidence when checking for feasibility. For instance, *Mondrian conformal prediction* [57, 60] enables conditional coverage guarantees for both ground-truth feasible points $\{x:h(x) \in \mathcal{Y}\}$ and infeasible points $\{x:h(x) \notin \mathcal{Y}\}$. This conditional coverage guarantee is formalized in the following Lemma, with the proof and details on Mondrian conformal sets deferred to Appendix B.1.

Lemma 3.1. (Ground-Truth Feasibility Conformal Coverage Guarantee): Assume that $(X_1, Y_1), \ldots, (X_{N+1}, Y_{N+1})$ are exchangeable random variables and that s(X, Y) is a symmetric score function. Mondrian conformal sets satisfy the conditional coverage guarantee at level α :

$$\mathbb{P}(Y_{N+1} \in \mathcal{C}_{\alpha}(X_{N+1}) \mid Y_{N+1} \in \mathcal{Y}) \ge 1 - \alpha,$$

$$\mathbb{P}(Y_{N+1} \in \mathcal{C}_{\alpha}(X_{N+1}) \mid Y_{N+1} \notin \mathcal{Y}) \ge 1 - \alpha.$$

Intuitively, this condition ensures that coverage holds simultaneously within both the feasible and infeasible regions, by constructing conformal sets calibrated separately for each group. This conditional guarantee is strictly stronger than (and implies) the marginal property from Theorem 3.1, preventing the procedure from achieving marginal coverage when high coverage on one subset compensates for poor coverage on the other.

4 Conformal Mixed-Integer Constraint Learning

In this section, we introduce our proposed *Conformal Mixed-Integer Constraint Learning* (C-MICL) framework, integrating conformal prediction into MICL problems to provide probabilistic guarantees on the ground-truth feasibility of feasible solutions. The core idea of C-MICL is to conformalize the learned constraint $y = \hat{h}(x) \in \mathcal{Y}$, thereby constructing a statistically valid uncertainty set that the predictions must satisfy. The general C-MICL problem can be written as follows:

$$\begin{array}{ll} \min\limits_{x,z} & f(x,z) \\ \text{s.t.} & g(x,z) \leq 0 \\ & (x,z) \in \mathcal{X} \\ & \mathcal{C}_{\alpha}(x) \subseteq \mathcal{Y} \end{array} \tag{C-MICL}$$

where $\mathcal{C}_{\alpha}(x)$ is a conformal set for the decision variable x based on the learned model $\hat{h}(x)$, a target coverage level α and a calibration set \mathcal{D}_{cal} . In the following, we detail how to construct $\mathcal{C}_{\alpha}(x)$ for both regression and classification settings, and how the constraint $\mathcal{C}_{\alpha}(x) \subseteq \mathcal{Y}$ can be formulated using MIP. Importantly, our approach remains compatible with the existing definition of ground-truth feasibility provided in Definition 3.1, $h(x) \in \mathcal{Y}$, requiring no modifications to the underlying feasibility notion.

Remark 4.1. The feasible regions defined by the oracle constraint $h(x) \in \mathcal{Y}$ and the conformal set containment $C_{\alpha}(x) \subseteq \mathcal{Y}$ do not exhibit a clear inclusion relationship in general. In particular, neither region necessarily contains the other, and therefore one is not more conservative nor yields larger/smaller optimal values in general.

We assume that the model h(x) is trained on $\mathcal{D}_{\text{train}}$ and considered fixed thereafter. The conformal sets are then constructed using a disjoint calibration dataset \mathcal{D}_{cal} of size N. This approach preserves the theoretical coverage guarantees from Theorem 3.1 and Lemma 3.1 without requiring model retraining [57]. To establish formal guarantees for the C-MICL approach, we introduce the following assumption about the conformal coverage of the feasible region of the C-MICL problem.

Assumption 4.1. (Conditional Independence of Feasibility and Coverage): The events of C-MICL feasibility and conformal coverage are conditionally independent given ground-truth feasibility.

This assumption states that the C-MICL constraints do not systematically bias the coverage properties of conformal prediction sets, conditional on the *true* (unknown) feasibility status of h(x). That is, we assume that whether a point is C-MICL feasible, i.e., $(x', z') \in \mathcal{F}_N = \{(x, z) \in \mathcal{X} : g(x, z) \leq 0, \mathcal{C}_{\alpha}(x) \subseteq \mathcal{Y}\}$ does not affect the probability that its conformal set contains the true function value, $h(x') \in \mathcal{C}_{\alpha}(x')$, given its ground-truth feasibility $(h(x') \in \mathcal{Y})$ or infeasibility $(h(x') \notin \mathcal{Y})$. We provide a detailed discussion and interpretation of this assumption in Appendix B.3.

Under this condition and ground-truth feasibility conformal coverage guarantee (Lemma 3.1), we establish our main theoretical result in the following theorem, which we prove in Appendix B.2. This result provides a probabilistic certificate for the ground-truth feasibility of feasible solutions to the C-MICL problem.

Theorem 4.1. Let $\mathcal{F}_N = \{(x,z) \in \mathcal{X} : g(x,z) \leq 0, \mathcal{C}_{\alpha}(x) \subseteq \mathcal{Y}\}$ be the feasible region of the *C-MICL* problem. Under the same conditions of Lemma 3.1 and Assumption 4.1, if \mathcal{F}_N is nonempty, then any feasible solution $(x',z') \in \mathcal{F}_N$ is ground-truth feasible with probability at least $1-\alpha$:

$$\mathbb{P}(h(x') \in \mathcal{Y} \mid (x', z') \in \mathcal{F}_N) \ge 1 - \alpha. \tag{3}$$

4.1 Regression Conformal MICL

For the regression setting, we adopt a two-step procedure to construct a predictive model $\hat{h}(x)$ and then quantify its uncertainty through an auxiliary uncertainty function $\hat{u}(x)$. Specifically, in the first step we use $\mathcal{D}_{\text{train}}$ to train a predictive regression model $\hat{h}(x)$. Then, we use the set of absolute residuals of the model in the training set $\mathcal{D}_{\text{train}}^{\text{res}} = \{(x_i, |\hat{h}(x_i) - y_i|)\}_{i=1}^{N_{\text{train}}}$ to train the secondary regression model $\hat{u}(x)$, which estimates the prediction uncertainty of $\hat{h}(x)$ as a function of the input variable x [61]. While alternative one-step approaches exist for simultaneously estimating both prediction and uncertainty models (often by leveraging model-specific choices of $\hat{h}(x)$), we adopt this two-step procedure because it is simple, broadly applicable, and model-agnostic. For a discussion of joint training techniques, we refer the reader to Angelopoulos et al. [57]. Note that our framework imposes no restrictions on the choice of predictive model for $\hat{u}(x)$; however, to ensure integration with the optimization formulation, we assume that $\hat{u}(x)$ admits a MIP representation.

Although $\hat{u}(x)$ provides a heuristic estimate of uncertainty, it can still suffer from misspecification due to data noise, bias, or limited training. To address this and provide formal guarantees on predictive reliability, we apply conformal prediction to calibrate the uncertainty estimates. Specifically, for each point in the calibration dataset \mathcal{D}_{cal} , a conformal score is computed as $s(x_i, y_i) = \frac{|\hat{h}(x_i) - y_i|}{\hat{u}(x_i)}$, following the procedure in Lei et al. [62]. Then, as outlined in Section 3.2, the conformal quantile $\hat{q}_{1-\alpha} = \text{Quantile}\left(s(x_1, y_1), \ldots, s(x_N, y_N); (1-\alpha)(1+1/N)\right)$ is computed based on the set of conformal scores. Using this quantile the conformal set constraint can be expressed as:

$$C_{\alpha}(x) \subseteq \mathcal{Y} \iff [\hat{h}(x) \pm \hat{q}_{1-\alpha} \cdot \hat{u}(x)] \subseteq [y, \overline{y}]$$
(4)

This set is defined by a pair of algebraic expressions that further constrain the lower and upper bounds on y and can be directly embedded into the optimization formulation (see Appendix C). In summary, C-MICL in the regression setting involves (i) training an uncertainty model $\hat{u}(x)$, (ii) computing the conformal quantile offline (so that the formulation size remains independent of N), and (iii) expressing the valid region for y through MIP constraints in terms of $\hat{h}(x)$ and $\hat{u}(x)$.

4.2 Classification Conformal MICL

For the classification setting, C-MICL does not require integrating an additional uncertainty model into the formulation. Instead, score functions can be directly computed from the predictor $\hat{h}(x) = \hat{\mathbb{P}}(Y \mid X = x)$, which typically involves a nonlinear transformation such as softmax [55]. However, in MICL, it is common practice to remove the final nonlinear transformation and operate directly on the pre-activation values of the last layer, commonly referred to as logits, to preserve linearity [8]. This approach does not affect the classification result: since softmax and sigmoid functions are monotonic, the class with the highest logit still corresponds to the highest probability.

To enable integration into MIP, we propose a valid conformal score function that operates directly in logit space. Specifically, we define the conformal score as:

$$s(x_i, y_i) = -\sum_{k \in \mathcal{K}} y_i^k \cdot \hat{h}(x_i)^k$$
(5)

where we select the negative of the logit output corresponding to the correct class, assuming that y_i is one-hot encoded. This formulation preserves the validity of the conformal method while enabling exact and efficient representation within a linear optimization model.

Given $\hat{q}_{1-\alpha} = \text{Quantile}\left(s(x_1,y_1),\ldots,s(x_N,y_N);(1-\alpha)(1+1/N)\right)$ as defined in Section 3.2, we define the conformal sets in the classification setting as:

$$C_{\alpha}(x) \subseteq \mathcal{Y} \iff -\hat{h}(x)^k > \hat{q}_{1-\alpha} \quad \forall \ k \in \mathcal{K}^{\text{und}}$$
 (6)

which ensures that only desired classes $k \in \mathcal{K}^{\mathrm{des}}$ are included in the conformal prediction set, by retaining those whose conformal scores lie below the quantile threshold. This behavior can be captured within a MIP framework by introducing an indicator constraint of the form $\mathbb{I}\{-\hat{h}(x)^k \leq \hat{q}_{1-\alpha}\}$ for all $k \in \mathcal{K}$, and enforcing it to be active only for desired classes $k \in \mathcal{K}^{\mathrm{des}}$. This logical constraint can be further reformulated as a set of algebraic mixed-integer inequalities using standard big-M techniques; we refer the reader to Appendix C for further details on this reformulation.

Remark 4.2. The functional forms chosen for the score functions s(x) and feasibility sets \mathcal{Y} serve primarily to illustrate the algorithmic implementation of C-MICL across broad classes of regression and classification problems within optimization formulations. However, any alternative score functions (see [57]) and feasibility sets can be used, provided they admit MIP representations [41].

5 Computational Experiments

We empirically validate our Conformal MICL (C-MICL) approach in both regression and multi-class classification settings. In each case, we benchmark against standard single-model MICL methods as well as the ensemble-based heuristic W-MICL proposed by Maragno et al. [7]. We denote the latter as W-MICL(P), where P refers to the ensemble size. For both settings, we evaluate performance on 100 randomly generated optimization instances, each defined by a sampled cost vector. This approach considers the calibration set $\mathcal{D}_{\rm cal}$ as fixed and allows us to sample from the feasible region for each C-MICL formulation. This allows us to empirically validate the theoretical guarantee in Theorem 4.1, suggesting ground-truth feasibility of optimal solutions that are feasible by construction. Moreover, in the Appendix E we empirically verify that the assumptions from Theorem 4.1 approximately holds in our experimental settings for both regression and classification problems.

In the following, we report: (i) the empirical ground-truth feasibility rate, (ii) the relative distance of optimal value $\Delta\%_i = \frac{f_i^* - f_{\rm CMICL}^*}{f_{\rm CMICL}^*} \cdot 100\%$ for each baseline method i when compared to ours, and (iii) the average CPU time required to solve each instance. The results shown in the main text correspond to a target coverage level of $\alpha=10\%$ across 100 runs with 95% confidence intervals (CIs) detailed in Appendix E. Additional results for $\alpha=5\%$ are available in Appendix E.1 and E.2. Key architectural and hyperparameter choices are summarized in the main text, while full details on hyperparameter selection, training, and optimization formulations are provided in Appendix D.

We focus our case studies on mixed-integer linear programming (MILP) examples, which are faster to solve in practice. This choice is made without loss of generality, as all proposed formulations and theoretical guarantees remain valid for both mixed-integer linear and nonlinear programs. All optimization problems were solved using the MILP solver Gurobi v12.0.1 with a relative optimality gap of 1% [63]. Machine learning models were implemented using scikit-learn and PyTorch, and subsequently integrated into Pyomo-based [64] optimization formulations via the open-source library OMLT [65].

5.1 Regression setting

We consider the optimal design and operation of a membrane reactor for methane aromatization and hydrogen production following Carrasco and Lima [66]. The goal is to determine five input variables

³Specifically, we implement the formulation proposed in Section 3.1 of Maragno et al. [7] using a bootstrap sample proportion of 0.5 as done by the authors.

(reactant flows, operating temperature, and reactor dimensions) to minimize operational cost while ensuring a product flow ≥ 50 , defining the feasible region $\mathcal Y$. The true system behavior is modeled by a set of ordinary differential equations serving as the oracle h(X). In this setting, $y \in [0, 100]$ defines an underlying regression problem. To generate training data, we sample 1,000 input points from $\mathcal X$ and evaluate y via the discretized ODEs with added Gaussian noise to simulate measurement uncertainty [67]. See Appendix D for additional details.

We train baseline MICL models using the full dataset, following established methodologies for different predictor types: ReLU NNs [9], GBTs [11], RFs [16], and LMDTs [17]. To evaluate the performance of the robust ensemble heuristic W-MICL, we follow the approach of Maragno et al. [7], training ensembles of size $P \in \{5, 10, 25, 50\}$ using the same cross-validated hyperparameters as their corresponding single models as described in Appendix D.1. Finally, we implement our proposed C-MICL method, which only requires training two models: $\hat{h}(x)$ and $\hat{u}(x)$. We use the same model architectures and hyperparameters as the baselines, allocating 80% of the data for training, and the remaining 20% for conformal calibration. All base models shared a common architecture for the uncertainty model: a ReLU NN with two hidden layers, each with 32 units.

Figure 1 presents the main results, comparing the ground truth feasibility coverage of the various methods across these instances for $\alpha=10\%$. The figure shows that, across all underlying models, our proposed C-MICL method consistently achieves the target ground-truth feasibility rate, as guaranteed by Theorem 4.1. In sharp contrast, none of the baseline MICL methods reach the desired coverage level, regardless of the machine learning model employed. Moreover, the W-MICL heuristic from Maragno et al. [7] only empirically approaches the target coverage when using neural networks. These results highlight the importance of theoretical guarantees: in settings where no reliable machine learning model is available to support MICL or W-MICL approaches, our Conformal MICL framework ensures ground-truth feasibility by construction. Furthermore, the level of probabilistic guarantee can be flexibly controlled by the parameter α , allowing practitioners to tune the method to any desired confidence level in (0,1) (with $\alpha=0.05$ included in Appendix E). In the same appendix, violin plots of the true constraint values h(x) for both α levels illustrate the magnitude of violations, showing that C-MICL not only incurs the fewest violations but also achieves the smallest violation magnitudes when they occur.

Figure 2 presents the relative distance in optimal objective value between each method and our proposed C-MICL approach, while Figure 3 reports the average CPU time (in seconds) required to solve the 100 optimization instances. In both figures, lighter-colored bars denote methods that failed to meet the desired feasibility rate shown in Figure 1. Results correspond to ReLU NNs, which exhibited the strongest empirical performance in terms of feasibility.

Figure 2 shows that our proposed C-MICL method achieves the desired feasibility guarantees with only a modest average reduction of 7% in objective value compared to heuristic approaches. This small trade-off highlights that C-MICL remains highly competitive while offering formal probabilistic guarantees on feasibility. Figure 3 further highlights the computational advantages of our approach, showing that it achieves optimal solutions several orders of magnitude faster than the baseline

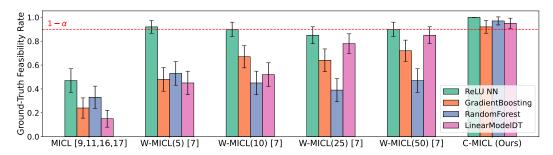
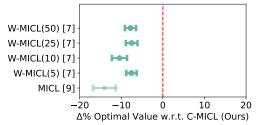
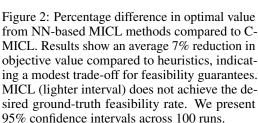


Figure 1: Empirical ground-truth feasibility rate of optimal solution across MICL methods on 100 optimization problem instances. Our Conformal MICL approach (rightmost bars) consistently achieves the target ground-truth feasibility rate $\geq 90\%$ regardless of the underlying base model used. Previous methods demonstrate variable performance below this theoretical guarantee. Error bars correspond to 95% confidence intervals for ground-truth feasibility rates across 100 runs.





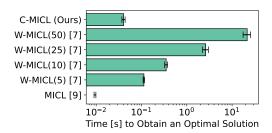


Figure 3: Average computational time to find optimal solution from NN-based MICL methods. Our C-MICL approach (top) is two orders of magnitude faster than previous methods. MICL (lighter bar) does not achieve the desired empirical feasibility rate. Error bars represent 95% confidence intervals across 100 runs.

methods. In summary, our method offers a compelling combination of computational speedup and theoretical probabilistic guarantees, while preserving comparable objective performance.

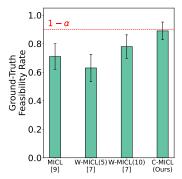
5.2 Multi-classification setting

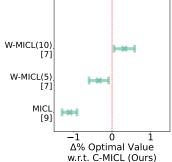
We consider a food basket design problem aimed at minimizing the cost of a basket composed of 25 commodities, subject to nutritional constraints covering 12 nutrients and a learned palatability requirement [68]. Palatability, reflecting the appeal or acceptability of the basket to recipients, is inferred from historical data on previously deployed baskets and community feedback. Following the formulation in Fajemisin et al. [1], we impose a minimum palatability threshold of 0.5 (consistent with the value proposed in Maragno et al. [7]) to ensure that the basket is not only nutritionally adequate and cost-effective, but also socially and culturally appropriate.

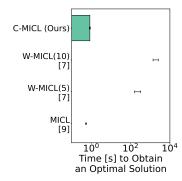
We use the dataset of 5,000 food baskets from Maragno et al. [7], where palatability scores range continuously in [0, 1]. To frame this as a classification task, we discretize the scores into four categories (*bad*, *regular*, *good*, and *very good*), using thresholds at 0.25, 0.5, and 0.75, respectively. Consistent with Maragno et al. [7], we constrain the optimization problem to only include baskets labeled *good* or *very good*. As no mechanistic oracle exists for palatability, we train a regression neural network on all 5,000 samples to approximate this function. The model's predictions are thresholded to produce the categorical ground truth used in our experiments. For a consistent comparison across MICL, W-MICL, and C-MICL, we benchmark all methods on a subset of 2,500 baskets data points.

We evaluate classification performance using ReLU NNs [9], the only model type readily supported by OMLT for classification tasks [65]. Despite this, our method remains model-agnostic, as demonstrated in the regression case study. For MICL and W-MICL (with ensemble sizes $P = \{5, 10\}$), we train on all 2,500 data points. For our proposed C-MICL approach, which requires only one trained model, we use 2,300 points for training and reserved 200 for conformal calibration. Full details on network architecture and training hyperparameters are provided in Appendix D.2.

As shown in Figure 4, only our C-MICL method achieves the desired empirical coverage at $\alpha=10\%$, in line with Theorem 4.1. This highlights the importance of formal statistical guarantees, especially under data constraints and high-dimensional inputs. Figure 5 shows the relative optimality gaps across methods, with differences averaging around 1%. However, none of the baselines yield implementable solutions, as only our method satisfies the desired empirical coverage guarantee. Figure 6 shows that C-MICL achieves solution times comparable to solving a single MICL model and is significantly faster than the W-MICL heuristics. In both figures, lighter bars indicate methods that do not meet the target empirical coverage. In summary, our method is the only one to guarantee ground-truth feasibility, preserves objective performance, and scales efficiently in classification tasks.







bars represent 95% CIs.

Figure 4: Empirical ground- Figure 5: Percentage difference truth feasibility rate of optimal in optimal value from NN-based solution across MICL methods. MICL methods compared to C-Our C-MICL approach (right- MICL. Difference is around 1% most bar) is the *only* method however all baseline approaches that consistently satisfies the tar- (lighter intervals) fail to achieve get ground-truth feasibility rate the desired feasibility guarantee. $\geq 90\%$ across 100 runs. Error Results across 100 runs, with 95% CIs.

Figure 6: Average computational time to find optimal solution from NN-based MICL methods, across 100 runs. C-MICL achieves solution times comparable to single-model MICL, and is significantly faster than the W-MICL heuristics. Error bars represent 95% CIs.

Discussion

Our work introduces C-MICL, a model-agnostic framework for incorporating learned constraints into optimization problems with probabilistic feasibility guarantees. C-MICL supports both regression and classification tasks, leverages conformal prediction to construct uncertainty sets, and encodes these into MIP-compatible formulations. It avoids the scalability bottlenecks of ensemble-based methods by requiring training at most two models, achieving substantial computational gains while maintaining comparable high-quality solutions. Furthermore, since our proposed method separates training from optimization, the size of the resulting optimization formulation is independent of the amount of data used to train the models, unlike previous end-to-end approaches (e.g., Zhao et al. [48]). Similarly, C-MICL can be integrated with existing methodologies aimed at improving the solution quality of standard MICL and W-MICL formulations, such as the enlarged convex hull trust-region approach proposed in Maragno et al. [7] or the trust-region-filter framework in Biegler [18]. C-MICL extends naturally to learning objective functions, providing probabilistic guarantees on true objective values, or to multiple learned constraints, representing a promising direction for future research on conformal-enhanced optimization.

One limitation of our work, common to general conformal prediction methods, is that when the underlying machine learning model performs poorly, the resulting conformal sets can become too wide and uninformative, potentially rendering the C-MICL problem infeasible. Moreover, while we believe the conditional independence assumption of feasibility and coverage is reasonable in our setting, it could potentially fail in scenarios where the feasible region is biased towards data regions where the conformal coverage does not hold. Nevertheless, an assumption of this nature is necessary to address the inherent dependency between the feasible solutions of the C-MICL problem and \mathcal{D}_{cal} , as this calibration set defines the feasible region, which would otherwise invalidate standard conformal guarantees. This represents a theoretical compromise when integrating conformal prediction with constrained optimization settings that enables us to provide statistical guarantees that were previously unattainable in MICL approaches.

This work contributes to the integration of machine learning and mathematical optimization by providing a principled framework for embedding data-driven constraints into decision-making processes with statistical guarantees. By ensuring marginal feasibility without access to the true constraint function, C-MICL advances the state of constraint learning, extends predict-then-optimize pipelines, and offers a tractable approach to optimizing over uncertain, learned constraints in real-world applications.

References

- [1] Adejuyigbe O Fajemisin, Donato Maragno, and Dick den Hertog. Optimization with constraint learning: A framework and survey. *European Journal of Operational Research*, 314(1):1–14, 2024.
- [2] Zhenan Fan, Bissan Ghaddar, Xinglu Wang, Linzi Xing, Yong Zhang, and Zirui Zhou. Artificial intelligence for operations research: Revolutionizing the operations research process. *arXiv* preprint arXiv:2401.03244, 2024.
- [3] Francisco Javier López-Flores, César Ramírez-Márquez, and José María Ponce-Ortega. Process Systems Engineering Tools for Optimization of Trained Machine Learning Models: Comparative and Perspective. *Industrial & Engineering Chemistry Research*, 63(32):13966–13979, 2024.
- [4] Atharv Bhosekar and Marianthi Ierapetritou. Advances in surrogate based modeling, feasibility analysis, and optimization: A review. *Computers & Chemical Engineering*, 108:250–267, 2018.
- [5] Sun Hye Kim and Fani Boukouvala. Machine learning-based surrogate modeling for data-driven optimization: a comparison of subset selection for regression techniques. *Optimization Letters*, 14(4):989–1010, 2020.
- [6] Huayu Tian and Marianthi G Ierapetritou. A surrogate-based framework for feasibility-driven optimization of expensive simulations. *AIChE Journal*, 70(5):e18364, 2024.
- [7] Donato Maragno, Holly Wiberg, Dimitris Bertsimas, Ş İlker Birbil, Dick den Hertog, and Adejuyigbe O Fajemisin. Mixed-integer optimization with constraint learning. *Operations Research*, 73(2):1011–1028, 2025.
- [8] Matteo Fischetti and Jason Jo. Deep neural networks and mixed integer linear optimization. *Constraints*, 23(3):296–309, 2018.
- [9] Bjarne Grimstad and Henrik Andersson. ReLU networks as surrogate models in mixed-integer linear programs. *Computers & Chemical Engineering*, 131:106580, 2019.
- [10] Ross Anderson, Joey Huchette, Will Ma, Christian Tjandraatmadja, and Juan Pablo Vielma. Strong mixed-integer programming formulations for trained neural networks. *Mathematical Programming*, 183(1):3–39, 2020.
- [11] Miten Mistry, Dimitrios Letsios, Gerhard Krennrich, Robert M Lee, and Ruth Misener. Mixed-integer convex nonlinear optimization with gradient-boosted trees embedded. *INFORMS Journal on Computing*, 33(3):1103–1119, 2021.
- [12] Alexander Thebelt, Calvin Tsay, Robert Lee, Nathan Sudermann-Merx, David Walz, Behrang Shafei, and Ruth Misener. Tree ensemble kernels for Bayesian optimization with known constraints over mixed-feature spaces. *Advances in Neural Information Processing Systems*, 35: 37401–37415, 2022.
- [13] Dimitris Bertsimas, Jack Dunn, and Aris Paschalidis. Regression and classification using optimal decision trees. In 2017 IEEE MIT undergraduate research technology conference (URTC), pages 1–4. IEEE, 2017.
- [14] Alessio Bonfietti, Michele Lombardi, and Michela Milano. Embedding decision trees and random forests in constraint programming. In *Integration of AI and OR Techniques in Constraint Programming: 12th International Conference, CPAIOR 2015, Barcelona, Spain, May 18-22, 2015, Proceedings 12*, pages 74–90. Springer, 2015.
- [15] Max Biggs, Rim Hariss, and Georgia Perakis. Optimizing objective functions determined from random forests. Available at SSRN 2986630, 2017.
- [16] Velibor V Mišić. Optimization of tree ensembles. Operations Research, 68(5):1605–1624, 2020.

- [17] Bashar L Ammari, Emma S Johnson, Georgia Stinchfield, Taehun Kim, Michael Bynum, William E Hart, Joshua Pulsipher, and Carl D Laird. Linear model decision trees as surrogates in optimization of engineering applications. *Computers & Chemical Engineering*, 178:108347, 2023.
- [18] Lorenz T Biegler. The trust region filter strategy: Survey of a rigorous approach for optimization with surrogate models. *Digital Chemical Engineering*, 13:100197, 2024.
- [19] Keliang Wang, Leonardo Lozano, Carlos Cardonha, and David Bergman. Optimizing over an ensemble of trained neural networks. *INFORMS Journal on Computing*, 35(3):652–674, 2023.
- [20] Rudy R Bunel, Ilker Turkaslan, Philip Torr, Pushmeet Kohli, and Pawan K Mudigonda. A unified view of piecewise linear neural network verification. Advances in neural information processing systems, 31, 2018.
- [21] Vincent Tjeng, Kai Xiao, and Russ Tedrake. Evaluating robustness of neural networks with mixed integer programming. *arXiv preprint arXiv:1711.07356*, 2017.
- [22] Chih-Hong Cheng, Georg Nührenberg, and Harald Ruess. Maximum resilience of artificial neural networks. In *Automated Technology for Verification and Analysis: 15th International Symposium, ATVA 2017, Pune, India, October 3–6, 2017, Proceedings 15*, pages 251–268. Springer, 2017.
- [23] Alessio Lomuscio and Lalit Maganti. An approach to reachability analysis for feed-forward ReLU neural networks. *arXiv preprint arXiv:1706.07351*, 2017.
- [24] Elena Botoeva, Panagiotis Kouvaros, Jan Kronqvist, Alessio Lomuscio, and Ruth Misener. Efficient verification of ReLU-based neural networks via dependency analysis. In *Proceedings* of the AAAI Conference on Artificial Intelligence, volume 34, pages 3291–3299, 2020.
- [25] Kaidi Xu, Huan Zhang, Shiqi Wang, Yihan Wang, Suman Jana, Xue Lin, and Cho-Jui Hsieh. Fast and complete: Enabling complete neural network verification with rapid and massively parallel incomplete verifiers. *arXiv* preprint arXiv:2011.13824, 2020.
- [26] Panagiotis Kouvaros and Alessio Lomuscio. Towards Scalable Complete Verification of ReLU Neural Networks via Dependency-based Branching. 30th International Joint Conference on Artificial Intelligence, pages 2643–2650, 2021.
- [27] Shiqi Wang, Huan Zhang, Kaidi Xu, Xue Lin, Suman Jana, Cho-Jui Hsieh, and J Zico Kolter. Beta-crown: Efficient bound propagation with per-neuron split constraints for neural network robustness verification. *Advances in neural information processing systems*, 34:29909–29921, 2021.
- [28] Pin-Yu Chen, Yash Sharma, Huan Zhang, Jinfeng Yi, and Cho-Jui Hsieh. Ead: elastic-net attacks to deep neural networks via adversarial examples. *Proceedings of the AAAI conference on artificial intelligence*, 32(1), 2018.
- [29] Francesco Croce and Matthias Hein. Minimally distorted adversarial examples with a fast adaptive boundary attack. In *International conference on machine learning*, pages 2196–2205. PMLR, 2020.
- [30] Calvin Tsay, Jan Kronqvist, Alexander Thebelt, and Ruth Misener. Partition-based formulations for mixed-integer optimization of trained ReLU neural networks. *Advances in neural information processing systems*, 34:3068–3080, 2021.
- [31] Buser Say, Ga Wu, Yu Qing Zhou, and Scott Sanner. Nonlinear Hybrid Planning with deep net Itearned transition models and mixed-integer linear programming. In *International Joint Conference on Artificial Intelligence 2017*, pages 750–756. Association for the Advancement of Artificial Intelligence (AAAI), 2017.
- [32] Moonkyung Ryu, Yinlam Chow, Ross Anderson, Christian Tjandraatmadja, and Craig Boutilier. CAQL: Continuous action Q-learning. *arXiv preprint arXiv:1909.12397*, 2019.

- [33] Arthur Delarue, Ross Anderson, and Christian Tjandraatmadja. Reinforcement learning with combinatorial actions: An application to vehicle routing. Advances in Neural Information Processing Systems, 33:609–620, 2020.
- [34] Jochen L Cremer, Ioannis Konstantelos, Simon H Tindemans, and Goran Strbac. Data-driven power system operation: Exploring the balance between cost and risk. *IEEE Transactions on Power Systems*, 34(1):791–801, 2018.
- [35] Spyros Chatzivasileiadis. From decision trees and neural networks to MILP: Power system optimization considering dynamic stability constraints. In 2020 European control conference (ECC), page 594. IEEE, 2020.
- [36] Dimitris Bertsimas, Allison O'Hair, Stephen Relyea, and John Silberholz. An analytics approach to designing combination chemotherapy regimens for cancer. *Management Science*, 62(5): 1511–1531, 2016.
- [37] Santiago Cortes-Gomez, Mateo Dulce, Carlos Patino, and Bryan Wilder. Statistical inference under constrained selection bias. In *Proceedings of the 41st International Conference on Machine Learning*, pages 9361–9379, 2024.
- [38] Oluwadare Badejo and Marianthi Ierapetritou. Integrating tactical planning, operational planning and scheduling using data-driven feasibility analysis. *Computers & Chemical Engineering*, 161: 107759, 2022.
- [39] Daniel Ovalle, Javal Vyas, Carl D Laird, and Ignacio E Grossmann. Integration of plant scheduling feasibility with supply chain network under disruptions using machine learning surrogates. In *Computer aided chemical engineering*, volume 53, pages 1489–1494. Elsevier, 2024.
- [40] Georgia Stinchfield, Natali Khalife, Bashar L Ammari, Joshua C Morgan, Miguel Zamarripa, and Carl D Laird. Mixed-integer linear programming formulation with embedded machine learning surrogates for the design of chemical process families. *Industrial & Engineering Chemistry Research*, 64(16):8299–8311, 2025.
- [41] Ruth Misener and Lorenz Biegler. Formulating data-driven surrogate models for process optimization. *Computers & Chemical Engineering*, 179:108411, 2023.
- [42] John P Eason and Lorenz T Biegler. Reduced model trust region methods for embedding complex simulations in optimization problems. In *Computer Aided Chemical Engineering*, volume 37, pages 773–778. Elsevier, 2015.
- [43] John P Eason and Lorenz T Biegler. Advanced trust region optimization strategies for glass box/black box models. AIChE Journal, 64(11):3934–3943, 2018.
- [44] Chirag Gupta, Arun K Kuchibhotla, and Aaditya Ramdas. Nested conformal prediction and quantile out-of-bag ensemble methods. *Pattern Recognition*, 127:108496, 2022.
- [45] Rina Foygel Barber, Emmanuel J. Candès, Aaditya Ramdas, and Ryan J. Tibshirani. Conformal Prediction Beyond Exchangeability. *The Annals of Statistics*, 51(2):816–845, 2023.
- [46] Ran Xie, Rina Barber, and Emmanuel Candes. Boosted conformal prediction intervals. *Advances in Neural Information Processing Systems*, 37:71868–71899, 2024.
- [47] Bo Lin, Erick Delage, and Timothy Chan. Conformal inverse optimization. *Advances in Neural Information Processing Systems*, 37:63534–63564, 2024.
- [48] Yiqi Zhao, Xinyi Yu, Jyotirmoy V Deshmukh, and Lars Lindemann. Conformal predictive programming for chance constrained optimization. *arXiv preprint arXiv:2402.07407*, 2024.
- [49] Chancellor Johnstone and Bruce Cox. Conformal uncertainty sets for robust optimization. *Conformal and Probabilistic Prediction and Applications*, pages 72–90, 2021.
- [50] Christopher Yeh, Nicolas Christianson, Alan Wu, Adam Wierman, and Yisong Yue. End-to-end conformal calibration for optimization under uncertainty. arXiv preprint arXiv:2409.20534, 2024.

- [51] Shayan Kiyani, George Pappas, Aaron Roth, and Hamed Hassani. Decision theoretic foundations for conformal prediction: Optimal uncertainty quantification for risk-averse agents. *arXiv* preprint arXiv:2502.02561, 2025.
- [52] Yash P Patel, Sahana Rayan, and Ambuj Tewari. Conformal contextual robust optimization. In *International Conference on Artificial Intelligence and Statistics*, pages 2485–2493. PMLR, 2024.
- [53] Juan Pablo Vielma. Mixed integer linear programming formulation techniques. *Siam Review*, 57(1):3–57, 2015.
- [54] Mustafa R Kılınç and Nikolaos V Sahinidis. Exploiting integrality in the global optimization of mixed-integer nonlinear programming problems with BARON. *Optimization Methods and Software*, 33(3):540–562, 2018.
- [55] Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. Conformal prediction: Classification and general case. In *Algorithmic Learning in a Random World*, pages 71–106. Springer, 2022.
- [56] Ryan J Tibshirani, Rina Foygel Barber, Emmanuel Candes, and Aaditya Ramdas. Conformal prediction under covariate shift. Advances in neural information processing systems, 32, 2019.
- [57] Anastasios N Angelopoulos, Rina Foygel Barber, and Stephen Bates. Theoretical foundations of conformal prediction. *arXiv preprint arXiv:2411.11824*, 2024.
- [58] Yaniv Romano, Matteo Sesia, and Emmanuel Candes. Classification with valid and adaptive coverage. *Advances in neural information processing systems*, 33:3581–3591, 2020.
- [59] Matteo Sesia and Yaniv Romano. Conformal prediction using conditional histograms. *Advances in Neural Information Processing Systems*, 34:6304–6315, 2021.
- [60] Vladimir Vovk, David Lindsay, Ilia Nouretdinov, and Alex Gammerman. Mondrian confidence machine. *Technical Report*, 2003.
- [61] Harris Papadopoulos, Kostas Proedrou, Volodya Vovk, and Alex Gammerman. Inductive confidence machines for regression. In *Machine learning: ECML 2002: 13th European* conference on machine learning Helsinki, Finland, August 19–23, 2002 proceedings 13, pages 345–356. Springer, 2002.
- [62] Jing Lei, Max G'Sell, Alessandro Rinaldo, Ryan J Tibshirani, and Larry Wasserman. Distribution-free predictive inference for regression. *Journal of the American Statistical Association*, 113(523):1094–1111, 2018.
- [63] Gurobi Optimization, LLC. Gurobi Optimizer Reference Manual, 2024. URL https://www.gurobi.com.
- [64] Michael L Bynum, Gabriel A Hackebeil, William E Hart, Carl D Laird, Bethany L Nicholson, John D Siirola, Jean-Paul Watson, David L Woodruff, et al. *Pyomo-optimization modeling in python*, volume 67. Springer, 2021.
- [65] Francesco Ceccon, Jordan Jalving, Joshua Haddad, Alexander Thebelt, Calvin Tsay, Carl D Laird, and Ruth Misener. OMLT: Optimization & machine learning toolkit. *Journal of Machine Learning Research*, 23(349):1–8, 2022.
- [66] Juan C Carrasco and Fernando V Lima. Nonlinear operability of a membrane reactor for direct methane aromatization. *IFAC-PapersOnLine*, 48(8):728–733, 2015.
- [67] Juan C Carrasco and Fernando V Lima. Novel operability-based approach for process design and intensification: Application to a membrane reactor for direct methane aromatization. *AIChE Journal*, 63(3):975–983, 2017.
- [68] Koen Peters, Sérgio Silva, Rui Gonçalves, Mirjana Kavelj, Hein Fleuren, Dick Den Hertog, Ozlem Ergun, and Mallory Freeman. The nutritious supply chain: optimizing humanitarian food assistance. *INFORMS Journal on Optimization*, 3(2):200–226, 2021.

- [69] Joey Huchette, Gonzalo Muñoz, Thiago Serra, and Calvin Tsay. When deep learning meets polyhedral theory: A survey. *arXiv preprint arXiv:2305.00241*, 2023.
- [70] Francisco Trespalacios and Ignacio E Grossmann. Improved Big-M reformulation for generalized disjunctive programs. *Computers & Chemical Engineering*, 76:98–103, 2015.
- [71] Victor Alves, San Dinh, John R Kitchin, Vitor Gazzaneo, Juan C Carrasco, and Fernando V Lima. Opyrability: A Python package for process operability analysis. *Journal of Open Source Software*, 9(94):5966, 2024.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The abstract and introduction clearly state the paper's claims, with the contributions explicitly listed as bullet points in the introduction. These claims are well-supported by both theoretical developments and experimental results, and their scope is appropriately framed to reflect the work's actual contributions.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals
 are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We explicitly discuss the limitations of our work in the Discussion section and throughout the main body of the article.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: All theorems and formulas are numbered and cross-referenced. The main theoretical result and its assumptions are clearly stated in the core text, while the complete proof and justification of the assumptions are provided in the appendix.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The paper clearly describes all steps needed to reproduce the main experimental results, with the proposed optimization formulation fully detailed in the main text. To further support reproducibility, we include an appendix that illustrates how to implement the formulation constraint by constraint, although the main description is already sufficient. All datasets used are either publicly available and properly referenced or included with the supplementary code, which is provided to facilitate replication.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived
 well by the reviewers: Making the paper reproducible is important, regardless of
 whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
- (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).

(d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We provide code that demonstrates how to train the models, embed them into optimization formulations, and solve the resulting problems. Running the code requires a Gurobi license, which is freely available for academic use. The dataset for the basket case study is publicly available and properly referenced, while the dataset for the reactor case study is included in the supplementary material. Clear instructions are provided to ensure reproducibility of the main experimental results.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how
 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: The main text specifies the data splits, key hyperparameters, and optimization solver configurations necessary to understand the results. Additional implementation details, including cross-validation setup, are provided in an appendix to ensure completeness without detracting from the clarity of the core paper.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: All figures report 95% confidence intervals, with variability arising from solving 100 optimization problems. The method for computing the intervals is explained, and the closed-form formulas used are provided in the appendix to ensure clarity and reproducibility.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
 of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: A key aspect of the results is the CPU time required to obtain an optimal solution, which is clearly expressed in both the results section and the Complementary Results appendix. The main paper specifies the optimization solver version and key configuration settings. Additionally, the Implementation Details appendix provides further information about the computational environment, including the number and type of CPUs, available memory, and thread settings used during the experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: NeurIPS Code of Ethics was reviewed and followed while conducting research for this paper.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: Our contribution aims to support the development of robust and implementable machine learning solutions for real-world decision-making, where "implementable" means feasible with respect to the ground truth systems. This can positively impact domains such as infrastructure planning, resource allocation, and operational efficiency. While our work is primarily methodological, we recognize that improved optimization-based ML pipelines could influence high-stakes applications.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper poses no such risks. No models are being released, and all datasets used in the experiments are publicly available and do not involve scraped or sensitive content. Therefore, safeguards for responsible release are not applicable in this case.

Guidelines:

• The answer NA means that the paper poses no such risks.

- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The only existing asset used in this work is the publicly available dataset for the food basket example, which is properly cited in the paper. The dataset is hosted on Github and is released under the MIT license. This license explicitly permits commercial use, modification, distribution, and private use. Given these terms, appropriate attribution is provided in the paper, and the use of this asset complies with the license requirements (i.e., CC-BY 4.0).

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- · For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: The code and dataset released with this paper are intended to support full reproducibility of the experiments. These assets are openly licensed under the CC-BY 4.0 license, allowing for broad reuse with proper attribution. Required documentation is provided alongside the released assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects. Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: LLMs were only used for formatting purposes.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.

A Machine Learning Models as MIPs

In the following we provide an overview on how to formulate off-the-shelf *pre-trained* machine learning models as mixed-integer programs (MIP) into constraint learning formulations.

A.1 ReLU Neural Networks

We consider feed-forward neural networks composed of fully connected layers with ReLU activations, which can be exactly represented using MIP [9, 10]. Let x denote the input variables to the network and y its output variables. All other internal quantities (preactivations, activations, and indicator variables) are distinct auxiliary decision variables introduced solely for encoding the network structure. Therefore, each neuron computes a transformation of the form:

$$a = \max\left(0, w^{\top} s + b\right),\tag{7}$$

where s are the activations from the previous layer, and (w,b) are fixed, trained parameters. The ReLU nonlinearity is encoded via the so called big-M formulation:

$$a \ge w^{\top} s + b$$

$$a \le w^{\top} s + b - (1 - \delta)L$$

$$a \le \delta U$$

$$\delta \in \{0, 1\}$$
(8)

where a is the activation output, δ is a binary indicator for whether the neuron is active, and L, U are tight lower and upper bounds on the pre-activation $w^{\top}s + b$. The full network is encoded layer by layer starting from the inputs x, recursively defining internal pre-activations and activations, until reaching the output layer y.

For further details and alternatives on MIP reformulations of ReLU neural networks, we refer readers to a recent review by Huchette et al. [69]. While multiple equivalent encodings exist, they yield the same global solution [53], differing only in computational efficiency, and thus do not affect our discussion on ground-truth feasibility.

A.2 Tree Ensembles and Linear-Model Decision Trees

We encode trained tree ensembles (including random forests and gradient-boosted trees) as well as linear-model decision trees using a unified MIP formulation based on leaf selection [11, 16, 17]. Let \mathcal{T} denote the set of decision trees in the ensemble, and \mathcal{L}_t the set of leaf nodes in tree $t \in \mathcal{T}$. Each tree partitions the input space into disjoint regions (leaves), each associated with a constant prediction value. We introduce binary variables $r_{t,\ell} \in \{0,1\}$ indicating whether leaf $\ell \in \mathcal{L}_t$ of tree t is selected. A valid configuration activates exactly one leaf per tree:

$$\sum_{\ell \in \mathcal{L}_t} r_{t,\ell} = 1 \quad \forall \ t \in \mathcal{T}$$

$$\tag{9}$$

Each internal node in a tree performs a threshold split of the form $x_i < v_{i,j}$, where $v_{i,j}$ is a threshold value for input feature x_i . We introduce binary variables $w_{i,j} \in \{0,1\}$ to model these comparisons, where $w_{i,j} = 1$ if, and only if, the condition is satisfied. Monotonicity constraints $w_{i,j} \le w_{i,j+1}$ ensure consistency across thresholds, reflecting the ordered structure of decision splits.

Each leaf ℓ in tree t is reachable only if all the splitting conditions on the path to that leaf are satisfied. For each split node s in tree t, let i(s) be the split variable and j(s) the corresponding threshold index. Let Left t,s and Right t,s denote the sets of leaf nodes in the left and right subtrees of s. We enforce:

$$\sum_{\ell \in \text{Left}_{t,s}} r_{t,\ell} \le w_{i(s),j(s)}, \quad \sum_{\ell \in \text{Right}_{t,s}} r_{t,\ell} \le 1 - w_{i(s),j(s)} \quad \forall \ t \in \mathcal{T}, s \in \text{Splits}(t)$$
 (10)

The final prediction is computed as a weighted sum of the selected leaves: the weights are stage-dependent for gradient-boosted trees [11], uniform across trees for random forests [16], or equal to a linear function of the input features within the selected leaf for linear-model decision trees [17].

B Proof of Main Results

B.1 Mondrian Conformal Prediction

Given calibration data $\{(X_1,Y_1),\ldots,(X_N,Y_N)\}$, Mondrian conformal prediction is a general framework that extends full conformal prediction by enforcing *conditional coverage guarantees* across a pre-defined set of *finite* groups defined in terms of both the features variables X and the outcome variable Y [57, 60]. Formally, for a group-indicator function $g(X,Y) \in \{1,\ldots,K\}$, Mondrian conformal prediction construct conformal sets $\mathcal{C}_{\alpha}(X_{N+1})$ such that

$$\mathbb{P}(Y_{n+1} \in \mathcal{C}_{\alpha}(X_{N+1}) \mid g(X_{N+1}, Y_{N+1}) = k) \ge 1 - \alpha,$$

for all groups $k \in [K] = \{1, \dots K\}$.

The core insight of Mondrian conformal prediction lies in computing different conformal quantiles for each group $k \in [K]$. Then, for a test point X_{n+1} , we generate *hypothetical* points (X_{N+1}, y) for each possible value of the outcome variable y, and compare the conformal score $s(X_{N+1}, y)$ to the conformal scores $\{s(X_i, Y_i)\}_{i \in \mathcal{I}_{g(X_{N+1}, y)}}$ for all calibration data points in the same group as (X_{N+1}, y) . That is, the group-specific comparison is performed for $i \in \mathcal{I}_{g(X_{N+1}, y)} = \{(i \in [N] : g(X_i, Y_i) = g(X_{N+1}, y)\}.$

This approach yields the following Mondrian conformal sets

$$C_{\alpha}(X_{N+1}) = \{ y : s(X_{N+1}, y) \le \hat{q}_{1-\alpha}^y \}, \tag{11}$$

where $\hat{q}_{1-\alpha}^{y}$ is the group-specific quantile

$$\hat{q}_{1-\alpha}^y = \operatorname{Quantile}\left(\{s(X_i,Y_i)\}_{i\in\mathcal{I}_{g(X_{N+1},y)}}; (1-\alpha)(1+1/|\mathcal{I}_{g(X_{N+1},y)}|)\right).$$

The conditional coverage guarantee is formalized in the following theorem:

Theorem B.1 (Theorem 4.11 in [57]). Assume that $(X_1, Y_1), \ldots, (X_{N+1}, Y_{N+1})$ are exchangeable random variables and s is a symmetric score function. Then, the Mondrian conformal set $C_{\alpha}(X_{N+1})$ from eq. (11) satisfies the conditional coverage guarantee

$$\mathbb{P}(Y_{n+1} \in \mathcal{C}_{\alpha}(X_{N+1}) \mid g(X_{N+1}, Y_{N+1}) = k) \ge 1 - \alpha \tag{12}$$

for all groups $k \in \{1, ... K\}$ with $\mathbb{P}(g(X_{N+1}, Y_{N+1}) = k) > 0$, for any user-specified coverage level $\alpha \in (0, 1)$.

Thus, Lemma 3.1 directly applies Theorem B.1 using the group identifier function $g(x,y) = \mathbb{I}(y \in \mathcal{Y})$, which separates the data according to the ground-truth feasibility condition. Note that the condition $\mathbb{P}(g(X_{N+1},Y_{N+1})=k)>0$ from Theorem B.1 assumes that both feasible and infeasible points can be sampled from the underlying distribution \mathcal{P}_{XY} , which is given by construction of the constraint learning optimization problem.

B.2 Proof of Theorem 4.1

Proof. We first show that Lemma 3.1 and Assumption 4.1 imply that we have appropriate coverage for the feasible points of the C-MICL problem $\mathcal{F}_N = \{(x,z) \in \mathcal{X} : g(x,z) \leq 0, \mathcal{C}_{\alpha}(x) \subseteq \mathcal{Y}\}.$

By the law of total probability,

$$\mathbb{P}(h(x) \in \mathcal{C}_{\alpha}(x) \mid (x, z) \in \mathcal{F}_{N}) = \mathbb{P}(h(x) \in \mathcal{C}_{\alpha}(x) \mid (x, z) \in \mathcal{F}_{N}, h(x) \in \mathcal{Y}) \mathbb{P}(h(x) \in \mathcal{Y} \mid (x, z) \in \mathcal{F}_{N}) \\
+ \mathbb{P}(h(x) \in \mathcal{C}_{\alpha}(x) \mid (x, z) \in \mathcal{F}_{N}, h(x) \notin \mathcal{Y}) \mathbb{P}(h(x) \notin \mathcal{Y} \mid (x, z) \in \mathcal{F}_{N}) \\
= \mathbb{P}(h(x) \in \mathcal{C}_{\alpha}(x) \mid h(x) \in \mathcal{Y}) \mathbb{P}(h(x) \in \mathcal{Y} \mid (x, z) \in \mathcal{F}_{N}) \\
+ \mathbb{P}(h(x) \in \mathcal{C}_{\alpha}(x) \mid h(x) \notin \mathcal{Y}) \mathbb{P}(h(x) \notin \mathcal{Y} \mid (x, z) \in \mathcal{F}_{N}) \\
\geq (1 - \alpha) \mathbb{P}(h(x) \in \mathcal{Y} \mid (x, z) \in \mathcal{F}_{N}) + (1 - \alpha) \mathbb{P}(h(x) \notin \mathcal{Y} \mid (x, z) \in \mathcal{F}_{N}) \\
= (1 - \alpha) [\mathbb{P}(h(x) \in \mathcal{Y} \mid (x, z) \in \mathcal{F}_{N}) + \mathbb{P}(h(x) \notin \mathcal{Y} \mid (x, z) \in \mathcal{F}_{N})] \\
= (1 - \alpha).$$

where the second equality follows from the conditional independence assumption 4.1

$$\mathbb{P}(h(x) \in \mathcal{C}_{\alpha}(x) \mid (x, z) \in \mathcal{F}_{N}, h(x) \in \mathcal{Y}) = \mathbb{P}(h(x) \in \mathcal{C}_{\alpha}(x) \mid h(x) \in \mathcal{Y}),$$

$$\mathbb{P}(h(x) \in \mathcal{C}_{\alpha}(x) \mid (x, z) \in \mathcal{F}_{N}, h(x) \notin \mathcal{Y}) = \mathbb{P}(h(x) \in \mathcal{C}_{\alpha}(x) \mid h(x) \notin \mathcal{Y}).$$

Moreover, the first inequality follows from the ground-truth feasibility conformal coverage guarantee in Lemma 3.1 at level α

$$\min \left(\mathbb{P}(h(x) \in \mathcal{C}_{\alpha}(x) \mid h(x) \in \mathcal{Y}), \mathbb{P}(h(x) \in \mathcal{C}_{\alpha}(x) \mid h(x) \notin \mathcal{Y}) \right) \ge 1 - \alpha,$$

where the probabilities are taken with respect to both the calibration data \mathcal{D}_{cal} and the test point.

The previous result says that when a solution $(x',z') \in \mathcal{F}_N$ is feasible, its true function value h(x') falls within $\mathcal{C}_{\alpha}(x')$ with probability at least $1-\alpha$. Since $\mathcal{C}_{\alpha}(x')$ is a subset of \mathcal{Y} for feasible solutions to the C-MICL problem, whenever h(x') is in $\mathcal{C}_{\alpha}(x')$, it must also be in \mathcal{Y} . Therefore,

$$\mathbb{P}(h(x') \in \mathcal{Y} \mid (x', z') \in \mathcal{F}_N) \ge \mathbb{P}(h(x') \in \mathcal{C}_{\alpha}(x') \mid (x', z') \in \mathcal{F}_N) \ge 1 - \alpha,$$

which proves the ground-truth feasibility guarantee.

B.3 Discussion on the Conditional Independence of Feasibility and Coverage Assumption

Assumption 4.1 states that, conditional on ground-truth feasibility $(h(x) \in \mathcal{Y} \text{ or } h(x) \notin \mathcal{Y})$, the event of C-MICL feasibility $((x,z) \in \mathcal{F}_N = \{(x,z) \in \mathcal{X} : g(x,z) \leq 0, \mathcal{C}_\alpha(x) \subseteq \mathcal{Y}\})$ is independent of whether the conformal set contains the true function value $(h(x) \in \mathcal{C}_\alpha(x))$. To motivate Assumption 4.1 and clarify when it is plausible, we first emphasize that the ground-truth feasibility (GTF) conditional coverage guarantee from Lemma 3.1 can be achieved in a fully data-driven way (e.g., using Mondrian conformal prediction or other label-conditional conformal methods):

$$\mathbb{P}(h(x) \in \mathcal{C}_{\alpha}(x) \mid h(x) \in \mathcal{Y}) \ge 1 - \alpha$$

$$\mathbb{P}(h(x) \in \mathcal{C}_{\alpha}(x) \mid h(x) \notin \mathcal{Y}) \ge 1 - \alpha$$

However, in C-MICL we aim to guarantee coverage over the feasible region of the optimization problem $\mathcal{F}_N = \{(x, z) \in \mathcal{X} : g(x, z) \leq 0, \mathcal{C}_{\alpha}(x) \subseteq \mathcal{Y}\}$, i.e.,

$$\mathbb{P}\left(h(x) \in \mathcal{C}_{\alpha}(x) \mid (x, z) \in \mathcal{F}_{N}\right) \geq 1 - \alpha$$

If the predictive model $\widehat{h}(x)$ were perfect (i.e., $\widehat{h}(x) = h(x)$), then the regions \mathcal{F}_N and the ground-truth feasible region \mathcal{F} would coincide, and the coverage guarantee would transfer directly. However, since we are interested in the more realistic case where $\widehat{h}(x)$ is imperfect, \mathcal{F}_N and \mathcal{F} differ in a data-dependent way. In this case, since the feasible region \mathcal{F}_N is implicitly shaped by the calibration set (via $\mathcal{C}_{\alpha}(x)$), there is a natural dependency between the feasible solutions of the C-MICL problem and the calibration data, which invalidates standard conformal guarantees relying on exchangeable calibration and test data. Assumption 4.1 precisely seeks to decouple this dependency: it allows us to approximate conformal coverage within \mathcal{F}_N by assuming that feasibility does not systematically bias conformal validity, once conditioned on ground-truth feasibility.

To build intuition, consider partitioning \mathcal{F}_N into two disjoint subsets: $\mathcal{F}_N \cap \mathcal{F}$ and $\mathcal{F}_N \cap \mathcal{F}^c$. Then, Assumption 4.1 implies that $\mathcal{F}_N \cap \mathcal{F}$ (respectively $\mathcal{F}_N \cap \mathcal{F}^c$) is not systematically biased towards a region of $\mathcal{F}(\mathcal{F}^c)$ that is miscalibrated. Mathematically,

$$\mathbb{P}(h(x) \in \mathcal{C}_{\alpha}(x) \mid (x, z) \in \mathcal{F}_{N} \cap \mathcal{F}) = \mathbb{P}(h(x) \in \mathcal{C}_{\alpha}(x) \mid (x, z) \in \mathcal{F}_{N}, h(x) \in \mathcal{Y})$$

$$\approx \mathbb{P}(h(x) \in \mathcal{C}_{\alpha}(x) \mid h(x) \in \mathcal{Y})$$

$$\mathbb{P}(h(x) \in \mathcal{C}_{\alpha}(x) \mid (x, z) \in \mathcal{F}_{N} \cap \mathcal{F}^{c}) = \mathbb{P}(h(x) \in \mathcal{C}_{\alpha}(x) \mid (x, z) \in \mathcal{F}_{N}, h(x) \notin \mathcal{Y})$$

$$\approx \mathbb{P}(h(x) \in \mathcal{C}_{\alpha}(x) \mid h(x) \notin \mathcal{Y})$$

These enable us to translate the conformal coverage guarantees from the ground-truth feasible region \mathcal{F} to the feasible set \mathcal{F}_N used in the optimization. Assumption 4.1 is therefore reasonable when the calibration data adequately covers the parts of the input space that intersect the feasible region \mathcal{F}_N , both within the ground-truth feasible region \mathcal{F} and its complement \mathcal{F}^c . In this sense, it aligns with standard generalization assumptions that require the training and calibration data to be representative of the regions where predictions are deployed. In our experimental settings, we observe good empirical alignment between target and achieved coverage (Appendix E), suggesting that Assumption 4.1 holds reasonably well in practice in realistic data scenarios.

Alternatively, Assumption 4.1 can be approximated using more granular conditional conformal methods, by partitioning the optimization region into finer subregions and enforcing local coverage guarantees within each, which can then be translated to the feasible set \mathcal{F}_N . However, the assumption may break down if the feasible region \mathcal{F}_N is heavily concentrated in areas where the calibration set is sparse or systematically miscalibrated.

C Conformal Set MIP Reformulations

The reformulation of the conformal set in the regression setting is straightforward, as the conformal prediction interval is already algebraic in nature and can be directly expressed as:

$$C_{\alpha}(x) \subseteq \mathcal{Y} \iff \left[\hat{h}(x) \pm \hat{q}_{1-\alpha} \cdot \hat{u}(x)\right] \subseteq [\underline{y}, \overline{y}]$$
 (13)

$$\iff \begin{cases} \hat{h}(x) + \hat{q}_{1-\alpha} \cdot \hat{u}(x) \leq \overline{y}, \\ \hat{h}(x) - \hat{q}_{1-\alpha} \cdot \hat{u}(x) \geq y \end{cases}$$
 (14)

Here, we include only the two non-dominated constraints i.e., those constraints not already implied by the others, to avoid redundancy. Note that these are linear constraints in terms of the estimated $\hat{h}(x)$ and $\hat{u}(x)$, which are assumed to be MIP-representable, and a constant $\hat{q}_{1-\alpha}$ computed offline during the conformalization procedure.

The classification setting introduces additional modeling complexity, as the inclusion of classes in the prediction set must be encoded using mixed-integer constraints. Specifically, the condition

$$\mathbb{1}\{-\hat{h}(x)^k \le \hat{q}_{1-\alpha}\} \quad \forall \ k \in \mathcal{K},\tag{15}$$

represents an indicator that is activated (i.e., equals one) if and only if class k is included in the prediction set. To enforce correct classification behavior, we require that at least one desired class $k \in \mathcal{K}^{\text{des}}$ is predicted (i.e., its indicator is activated), and no undesired class $k \in \mathcal{K}^{\text{und}}$ is allowed in the prediction set. To model this, we introduce a binary decision variable for each class $k \in \mathcal{K}$:

$$w_k = \begin{cases} 1 & \text{if class } k \text{ is included in the prediction set.} \\ 0 & \text{otherwise.} \end{cases} \quad \forall k \in \mathcal{K}$$
 (16)

Using these variables, the desired classification logic can be encoded via the following mixed-integer constraints:

$$C_{\alpha}(x) \subseteq \mathcal{Y} \iff -\hat{h}(x)^{k} > \hat{q}_{1-\alpha} \quad \forall \ k \in \mathcal{K}^{\text{und}}$$

$$\iff \begin{cases} -\hat{h}(x)^{k} - \hat{q}_{1-\alpha} \leq M(1 - w_{k}) & \forall \ k \in \mathcal{K} \\ \hat{h}(x)^{k} + \hat{q}_{1-\alpha} + \epsilon \leq Mw_{k} & \forall \ k \in \mathcal{K} \end{cases}$$

$$\iff \begin{cases} \sum_{k \in \mathcal{K}^{\text{des}}} w_{k} \geq 1 \\ w_{k} = 0 & \forall \ k \in \mathcal{K}^{\text{und}} \end{cases}$$

$$(17)$$

Here, ϵ is a small numerical tolerance (e.g., 10^{-6}) to ensure strict inequality, and M is a large positive constant used in the big-M reformulation. A practical choice for M is the maximum absolute value of the predicted logits observed in the calibration dataset, scaled by a safety factor (e.g., 4):

$$M = 4 \cdot \max_{i} \left| \hat{h}\left(x_{i}\right) \right| \tag{19}$$

For detailed discussion on selecting valid big-M values and their impact on computational performance, we refer the reader to Trespalacios and Grossmann [70].

Note that the conformal prediction sets $C_{\alpha}(x)$ used in our formulation of C-MICL have predictable, well-structured forms that are always MIP-representable, making the containment constraint $C_{\alpha}(x) \subseteq \mathcal{Y}$ tractable regardless of the underlying prediction model.

D Implementation Details

All computations were performed on a Linux machine running Ubuntu, equipped with eight Intel®, Xeon®, Gold 6234 CPUs (3.30 GHz) and 1 TB of RAM, utilizing a total of eight hardware threads.

D.1 Regression

This case study focuses on the optimal design and operation of a membrane reactor system used for the direct aromatization of methane. This integrated unit enables the conversion of methane into hydrogen and benzene, achieving simultaneous chemical reaction and product separation. By selectively removing hydrogen through a membrane, the reactor leverages Le Chatelier's principle to drive the equilibrium forward, resulting in higher methane conversion rates [66].

The optimization problem involves five key decision variables: the inlet volumetric flow rate of methane (v_0) , the inlet flow rate of sweep gas $(v_{\rm He})$, the operating temperature (T), the tube diameter (d_t) , and the reactor length (L). These variables were sampled uniformly within physically reasonable bounds, as shown below:

- $v_0 \sim \mathcal{U}(450, 1500) \text{ cm}^3/\text{h}$
- $v_{\rm He} \sim \mathcal{U}(450, 1500) \, {\rm cm}^3/{\rm h}$
- $d_t \sim \mathcal{U}(0.5, 2.0) \text{ cm}$
- $L \sim \mathcal{U}(10, 100) \text{ cm}$
- $T \sim \mathcal{U}(997.18, 1348.12) \text{ K}$

The resulting samples were used as initial conditions for solving the system of ordinary differential equations governing the reactor, enabling the computation of the outlet benzene flow ($F_{C_6H_6}$). The numerical integration was performed using code available in the opyrability repository [71]. To simulate measurement uncertainty, Gaussian noise was added to the computed benzene flows.

Table 1 shows an illustrative subset of the 1,000 data-points generated:

Table 1: Sampled decision variables and simulated benzene outlet flow.

v_0	(cm ³ /h)	$v_{\rm He}~({\rm cm^3/h})$	T(K)	d_t (cm)	L (cm)	$F_{\mathrm{C_6H_6}}$ (mol/h)
1	157.75	845.25	1272.93	0.57	30.67	37.45
7	73.89	1319.69	998.56	0.75	88.10	42.57
14	484.77	476.68	1201.21	1.64	77.38	40.70
8	17.68	536.65	1123.28	0.72	86.07	28.26
11	162.73	1262.03	1256.59	1.95	92.29	36.58

All hyperparameters used across the single-model baselines, wrapped approaches, and our proposed methods were kept consistent to ensure a fair comparison. For the Linear-Model Decision Tree, we set the maximum depth to five, the minimum number of samples required to split an internal node to ten, and the number of bins used for discretization to forty. For the Random Forest, we used fifteen estimators, a maximum depth of five, a minimum samples split of three, and considered sixty percent of the features when looking for the best split. The Gradient Boosting Tree model was configured with fifteen estimators, a learning rate of 0.2, a maximum depth of five, a minimum of five samples per split, and sixty percent of the features considered at each split. Lastly, the ReLU Neural Network was set up with two hidden layers of 32 units each, an L2 regularization strength of 0.01, and trained for 2000 epochs using Adam optimizer.

All models were cross-validated using a 5-fold split of their respective datasets. For the single base models, the complete 1,000 point dataset was used for training and evaluation. In the case of the ensemble methods, bootstrapping the 1,000 data-point dataset was employed to generate the required 500 data-point (half sized) subsets for training following [7]. For our conformal method, a split of the training data (800 data points) was used, ensuring that the calibration data was kept separate (200 data points). The Table 2 presents the average Mean Squared Error (MSE) across all folds for each of the methods, considering that the variance of the output for the 1,000 data-points is 1.2871.

Table 2: Average MSE across folds for all Methods for reactor model.

Predictive I	Method	
Model	Approach	MSE
RandomForest	MICL	0.1748
ReLU NN	MICL	0.0597
GradientBoosting	MICL	0.1192
LinearModelDT	MICL	0.0634
RandomForest	W-MICL(5)	0.1712
ReLU NN	W-MICL(5)	0.070
GradientBoosting	W-MICL(5)	0.1298
LinearModelDT	W-MICL(5)	0.0594
RandomForest	W-MICL(10)	0.1711
ReLU NN	W-MICL(10)	0.0699
GradientBoosting	W-MICL(10)	0.1239
LinearModelDT	W-MICL(10)	0.0568
RandomForest	W-MICL(25)	0.1696
ReLU NN	W-MICL(25)	0.0689
GradientBoosting	W-MICL(25)	0.1217
LinearModelDT	W-MICL(25)	0.0568
RandomForest	W-MICL(50)	0.1665
ReLU NN	W-MICL(50)	0.0685
GradientBoosting	W-MICL(50)	0.1194
LinearModelDT	W-MICL(50)	0.0559
RandomForest	C-MICL	0.1847
ReLU NN	C-MICL	0.0841
GradientBoosting	C-MICL	0.1402
LinearModelDT	C-MICL	0.0696

For all uncertainty models $\hat{u}(x)$, we employed ReLU-based Neural Networks with a shared architecture and regularization setup to maintain consistency across approaches. Specifically, each model used two hidden layers with 32 units each and an L_2 regularization (weight decay) coefficient of 0.001. The only hyperparameter that varied during cross-validation was the number of training epochs, which does not influence the MIP optimization outcomes but can affect the quality of the uncertainty estimates. The optimal number of training epochs identified through cross-validation were as follows: 1000 epochs for the model paired with Gradient Boosting, 3000 epochs for the one used with the ReLU Neural Network (using Adam optimizer), and 2000 epochs for both the Random Forest and Linear-Model Decision Tree variants. Table 3 displays the average MSE for the uncertainty model of each base predictor across all folds.

Table 3: Uncertainty average MSE for each base model.

Base Model	$\hat{u}(x)$ MSE
GradientBoosting	0.0142
ReLU NN	0.0247
RandomForest	0.0408
LinearModelDT	0.0225

Tables 4, 5, and 6 present the sets, parameters and decision variables of the problem, respectively.

Table 4: Sets used in the reactor optimization model.

Symbol	Description
\mathcal{I}	Set of decision variables
	$(e.g., \mathcal{I} = \{v0, v_He, T, dt, L\})$

Table 5: Parameters used in the reactor optimization model.

Symbol	Description
c_i	Operational or design cost coefficient associated with variable $i \in \mathcal{I}$

Table 6: Decision variables in the reactor design model.

Symbol	Description
x_i y	Value of design variable $i \in \mathcal{I}$, where $\mathcal{I} = \{ \texttt{v0, v_He, T, dt, L} \}$ Predicted outlet flow of C_6H_6 from the surrogate model

Formulation.

$$\min_{\{x_i\}_{i\in\mathcal{I}},y} \quad \sum_{i\in\mathcal{I}} c_i x_i \tag{D.1a}$$

s.t.
$$10 \le \frac{x_{\rm L}}{x_{\rm dt}} \le 150$$
 (D.1b)

$$0.75 \le \frac{x_{\text{v0}}}{x_{\text{v_He}}} \le 3.0$$
 (D.1c)

$$x_{\text{v_He}}$$
 $20 \le \frac{x_{\text{v0}}}{x_{\text{L}}} \le 120$ (D.1d)

$$x_{\text{v0}} \le 1.1 \cdot x_{\text{T}} \tag{D.1e}$$

$$\hat{h}(\boldsymbol{x}) = y \tag{D.1f}$$

$$y \ge 50 \tag{D.1g}$$

$$x_i \ge 0 \quad \forall i \in \mathcal{I}$$
 (D.1h)

Explanation of Constraints.

- (D.1a): Minimize the operating cost as a linear function of the decision variables.
- (D.1b): Enforce physical bounds on the tube length-to-diameter ratio.
- (D.1c): Maintain a suitable gas feed ratio between CH₄ and He.
- (D.1d): Control the residence time via the CH₄ flow and reactor length.
- (D.1e): Limits methane flow rate based on temperature to ensure safe and stable operation.
- (D.1f): Enforces that the surrogate model output y, predicting C_6H_6 flow, is computed from the design and operation variables x.
- (D.1g): Requires that the predicted C_6H_6 outlet flow (i.e., product quality) meets or exceeds the target value of 50.
- (D.1h): Enforces nonnegativity for all design variables.

D.2 Classification

We now examine the food basket optimization problem introduced earlier and grounded in the work of Peters et al. [68]. Our analysis is based on the model developed by Fajemisin et al. [1], which aims to minimize the cost of assembling a basket consisting of 25 different commodities while satisfying nutritional requirements across 12 key nutrients. An additional constraint is placed on the palatability of the basket, which reflects how acceptable or appealing the food is to the target population. In line with the case study by Maragno et al. [7], we require a minimum palatability score of t=0.5 to ensure that the resulting food baskets are not only affordable and nutritionally adequate but also culturally and socially acceptable. The dataset used in this analysis is publicly available at and published by Maragno et al. [7] here.

The palatability scores initially range continuously from 0 to 1. To transform this into a classification problem, we discretize these scores into four distinct categories: *bad*, *regular*, *good*, and *very good*. This discretization is achieved by setting thresholds at 0.25, 0.5, and 0.75. Consequently, the problem

becomes a multi-class classification task, where each food basket is assigned one of these categorical labels based on its palatability score. We restrict the optimization problem to include only food baskets that fall within the good and $very\ good$ categories. This approach is inspired by the work of Maragno et al. [7], where a similar threshold of t=0.5 was used, but applied in a categorical context. Table 7 displays a sample of the dataset available.

Table 7: Sample of food basket dataset with commodity quantities and palatability score. Full dataset available here.

Beans	Bulgur	Cheese	Fish	Meat	CSB	Dates	DSM	 Palatabi	lity Class
0.7226	0	0	0	0	0	0	0.5987	 0.199	Bad
0.7860	0	0	0	0	0.0419	0	0.3705	 0.8049	Very Good
0.4856	0	0	0	0	0	0	0.2696	 0.6517	Good
0	0	0.5734	0	0	0	0	0.0025	 0.3220	Regular

In this case, we exclusively trained ReLU-based neural networks across all methods (single-model baseline, wrapped approach, and our proposed method) ensuring that all hyperparameters remained identical for a fair comparison. For the predictors, the network architecture consisted of three hidden layers with 64 units each, trained for 500 epochs with a weight decay (L2 regularization) of 0.01. Additionally, we trained an oracle model using all available data. This oracle network was configured with five hidden layers of 256 units each, trained for 1000 epochs, also with a weight decay of 0.01.

Table 8 displays the average accuracy across validation folds on the data available for each approach.

Table 8: Average accuracy across folds for all methods for the basket model.

Predict		
Model	Approach	Accuracy [%]
ReLU NN ReLU NN ReLU NN	MICL W-MICL(5) W-MICL(10)	84.32 78.24 76.88
ReLU NN	C-MICL	83.30

Tables 9, 10, and 11 present the sets, parameters and decision variables of the problem, respectively.

Table 9: Sets used in the basket model.

Symbol	Description
\mathcal{M} \mathcal{L} $\mathcal{K}^{\mathrm{des}}$ $\mathcal{K}^{\mathrm{und}}$	Set of commodities (e.g., $\mathcal{M} = \{\text{rice, beans, salt, sugar, } \ldots\}$) Set of nutrients (e.g., $\mathcal{L} = \{\text{protein, iron, calories, } \ldots\}$) Set of desired categories (e.g., $\mathcal{K}^{\text{des}} = \{\text{good, very good}\}$) Set of undesired categories (e.g., $\mathcal{K}^{\text{und}} = \{\text{bad, regular}\}$)

Table 10: Parameters used in the basket model.

Symbol	Description
$\begin{array}{c} \text{Nutreq}_l \\ \text{Nutval}_{ml} \\ p_m \\ M \end{array}$	Nutritional requirement for nutrient $l \in \mathcal{L}$ (grams/person/day) Nutrient content of commodity m for nutrient l (grams per gram) Procurement cost of commodity m (in \$/metric ton) Large positive constant for big-M reformulation.

Table 11: Decision variables in the basket model.

Symbol	Description
$x_m \\ y \\ w_k$	Quantity of commodity $m \in \mathcal{M}$ in the food basket (grams) Palatability level of the food basket Binary indicator for predicting desired category $k \in \mathcal{K}^{\text{des}}$ ({0,1})

Formulation.

$$\min_{\{x_m\}_{m \in \mathcal{M}}, \ y} \quad \sum_{m \in \mathcal{M}} p_m x_m$$
 (D.2a) s.t.
$$\sum_{m \in \mathcal{M}} \text{Nutval}_{ml} \cdot x_m \ge \text{Nutreq}_l, \qquad \forall \ l \in \mathcal{L}$$
 (D.2b)

s.t.
$$\sum_{m \in \mathcal{M}} \text{Nutval}_{ml} \cdot x_m \ge \text{Nutreq}_l, \qquad \forall l \in \mathcal{L}$$
 (D.2b)

$$x_{\text{salt}} = 5 \tag{D.2c}$$

$$x_{\text{sugar}} = 20 \tag{D.2d}$$

$$\hat{h}(\boldsymbol{x}) = \boldsymbol{y} \tag{D.2e}$$

$$y_{k'} - y_k \le M(1 - w_k)$$
 $\forall k \in \mathcal{K}^{\text{des}}, k' \in \mathcal{K}^{\text{und}}$ (D.2f)

$$\sum_{k \in \mathcal{K}^{\text{des}}} w_k \ge 1 \tag{D.2g}$$

$$x_m \ge 0$$
 $\forall m \in \mathcal{M}$ (D.2h)

Explanation of Constraints.

- (D.2a): Objective function minimizing total procurement cost of the food basket.
- (D.2b): Nutritional constraints to ensure daily nutrient requirements are met.
- (D.2c)-(D.2d): Fixed amounts of salt and sugar imposed (e.g., due to guidelines).
- (D.2e): Palatability is computed as a function $\hat{h}(x)$ of the selected commodity quantities.
- (D.2f): Big-M constraint that activates an indicator if the logit of a desired class $k \in \mathcal{K}^{\text{des}}$ is higher than the logit of all undesired classes $k' \in \mathcal{K}^{und}$, representing desired prediction. M is calculated using the largest value in magnitude observed in calibration data times a enlarging safety factor (e.g., 4) as $M = 4 \cdot \max_{i} |\hat{h}(x_i)|$. For more information on how to calculate valid M values and their impacts in optimization we refer the readers to Trespalacios and Grossmann [70].
- (D.2g): Enforce at least one desired class to have larger logit that both undesired classes.
- (**D.2h**): Ensures quantity nonnegativity.

Complementary Results

To quantify uncertainty and support the statistical reliability of our results, we report 95% confidence intervals (Cls) for all key metrics: empirical feasibility rates, optimization times, relative differences in objective value, and true coverage. All uncertainty comes from 100 problem instances of optimization problems.

For feasibility rates and true coverage plots intervals are computed as the proportion of feasible solutions over meaning that we model each as a Bernoulli random variable and apply the standard error formula for proportions:

$$SEM = \sqrt{\frac{\bar{p}(1-\bar{p})}{n}},$$

where \bar{p} is the sample mean and n is the number of problem instances. The corresponding confidence interval is calculated using the Student's t-distribution:

$$CI_{95\%} = \bar{p} \pm t_{n-1,0.975} \cdot SEM$$

For optimization times and relative differences in objective values, we compute the sample mean, standard deviation, and standard error of the mean (SEM) for each method across instances. The 95% confidence intervals are then given by:

$$\text{CI}_{95\%} = \bar{x} \pm t_{n-1,0.975} \cdot \frac{s}{\sqrt{n}}$$

where s is the sample standard deviation and n=100 is the number of observations. This method assumes approximate normality of the sample means, which is reasonable due to the Central Limit Theorem given the sample size.

All error bars shown in plots correspond to these 95% confidence intervals. The factors of variability captured are due to the random generation of optimization instances and the resulting performance metrics across different methods. These intervals are reported directly in the results section and supporting figures to substantiate claims about statistical significance and performance differences.

E.1 Regression

The following results summarize performance on the reactor case study at $\alpha=10\%$, evaluated after solving 100 optimization instances for each method. More specifically, Figure 7 reports the average computational time required to obtain an optimal solution for each method at $\alpha=10\%$, while Figure 8 presents the relative difference in optimal objective value between baseline methods and our proposed C-MICL at the same confidence level. Lighter bars in both figures denote methods that failed to achieve the target empirical feasibility rate, as established in Figure 1. Additionally, Figure 9 shows the distribution of true constraint values h(x) for each methods for $\alpha=10\%$. Finally, Figure 10 reports the empirical coverage over 1,000 out-of-sample data points, stratified by deciles of the true output variable y. For each decile, we report the proportion of instances where the true value lies within the predicted interval, demonstrating strong coverage across the output space for all baseline models. Notably, our ReLU NN-based conformal sets empirically satisfy the ground-truth feasibility coverage guarantee from Lemma 3.1, achieving coverage rates of 89.31% and 95.38% on ground-truth infeasible and feasible samples, respectively.

We report here the full results for the reactor case study at $\alpha=5\%$, based on evaluations conducted over 100 solved optimization instances per method. Specifically, Figure 11 reports the empirical ground-truth feasibility rate achieved by each MICL approach. Figure 12 shows the average computational time needed to solve each method at $\alpha=5\%$, and Figure 13 displays the relative difference in optimal objective value between baseline methods and our C-MICL approach. Furthermore, Figure 14 presents the distribution of true constraint values h(x) for each method at $\alpha=5\%$. Finally, Figure 15 reports the empirical coverage of the models over 1,000 out-of-sample data points, grouped by deciles of the true output y. Coverage is measured as the proportion of points for which the true value of y falls within the predicted interval at a target level of $\alpha=5\%$. Our ReLU NN-based conformal sets demonstrate empirical compliance with the ground-truth feasibility coverage guarantee established in Lemma 3.1, yielding coverage rates of 96.92% for ground-truth feasible instances and 96.09% for ground-truth infeasible ones.

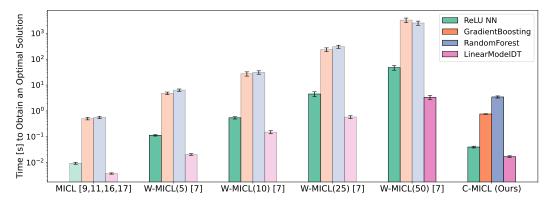


Figure 7: Average computational time to obtain an optimal solution for all methods at $\alpha=10\%$. Our proposed C-MICL approach is comparable to single-model baselines while being orders of magnitude faster than ensemble heuristics using the same underlying base model. Lighter bars indicate methods that failed to achieve the target empirical feasibility rate.

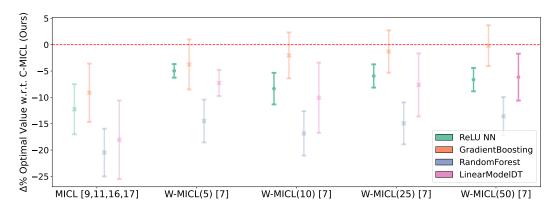


Figure 8: Relative difference in optimal objective value of baseline methods compared to our proposed C-MICL at $\alpha=10\%$. Our approach exhibits a small difference of around 8% compared to the method that achieved empirical coverage, indicating that C-MICL attains comparable solution quality relative to valid approaches. Lighter bars marked with an "x" denote methods that failed to meet the target empirical feasibility rate.

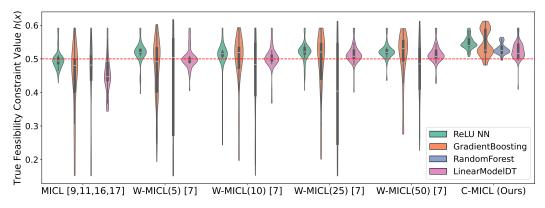


Figure 9: Distribution of true constraint values h(x) for each method for $\alpha=10\%$ The dotted line at 0.5 marks the lower bound imposed in the oracle constraint values below this line correspond to true infeasibilities. C-MICL not only yields the fewest violations but also achieves the smallest violation magnitudes when they occur.

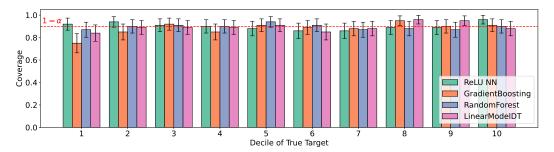


Figure 10: Empirical out-of-sample coverage across predictive models at $\alpha=10\%$, evaluated on 1,000 test points and stratified by deciles of the true output y. Our conformal sets achieve strong target coverage. Crucially, the ground-truth feasibility threshold ($y \geq 50\%$) lies within the 9th and 10th deciles, where all methods attain valid empirical coverage, empirically supporting the assumptions of Theorem 4.1.

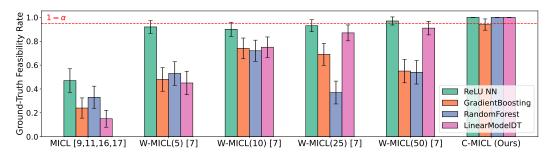


Figure 11: Empirical ground-truth feasibility rate of optimal solutions across MICL methods on 100 optimization problem instances for $\alpha=5\%$. Our Conformal MICL approach (rightmost bars) consistently meets the target feasibility rate of $\geq 95\%$ regardless of the underlying base model. In contrast, existing methods show inconsistent performance and often fall short of the theoretical guarantee.

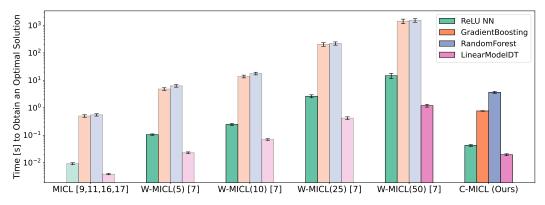


Figure 12: Average computational time to solve 100 optimization instances at $\alpha=5\%$. Our proposed C-MICL method matches the speed of single-model baselines and is orders of magnitude faster than ensemble-based heuristics using the same base model. Lighter bars denote methods that did not reach the target empirical feasibility threshold.

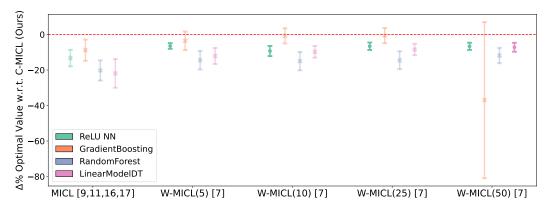


Figure 13: Relative difference in optimal objective value between baseline methods and our proposed C-MICL at $\alpha=5\%$ across 100 optimization problems. Our approach shows a modest 8% difference compared to the method that achieved empirical coverage, demonstrating that C-MICL provides solution quality on par with other valid approaches. Lighter bars with an "x" indicate methods that did not meet the empirical feasibility target. Five outliers were removed for the W-MICL(50) Gradient Boosted Tree model, though the comparison remains statistically insignificant with or without them.

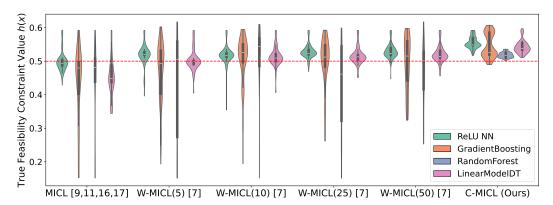


Figure 14: Violin plots of the true constraint values h(x) for each method at $\alpha=5\%$. The dotted line at 0.5 denotes the lower bound enforced in the oracle constraint where values falling below this line indicate true infeasibilities. At this stricter confidence level, C-MICL maintains the highest reliability, producing the fewest and smallest violations across all methods.

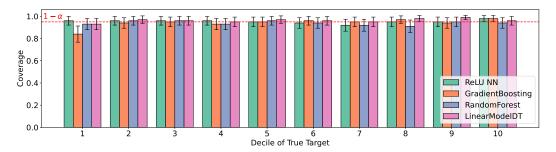
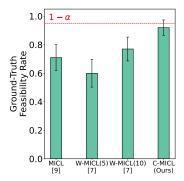
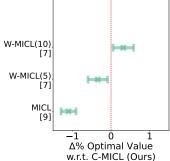
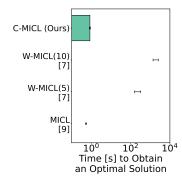


Figure 15: Empirical out-of-sample coverage at $\alpha=5\%$ for all predictive models, evaluated on 1,000 test points and stratified by deciles of the true output y. Our conformal prediction sets achieve the desired coverage level across deciles. Notably, the ground-truth feasibility threshold $(y\geq 50)$ lies in the 9th and 10th deciles, where all methods achieve valid empirical coverage, providing empirical support for the assumptions in Theorem 4.1.







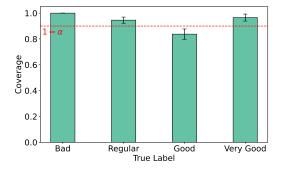
short of this benchmark.

Figure 16: Empirical ground- Figure 17: Relative difference truth feasibility rates across in optimal objective value be-MICL methods on 100 optimiza- tween baseline MICL methods tion problem instances at $\alpha = \text{ and C-MICL across } 100 \text{ opti-}$ 5%. C-MICL (rightmost bar) mization problem instances at is the only method that consis- $\alpha=5\%$. The differences avertently achieves the target feasi- age about 1%, with all methods bility threshold of at least 95%, achieving statistically similar soin line with theoretical guaran- lution quality. However, none tees. All baseline methods fall of the baselines produce implementable solutions, as indicated by the lighter bars, which represent methods that failed to meet the required empirical coverage guarantee. Only our method satisfies this guarantee.

Figure 18: Average time to compute an optimal solution for each MICL method on 100 optimization instances at $\alpha = 5\%$. C-MICL matches the runtime of single-model MICL and outperforms ensemble-based methods by a wide margin. Lighter bars represent approaches that did not achieve the target feasibility

E.2 Classification

The following three figures present the performance of different MICL methods applied to the food basket design problem under a classification setting, with $\alpha = 5\%$. These figures assess the feasibility, objective value, and computational efficiency of the methods across 100 problem instances, highlighting the strengths of C-MICL in comparison to baseline approaches. Specifically, Figure 16 illustrates the empirical ground-truth feasibility rate, Figure 17 shows the relative difference in optimal objective value, and Figure 18 reports the average computational time required to obtain an optimal solution. Finally, Figures 19 and 20 present the empirical coverage achieved by each model across 1,000 out-of-sample points, grouped by the label y category, for $\alpha = 10\%$ and $\alpha = 5\%$, respectively. These figures report the proportion of instances in which the true value of y falls within the predicted conformal set. Our conformal sets empirically satisfy the ground-truth feasibility coverage guarantee from Lemma 3.1, achieving coverage rates of 88.22% (feasible) and 96.19% (infeasible) for $\alpha = 10\%$, and 93.41% (feasible) and 98.79% (infeasible) for $\alpha = 5\%$.



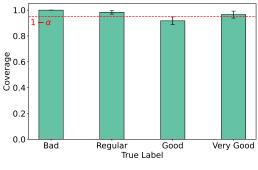


Figure 19: Empirical out-of-sample coverage, Figure 20: Evaluation of empirical coverage with valid empirical coverage, lending support to the viding evidence consistent with Theorem 4.1. assumptions in Theorem 4.1.

evaluated over 1,000 test points and stratified by ReLU NNs over 1,000 test samples, grouped by the true categorical label of y at $\alpha = 10\%$. Across the true y label, for a target level of at $\alpha = 5\%$. all categories, the predicted conformal sets cap- The model produces conformal prediction sets ture the true label at the desired coverage level. that contain the true class label at the expected Notably, for the "Good" and "Very Good" cate- rate. In particular, the "Good" and "Very Good" gories, (which correspond to the ground-truth fea- categories, those relevant to enforcing feasibility, sibility region) the models approximately achieve demonstrate approximate empirical coverage, pro-