

# Decision-Focused Learning for Complex System Identification: HVAC Management System Application

Pietro Favaro  
Jean-François Toubeau  
François Vallée

Power Systems and Markets Research Group  
University of Mons  
Mons, Belgium  
pietro.favaro@umons.ac.be  
jean-francois.toubeau@umons.ac.be  
francois.vallee@umons.ac.be

Yury Dvorkin

Ralph O'Connor Sustainable Energy Institute  
Department of Electrical and Computer Engineering  
Department of Civil and System Engineering  
Johns Hopkins University  
Baltimore, Maryland, US  
ydvorki1@jhu.edu

## Abstract

As opposed to conventional training methods tailored to minimize a given statistical metric or task-agnostic loss (e.g., mean squared error), Decision-Focused Learning (DFL) trains machine learning models for optimal performance in downstream decision-making tools. We argue that DFL can be leveraged to learn the parameters of system dynamics, expressed as constraint of the convex optimization control policy, while the system control signal is being optimized, thus creating an end-to-end learning framework. This is particularly relevant for systems in which behavior changes once the control policy is applied, hence rendering historical data less applicable. The proposed approach can perform system identification – i.e., determine appropriate parameters for the system analytical model – and control simultaneously to ensure that the model’s accuracy is focused on areas most relevant to control. Furthermore, because black-box systems are non-differentiable, we design a loss function that requires solely to measure the system response. We propose pre-training on historical data and constraint relaxation to stabilize the DFL and deal with potential infeasibilities in learning. We demonstrate the usefulness of the method on a building Heating, Ventilation, and Air Conditioning day-ahead management system for a realistic 15-zone building located in Denver, US. The results show that the conventional RC building model, with the parameters obtained from historical data using supervised learning, underestimates HVAC electrical power consumption. For our case study, the ex-post cost is on average six times higher than the expected one. Meanwhile, the same RC model with parameters obtained via DFL underestimates the ex-post cost only by 3%.

## Keywords

Decision-Focused Learning, Task-Aware Learning, End-to-End Learning, Online Learning, System Identification, Building Management System

## List of Acronyms

DFL	Decision-Focused Learning.
HVAC	Heating, Ventilation, and Air-Conditioning.
ITO	Identify-Then-Optimize.
LP	Linear Programs.

MAE	Mean Absolute Error.
ML	Machine Learning.
MSE	Mean Squared Error.
RC	Resistance-Capacitance.
RL	Reinforcement Learning.

## 1 Introduction

Optimization is used to devise control policies such as for Heating, Ventilation, and Air Conditioning (HVAC) in buildings [5], the guidance of vehicles [35], and the planning of complex dynamic systems such as day-ahead energy scheduling of active buildings [28], or energy management systems for micro-grid users [27]. The optimization problem solution is the actions or decisions to apply to the system. However, an analytical model of the system dynamics is always assumed to be known and formulated within an optimization problem, usually as constraints. The necessity to have an analytical model hinders control tasks where physics-based models are unpractical; either because the system is a black box, or physical models are far too complex. Data-driven surrogates become impractical when system dynamics change drastically due to a newly applied control policy, rendering prior historical observations incomplete and minimally informative. Therefore, the conventional two-stage Identify-Then-Optimize (ITO) approach is inefficient.

This paper addresses the ITO approach inefficiencies using the day-ahead HVAC scheduling in buildings as a use case. HVAC scheduling becomes increasingly important because of the roll-out of smart meters and dynamic electricity tariffs [40]. The underlying goal of such dynamic (e.g., time-of-use) tariffs is to prompt more flexibility at the consumption level to accommodate more renewable-based, and thus often weather-dependent and uncertain, generation. Unlike lightning or essential equipment that needs to be available on demand, HVAC can be scheduled and benefit from the building thermal inertia to shift the consumption across the day. Consequently, building managers can save money by optimizing HVAC operations. Moreover, the building sector accounts for 75% of the final electricity consumption in the USA, in which HVAC consumption contributes 40% [15, 30]. Therefore, if minimizing greenhouse gas emissions, HVAC scheduling can have a significant impact.

## 1.1 Building Thermodynamics Modeling

The goal of day-ahead HVAC scheduling is to find the optimal temperature set point profiles for each building thermal zone to minimize the electricity cost for the upcoming day, alleviate the potential strain on the power grid, and ensure thermal comfort. Therefore, building management systems rely on a thermodynamic model of the building that relates the HVAC operation to the indoor temperature.

Historically, the thermodynamic model has always been assumed to be known. Most HVAC models are physics-based (i.e., white-box) [2]. These models contain major assumptions. Moreover, some building parameters, such as the materials resistance and capacitance, are intrinsically unknown because of aging and installation process which may cause large error [13]. Last, it requires expertise in HVAC modeling, which may not be readily available at each building site.

On the other end of the spectrum lies fully data-driven models (i.e., black-box). Neural networks were widely used for system identification [12], but cannot be easily integrated into optimization models; otherwise, the problem becomes non-linear and there is no guarantee on the solution quality [17, 39]. Furthermore, data-driven models perform poorly out of their training zone [2]. This hinders day-ahead HVAC scheduling applications since their goal is to explore the whole feasible domain to find the best course of action. Reinforcement Learning (RL) has proven to be effective for real-time set point control [38], emission minimization [19], and peak demand reduction [32]. However, inputs and assumptions underlying RL methods are often impractical. For instance, the active RL algorithm in [18] assumes real-time measurement of the occupants' clothing value, number, and thermal comfort. A review concludes that model-based deep RL must be favored to leverage some prior knowledge [38]. Similarly, the review by [16] states that the grey-box models are the most promising.

Grey-box models involve physics-based formulation in which parameters are obtained through data-driven techniques. Grey-box models require less data than black-box models [1]. Among such models, the Resistance-Capacitance (RC) model (also named lumped capacitance model or network model), a reduced order physics model, has exhibited good performances for buildings where there is no important indoor air convection [22]. Furthermore, the RC model is linear, thus lending convexity to the optimization problem. However, estimating the parameters of the model is challenging. A first option to estimate the parameters is to analyze the building materials and select default tabular data. In addition to requiring building blueprints, this method does not account for manufacturing dispersion (i.e., variations in product dimensions, properties, or performance due to inconsistencies in the production process), the on-site installation process, and the aging of materials. Therefore, data-driven parameter identification has led to significantly improved RC models [7].

## 1.2 Decision-Focused Learning

In this paper, we adopt a classical RC model. Aware of the limitations in estimating the parameters on historical data, we learn the parameters in a task-aware manner. In [20], the authors were the first to suggest training a Machine Learning (ML) model in

a directed manner based on its impact on the downstream decisions. The weights of the ML model that is used to predict the uncertain parameters of the downstream optimization problem are learnt to minimize the decision error induced by the parameters misestimation. To that end, the authors established the gradient of the optimization problem solution with respect to the problem parameters. This enables backpropagation through the optimization problem. Only unconstrained quadratic problems were considered. By allowing loss functions that depend explicitly on downstream optimization decisions (and implicitly on the ML model output), the calculation of optimization problem gradient paves the way for Decision-Focused Learning (DFL). The framework was then extended to stochastic quadratic problems [11]. Because the gradient of unconstrained problems is easier to compute, the constraints were relaxed in the objective function.

Bounded Linear Programs (LP) and Combinatorial Program (CP) lead to zero-valued gradient almost everywhere, and undefined elsewhere. Geometrically, the solution of LP belongs to the set of the vertices of the feasible space polytope [6]. Consequently, an infinitesimal change in the parameter values has no impact on solution; i.e., the solution remain on the same vertex. However, as soon as parameter changes become large enough, the LP solution jumps to another vertex. This results in a piecewise constant objective function, leading to a gradient that is either zero or undefined, making it uninformative for gradient descent training. For CP, the discrete nature of the decision variables leads to the same observation as for LP. The first method to calculate the gradient of an LP is to calculate the gradient of its Karush-Kuhn-Tucker (KKT) conditions. The gradient of the KKT conditions is also null or undefined since the KKT conditions are an exact representation of the same LP. To address this issue, the objective function can be augmented by a  $L_2$  regularization of the decisions variables. This turns the LP into a strong concave or convex quadratic program if it is a maximization or minimization, respectively [36]. The second approach is to design a surrogate loss function. Smart "Predict, then Optimize" Plus (SPO+) is such a convex function [14]. Interestingly, SPO+ extend beyond LP and can be applied to any optimization problem with a linear objective function and linear, convex, and integer constraints. However, the uncertain parameters cannot appear in the constraints. Moreover, SPO+ requires the optimal decisions to be known, which may not always be the case. The third option for LP and CP is stochastic smoothing. A random perturbation is applied to uncertain parameters to smooth the loss function transition and, thus, create informative gradients from the otherwise null and undefined derivative [31]. Unlike SPO+, the uncertain parameters can appear in constraints.

Leveraging the differentiability of conic programs, Agrawal et al. demonstrated how controller parameters within the objective function can be learned assuming full knowledge and differentiability of the system, and observation of the system dynamics [4]. It was applied to the forecasting of electricity prices through wind power prediction [34]. It was also employed to split the demand for flexibility from the transmission grid towards the flexibility of the distribution grid assets [25].

To the best of the authors' knowledge, no prior work has simultaneously addressed system identification (i.e., learning constraint parameters, specifically the parameters of the dynamics model)

and optimization, nor has it considered scenarios where the controlled system is both unknown and non-differentiable. Tackling this challenge requires a profound understanding of optimization and decision-focused learning theory, combined with extensive domain-specific expertise. The only previous work that does not assume knowledge of the system dynamics and differentiability is [29] in which convex surrogate for both the optimization problem and the performance loss is constructed to train a ML forecaster. However, the construction of the surrogate model requires knowing the true value of the uncertain parameters.

### 1.3 Contributions

This paper presents three major contributions.

First, we leverage advances in decision-focused learning to perform complex system identification and control simultaneously. Compared to [4] and [14], uncertain parameters to be learned are within the constraints of the optimization control policy. Due to end-to-end learning, the system model parameters are optimized for operating zones relevant to the control policy. This framework bypasses the need for extensive high-quality historical database reflecting control data distributions that are rarely available in practice.

Second, we formulate the performance loss that is necessary to compute the quality of the decisions and backpropagate the gradient with respect to the dynamics model parameters. Unlike in [4], we do not assume that the system is differentiable. The only requirement is for the system state variables used in the system dynamics model to be observable. In addition, we introduce a hierarchical loss to inform the learning about the HVAC system structure and its impact on the objective value.

Third, training is made robust to infeasibilities by relaxing constraints on the dynamics model output. Moreover, feasibility in the early stages of DFL training is prompted by performing a warm-start with a model pre-trained on historical data. A small noise is added on the parameters after pre-training to get the model out of the local minimum while retaining a maximum of prior knowledge.

The performance and relevance of the method is illustrated on the day-ahead HVAC scheduling of a realistic 15-zone building located in Denver, CO, US [33]. Furthermore, the method's robustness is evaluated under a distribution shift in the input parameters. Specifically, we analyze the model's performance metrics on a dataset representing an exceptionally hot year, a scenario likely to become increasingly common due to global warming.

### 1.4 Outline

Section 2 describes the problem class, explains gradient computations necessary to learn the parameters, the loss function design for non-differentiable black-box systems, and training robustness enabled by constraint relaxation. Afterwards, section 3 presents the case study. In particular, sections 3.2 and 3.4 describe the specific optimization problem, and the loss function for the day-ahead HVAC scheduling, respectively. We report the results in section 3.5. Section 4 discusses the results and limitations of the case study. Section 5 concludes the paper and outlines future work directions.

## 2 A Method for Simultaneous System Identification and Control

The proposed framework performs system identification and control simultaneously. To that end, we start by describing the addressed problem class, then come the explanation about the gradient computation needed for updating the parameters. Afterwards, we propose a DFL supervised loss that bypasses the need for a differentiable system. Lastly, we address the infeasibility that may arise during training. Figure 1 provides an overview of the proposed method.

### 2.1 Problem Class

Given a system with dynamics described by the unknown state transition function:

$$x_{t+1} = f(x_t, u_t, w_t), \forall t, \quad (1)$$

the goal is to learn the parameters  $\theta$  of a proxy model  $\hat{f}(\hat{x}_t, u_t; \theta)$  that represents the system dynamics, while keeping the formulation of the control policy  $\phi$  convex (2). The known system state at a given time step  $t \in \mathcal{T}$  is given by  $x_t \in \mathbb{R}^n$ , the expected state by  $\hat{x}_t \in \mathbb{R}^n$ , and the command or action by  $u_t \in \mathbb{R}^m$ . The true state transition function can be affected by some noise  $w_t$ . We do not make any assumption on  $w_t$ .

A general formulation of the convex optimization control policy  $\phi$  is given by:

$$\begin{aligned} \phi(x_0) = \underset{u}{\operatorname{argmin}} f_o(\hat{x}, u) \\ \text{s.t. } g_i(\hat{x}, u) \leq 0 \quad i = 1, \dots, k, \forall t, \\ h_i(\hat{x}, u) = 0 \quad i = 1, \dots, l, \forall t, \\ \hat{x}_{t+1} = \hat{f}(\hat{x}_t, u_t; \theta), \forall t. \end{aligned} \quad (2)$$

To keep the problem convex, the equality constraints  $h_i$  and  $\hat{f}$  must be affine with respect to the optimization variables  $\hat{x}$  and  $u$  whereas  $g_i$  must be convex. The convexity of the problem is crucial for the gradient computation as explained in section 2.2. In addition, convexity guarantees that the solution is the global optimum and offers tractable problems well handled by off-the-shelf solvers.

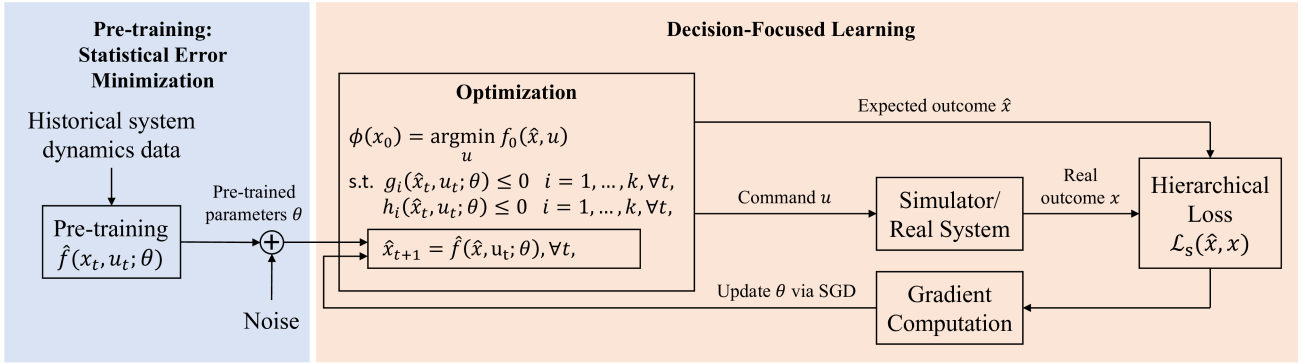
Even though we focus on control optimization policy problems and state transition functions, the proposed framework can be used to learn any parameters of a convex optimization problem.

### 2.2 Gradient Computation

Let us assume we have a scalar function  $\mathcal{L}$  that quantifies the loss of performance (i.e., the suboptimality) caused by the misestimated parameters  $\theta$ . The parameters  $\theta$  can be learned while controlling the system with the policy  $\phi$  by minimizing  $\mathcal{L}$ . This requires computing the gradient  $\frac{\partial \mathcal{L}}{\partial \theta}$ . Then,  $\theta$  can be updated iteratively by any gradient-based optimization algorithm. Misestimating  $\theta$  leads to suboptimal commands  $u$  that are responsible for the performance loss. By applying the chain rule, the gradient becomes:

$$\frac{\partial \mathcal{L}}{\partial \theta} = \frac{\partial \mathcal{L}}{\partial u} \frac{\partial u}{\partial \theta}. \quad (3)$$

The first factor accounts for the performance loss caused by a variation of the commands  $u$ . The second factor creates the need for differentiating through the optimization problem. More specifically,



**Figure 1: The proposed framework starts by pre-training the system dynamics model on historical data. Then, the pre-trained model parameters are fed into the convex optimization control policy where they keep being updated on representative scenarios to minimize the loss metric evaluating the mismatch between the expected and observed system states.**

the impact of a variation in the optimization parameters  $\theta$  on the optimization output  $u$  must be determined.

An optimization problem can be considered as a mapping  $O$  between its parameters and the optimal solution. Differentiating such a problem is challenging because of the complex and implicit mapping definition. Moreover, some classes of constrained programs, such as the linear programs, define piecewise-constant mappings in which the gradient is often null or nonexistent [23].

A general framework, named *Cvxpylayers*, is able to differentiate conic programs [3]. Since any convex program can be recast as a conic program [24], *Cvxpylayers* is very versatile. A conic program is of the form

$$\begin{aligned} \min_x \quad & c^T x \\ \text{s.t.} \quad & b - Ax \in \mathcal{K}, \end{aligned} \quad (4)$$

where  $(A, b, c) \in (\mathbb{R}^{m \times n}, \mathbb{R}^m, \mathbb{R}^n)$  are the parameters of the problem,  $x \in \mathbb{R}^n$  is the primal decision variable,  $\mathcal{K} \subseteq \mathbb{R}^m$  is a nonempty, closed, convex cone. Under the assumption that (4) has a unique solution, its gradient can be computed. Solving a conic program can be done in three steps.

- (i) The data parameters  $(A, b, c)$  are used to build the corresponding skew-symmetric matrix  $Q$  [26, 37].
- (ii) The self-homogeneous dual embedding problem – which reformulates the KKT conditions of the original conic program into a single system of equations and incorporates the matrix  $Q$  – is solved by finding a root (i.e., a zero) to its normalized residual map  $s$  [8]. The residual map is a function that quantifies the difference (or residual) between the current solutions of the primal and dual problems and the optimal solution. Therefore, it guides the search for the optimal solution, directing the optimization process.
- (iii) The solution  $z$  is retrieved by  $R$  and mapped to a valid solution of the original problem.

Therefore, the mapping  $O$  can be seen as the composition of three submappings  $R \circ s \circ Q$ . By the chain rule, we have

$$DO(A, b, c) = DR(z) Ds(Q) DQ(A, b, c) \quad (5)$$

with  $D$ , the derivative operator. The derivative of  $Q$  is straightforward since  $A$ ,  $b$ , and  $c$  appear explicitly in its formulation. The derivative of the root-finding problem  $s$  is obtained via implicit differentiation. The derivative of the retriever  $R$  relies on the derivative of the euclidean projection of the self-homogeneous dual embedding solution  $z$  onto the feasible cone  $\mathcal{K}^*$  associated with the dual of the original conic problem (4) [3].

### 2.3 Performance Loss for Non-Differentiable System

As established by (3), the gradient of the performance loss with respect to the parameters  $\theta$  must be defined and exist. The ideal loss prescribed in the literature is the regret  $\mathcal{L}_r$  (6) [14, 23]. The regret assesses the objective value loss at optimality (marked by  $*$ ) caused by the error between the actual  $\theta$  and its estimator  $\hat{\theta}$ :

$$\mathcal{L}_r = f_o^*(\hat{\theta}) - f_o^*(\theta). \quad (6)$$

However, the definition of regret assumes the actual value of  $\theta$  is known. Other loss functions have considered the transition state function  $f$  (1) as fully known and differentiable [4].

Many physical systems must be considered as black boxes, and thus non-differentiable, because of their inherent complexity and the unbearable cost to build an accurate system model. The gradient of black-box model can be estimated via simultaneous perturbation stochastic approximation [9]. However, it is unpractical if the system environment conditions are evolving (e.g. the weather) making the repetition of different commands in the same conditions impossible.

Consequently, we design a novel performance loss  $\mathcal{L}_{\text{DFL}}$  that does not assume any knowledge about the system:

$$\mathcal{L}_{\text{DFL}} = \mathcal{L}_s(\hat{x}, x), \quad (7)$$

where  $\mathcal{L}_s$  refers to any supervised loss such as the Mean Squared Error (MSE) loss or the Mean Absolute Error (MAE) loss. The loss  $\mathcal{L}_{\text{DFL}}$  computes a statistical metric on the error between the expected system state  $\hat{x}$ , as predicted by the optimization model, and the actual observed realisation  $x$ . Similar to supervised learning, there is a target value,  $x$ , to which no gradient is attached. We refer to  $\mathcal{L}_{\text{DFL}}$  as the supervised DFL loss. Unlike traditional methods

that rely solely on historical data and remain agnostic on the downstream control policy, our approach benefits of more informed sampling of the input space  $u$ . This allows the state-transition model to allocate its modeling resources more effectively. Consequently, input regions that are critical to the control policy are modeled with greater precision, enabling even simple models to achieve high accuracy.

Interestingly, the regret  $\mathcal{L}_r$  and the supervised DFL losses  $\mathcal{L}_{\text{DFL}}$  do not pursue the same objective. On the one hand, the regret  $\mathcal{L}_r$  trains the parameter  $\theta$  to bring the obtained objective value as close as possible to the objective value obtained with perfect knowledge of  $\theta$ . On the other hand, the supervised DFL loss aims at improving the ability of the state transition proxy model  $\hat{f}$  to match the true outcome. Notably, all performance losses would be the same if  $\hat{f}$  is an error-free approximation of  $f$ .

## 2.4 Robustness to infeasibility

The updates of the constraint parameters by the SGD may lead the optimization program to become infeasible. A small number of infeasible optimization problems over a full training epoch may not be an issue if the remaining solvable samples update the parameters in a way which restores feasibility. However, such an approach lacks robustness. Therefore, we relax the constraints bounding the output of the system state model and formulate them as a quadratic regularization term in the objective function. To that end, the state model output is associated with a target value  $X^g$ . The target value is the desired state for the system and is chosen by the control policy designer. For example, the target value for the indoor temperature might be set to 21°C (70°F). A weight  $w$  defines the importance of meeting the target value. In addition to preventing infeasibility, the quadratic regularization term can transform linear programs into strongly convex quadratic programs, which, unlike linear programs, provide a non-null (and thus informative) gradient:

$$\begin{aligned} \phi(x) &= \underset{u}{\operatorname{argmin}} f_o(x, u) + w * (x - X^g)^2 \\ \text{s.t. } f_i(u) &\leq 0 \quad i = 1, \dots, k \quad \forall t, \\ h_l(u) &= 0 \quad l = 1, \dots, l \quad \forall t, \\ \hat{x}_{t+1} &= \hat{f}(\hat{x}_t, u_t; \theta), \quad \forall t. \end{aligned} \quad (8)$$

In addition, we propose a pre-training stage where the dynamics model is trained on historical data. It provides the policy with an initial estimate of the dynamics model parameters, which improves feasibility in the early stages of training. However, minimizing a task-agnostic loss on historical data may lead the pre-training to converge to a local minimum of the supervised DFL loss. Escaping this local minimum might require using a large step size in the gradient descent algorithm, which could harm training convergence. To address this, we apply calibrated Gaussian noise to the pre-trained parameters. The intensity of the noise must be carefully selected to retain as much pre-training information as possible while helping the model escape the local minimum.

## 3 Building Thermodynamics Modeling and Day-ahead Control

We demonstrate the effectiveness of our method on the hourly day-ahead HVAC scheduling of a three-floor medium office building with 15 conditioned zones. We aim at learning the parameters of a multi-zone RC formulation modeling the building thermodynamics, while optimizing HVAC scheduling. In this section, we start by describing the building. The formulation of the optimization problem follows. Then, the data and their preprocessing are presented. Afterwards, we report and analyze the results of the proposed approach before comparing it to the conventional two-stage ITO performances. Finally, we assess the robustness of the model to distribution shift of the input data.

### 3.1 Building Description

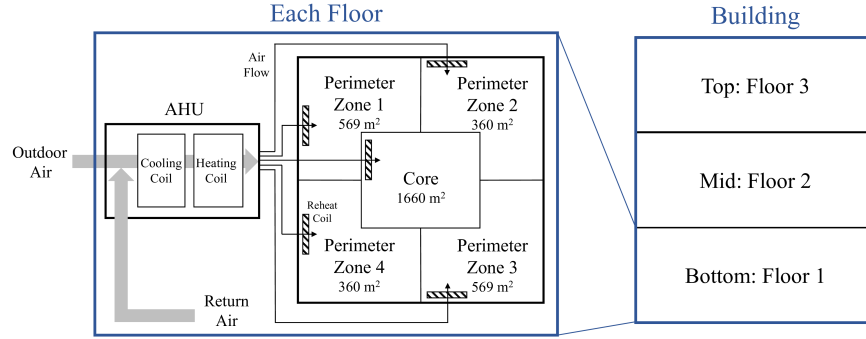
The building model is provided by the U.S. Department of Energy as part of EnergyPlus [10]. EnergyPlus is a state-of-the-art physics-based high-fidelity simulator for buildings. The building is an office building of 4928 m<sup>2</sup> (53,628 ft<sup>2</sup>) assumed to be located in Denver, CO. It comprises three floors of six zones each, only five being thermally controlled. The sixth zone is a plenum that contains the ducts transporting the conditioned air. Each floor is equipped with an Air Handling Unit (AHU). An AHU is an equipment that centrally conditions the air (either cooling or heating) before distribution to the zones via the duct network. The AHU is also responsible for the ventilation of the zones. The system in this building is a variable air volume system with reheat. The air flow rate reaching each zone can be modulated individually and, if necessary, the air can be heated just before entering the zone. Figure 2 provides a layout of the building floor and the AHU architecture. We make two assumptions with respect to the original template:

- (i) we replace the gas-fired heating coil of the AHU with an efficiency of 81 % by an electrical coil with an efficiency of 100 %;
- (ii) we divide the AHU electricity consumption proportionally to the air mass flow rate of each zone served by the AHU, while the reheat coil electricity consumption is associated with the zone it supplies.

The building features a conventional occupancy schedule for offices; most of the activity takes place during weekdays from 7am to 6pm.

### 3.2 Control Formulation

The control problem aims to find the temperature profile for each zone  $\tau_{t,z}^{\text{in}}$  that minimizes the electricity cost of HVAC while ensuring thermal comfort (9). The set of zones and time steps are  $\mathcal{Z}$  (with cardinality  $Z$ , index  $z$ ) and  $\mathcal{T}$  (with cardinality  $T$ , index  $t$ ), respectively. The electricity price consists of two terms. The first term (i) is the demand charge. The demand charge is proportional to the daily power consumption peak  $p^{\text{d}}$  and cost  $\lambda^{\text{d}}$ . The demand charge reflects the cost of network infrastructure to supply the requested power. The second term (ii) is the cost for the energy consumed over the time step  $p_t^1$  at price  $\lambda_t^1$  following the time-of-use tariff. Finally, the regularization term (iii) guarantees the problem feasibility, even during DFL. The regularization term penalizes the



**Figure 2: The building is made of three similar floors. Each floor is equipped with its own variable air volume air handling unit responsible for the thermal comfort of the five zones.**

deviation of the indoor temperature  $\tau_{t,z}^{\text{in}}$  from the target  $T_{t,z}^{\text{tgt}}$  with a quadratic term. The weight  $w_{t,z}$  adjusts the penalty incurred for thermal discomfort (i.e., the indoor temperature deviation from the target value) through the zones and time.

$$\min_{\tau^{\text{in}}} \underbrace{p^{\text{d}} \lambda^{\text{d}}}_{(i)} + \sum_{t=0}^T \left( \underbrace{p_t^{\text{i}} \lambda_t^{\text{i}}}_{(ii)} + \underbrace{\sum_{z=0}^Z w_{t,z} (\tau_{t,z}^{\text{in}} - T_{t,z}^{\text{tgt}})^2}_{(iii)} \right) \Delta t \quad (9)$$

The RC model in (10) accounts for multi-zonal building dynamics, where matrix  $\alpha \in \mathbb{R}^{Z \times Z}$  represents inter-zonal heat transfer, vectors  $R \in \mathbb{R}^Z$  and  $C \in \mathbb{R}^Z$  are zonal lumped resistances and capacitances, and vectors  $\eta^{\text{h}} \in \mathbb{R}^Z$  and  $\eta^{\text{c}} \in \mathbb{R}^Z$  are heating and cooling efficiencies, which include thermal losses of the duct system. For the 15-zone building, the RC model features 265 parameters. The parameters  $\theta = (\alpha, \eta^{\text{c}}, \eta^{\text{h}}, R, C)$  will be learned while scheduling the HVAC system.

$$\tau_{t+1}^{\text{in}} = \Delta t \left( \alpha \tau_t^{\text{in}} + \frac{(\eta^{\text{h}} p_t^{\text{h}} - \eta^{\text{c}} p_t^{\text{c}})}{C} + \frac{(\tau_t^{\text{amb}} - \tau_t^{\text{in}})}{RC} \right) \quad \forall t \quad (10)$$

The initial conditions are given by:

$$\tau_{0,z}^{\text{in}} = T_{0,z}^{\text{in}} \quad \forall z. \quad (11)$$

For each zone, the HVAC capacities for cooling  $\bar{P}_{t,z}^{\text{c}}$  and heating  $\bar{P}_{t,z}^{\text{h}}$  are determined based on historical data. The constraints are imposed at the zonal (12) and floor (13) levels. The floor level constraint is essential to capture the maximum consumption of each AHU while the zonal constraint limits the amount of energy that can be dedicated to a specific zone. The set of floor is  $\mathcal{F}$  with cardinality  $F$ , and index  $f$ . Each floor  $f$  is a set containing the zones  $z$  of that floor.

$$p_{t,z}^{\text{c}} \leq \bar{P}_{t,z}^{\text{c}}, p_{t,z}^{\text{h}} \leq \bar{P}_{t,z}^{\text{h}} \quad \forall t, z \quad (12)$$

$$\sum_{z \in f} p_{t,z}^{\text{c}} \leq \bar{P}_{t,f}^{\text{c}}, \sum_{z \in f} p_{t,z}^{\text{h}} \leq \bar{P}_{t,f}^{\text{h}} \quad \forall t, f \quad (13)$$

We define the HVAC electricity power as the sum of the cooling and heating electricity powers (14). Eq. (15) maintains the energy balance. In this reduced form, the HVAC electricity must be purchased from the grid.

$$p_{t,z}^{\text{hvac}} = p_{t,z}^{\text{h}} + p_{t,z}^{\text{c}} \quad \forall t, z \quad (14)$$

$$\sum_{z=0}^Z p_{t,z}^{\text{hvac}} = p_t^{\text{i}} \quad \forall t \quad (15)$$

Constraint (16) defines the power peak demand as being greater than all power demand. Because the peak demand directly increases the electricity cost, this constraint is equivalent to  $p^{\text{d}} = \max \{p_t^{\text{i}} : t \in T\}$ , but it avoids making the problem bilevel.

$$p^{\text{d}} \geq p_t^{\text{i}} \quad \forall t \quad (16)$$

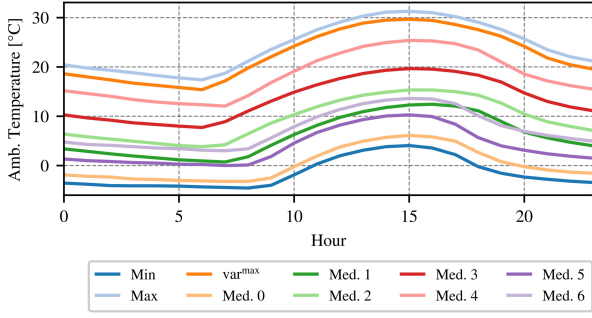
Finally, the power imported from the grid must be inferior to the line capacity  $\bar{P}^{\text{i}}$ . This is equivalent to imposing the peak power demand to be inferior to  $\bar{P}^{\text{i}}$  (17). All power levels must be positive (18).

$$p^{\text{d}} \leq \bar{P}^{\text{i}} \quad (17)$$

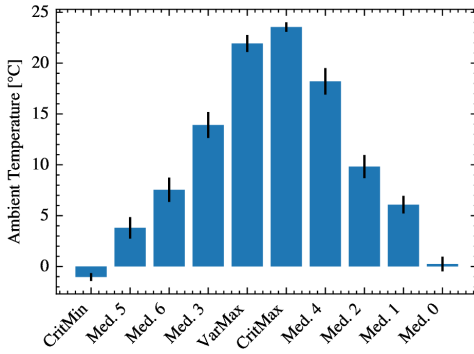
$$p_{t,z}^{\text{c}}, p_{t,z}^{\text{h}}, p_t^{\text{i}} \geq 0 \quad \forall t, z \quad (18)$$

### 3.3 Datasets and Clustering

To simulate building thermodynamics, two typical meteorological datasets are used. A typical meteorological dataset is a representative year of hourly weather data for a specific location, compiled from long-term observations. The first dataset is used as input to generate one year of historical data *without* any optimal scheduling. The second dataset provides weather scenarios for the scheduling optimization. Since optimizing and simulating the daily schedule for the 365 days of the year is burdensome, we perform a k-medoid clustering on the 365 days of weather data in the second dataset. The k-medoid algorithm is preferred to the k-mean algorithm to avoid the creation of fictitious average data. We focus on the ambient temperature because it is the only weather measurement necessary to the RC model (10). First, we identify three extreme days: the coldest, the hottest, and the day with the highest temperature variance. These three days are three fixed medoids of our medoid search. Then, we look for seven additional medoids to



**Figure 3: Ambient temperature profiles of the ten medoids. The three first labels are extreme scenarios.**



**Figure 4: Mean and standard deviation of the medoids. They are ordered to form a smooth cycle.**

obtain the best partition of the space. The resulting temperature profiles of the k-medoid algorithm are given in Figure 3. The cluster means and standard deviations are reported in Figure 4.

### 3.4 Hierarchical Performance Loss

As seen in section 2.3, we can define any supervised DFL loss for learning the parameters  $\theta = (\alpha, \eta^c, \eta^h, R, C)$  over the ten daily scenarios selected. In this specific case, the command is the indoor temperature profile for each zone  $\tau_{t,z}^{\text{in}}$  since it is the input of the thermostat. Therefore, the DFL loss focuses on the difference between the expected and the ex-post HVAC power consumptions. The underlying goal is to accurately capture the impact of HVAC power prediction errors on the objective value. Because the HVAC power  $p^{\text{hvac}}$  appears linearly in the objective function, we use a linear performance loss, the MAE. Furthermore, we weight the loss at each time step with the all-inclusive price coefficient  $\lambda_t^{\text{id}}$  that adds the peak demand price  $\lambda_t^{\text{d}}$  to the energy prices  $\lambda_t^{\text{i}}$  for the time step with the highest ex-post power. Otherwise,  $\lambda_t^{\text{id}}$  and  $\lambda_t^{\text{i}}$  are equal. Finally, we leverage hierarchical loss. By order of importance, we minimize the error for (i) the whole building to obtain an accurate expectation of the electricity bill; (ii) each AHU which means for the sum of the five zones at each floor  $f$ ; (iii) each zone  $z$ . The final loss function is given by (19). The number of zones at

each loss level is used to weight the importance. The building encompasses 15 zones ( $w_b = 15$ ) while each floor features five zones ( $w_f = 5 \forall f$ ).

$$L_{\text{DFL}} = \frac{1}{T} \sum_{t=0}^T \lambda_t^{\text{id}} \left( w_b \cdot \text{MAE} \left( \hat{p}_{t,b}^{\text{hvac}}, p_{t,b}^{\text{hvac}} \right) + w_f \sum_{f=0}^F \text{MAE} \left( \hat{p}_{t,f}^{\text{hvac}}, p_{t,f}^{\text{hvac}} \right) + \sum_{z=0}^Z \text{MAE} \left( \hat{p}_{t,z}^{\text{hvac}}, p_{t,z}^{\text{hvac}} \right) \right) \quad (19)$$

where  $\hat{p}_{t,f}^{\text{hvac}} = \sum_{z \in f} p_{t,z}^{\text{hvac}}$  is the HVAC power of floor  $f$  and  $\hat{p}_{t,b}^{\text{hvac}} = \sum_{z=0}^Z p_{t,z}^{\text{hvac}}$  is the HVAC power of the whole building.

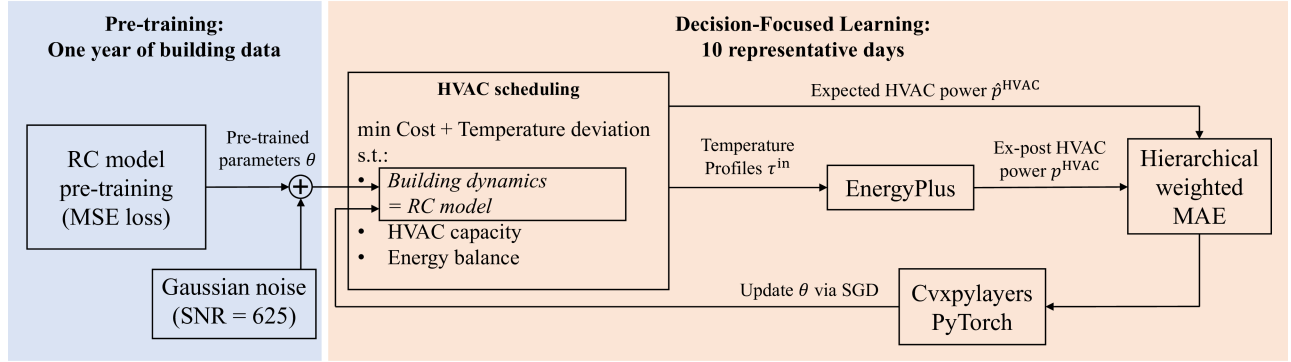
### 3.5 Results

Our goal is to learn the parameters  $\theta = (\alpha, \eta^c, \eta^h, R, C)$  that define the building thermodynamics (10). We initialize the algorithm by the pre-training. During this warm-start period, the parameters are trained over the 365 days of the first typical weather file using an MSE loss. Then, an element-wise Gaussian noise  $N \sim \mathcal{N}(0, \theta_i/25)$ , which corresponds to a Signal-to-Noise Ratio (SNR) of 625 (i.e.,  $25^2$ ), is applied independently on each element of the parameter vector. This operation aims at getting the RC model out of the pre-training local minimum while retaining a maximum of information.

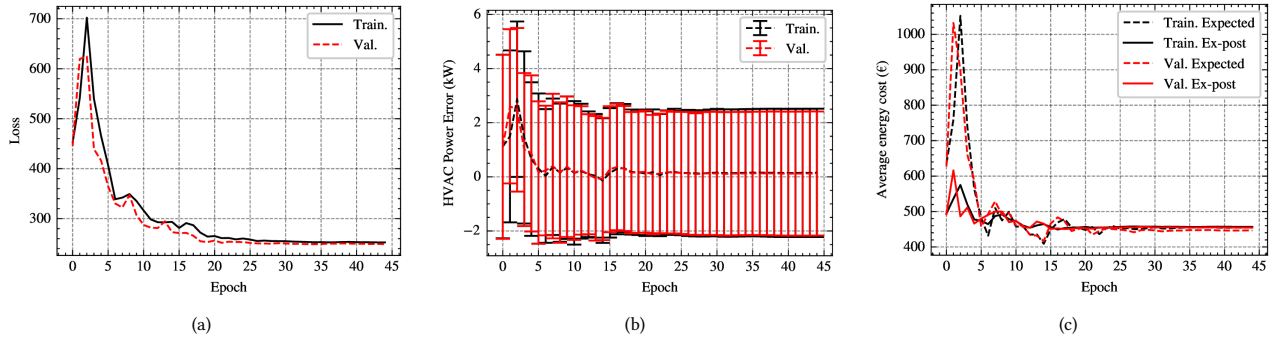
Afterwards, we enter the DFL stage. The DFL training set is made of ten representative days obtained via clustering. In addition, we sample randomly in each cluster two days to form a validation set and a test set of ten days. We consider solving and simulating the solution once for each of the ten training problems as a training epoch of ten samples. Each problem is solved using the ECOS solver since it was reported to be the best performing solver for conic DFL [34]. *Cvxpylayer* computes the gradient of the optimization problem and *PyTorch* performs the backpropagation. As soon as a sample (i.e., a day) has been solved and simulated, the parameters are updated (i.e., pure Stochastic Gradient Descent (SGD)). To facilitate learning, we order the days to form a smooth cycle when looping over the epochs as depicted by Figure 4. We initialize *Adam*, a gradient-based optimization algorithm, with a learning rate at 0.001 and a polynomial decay of  $\gamma = 0.9$  [21]. We set the maximum number of epochs to 50 with an early stopping featuring a 10-epoch patience. Figure 5 represents the whole process.

The time-of-use tariff is a rectangular function. From 7pm to 6am (excluded) the price is 0.3 €/kWh. It rises to 0.6 €/kWh from 6am to 7pm (excluded) as depicted by Figure 7. The increase highlights the higher demand during the day.

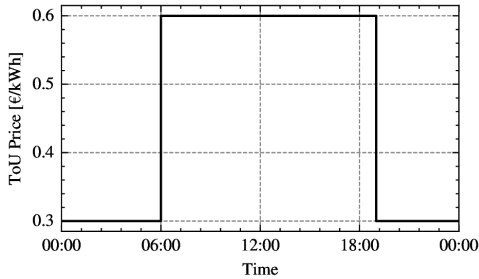
The performances of the DFL training are presented in Figure 6. Overall, the training and validation metrics follow very similar trends. This suggests that ten medoids are sufficient for the model to generalize effectively. The first plot displays the evolution of the loss over the training epochs. The training is interrupted after 44 epochs by the early stopping because the minimum value for the validation loss is reached at the 34<sup>th</sup> epoch. Both the training and



**Figure 5:** The building RC model is pre-trained on one year of historical data before a small noise is applied on the parameters. They are then updated to perform the best day-ahead HVAC scheduling over 10 representative days with respect to the task-aware MAE loss.



**Figure 6:** (a) The DFL supervised loss (hierarchical weighted MAE) converges at much lower levels than initially; (b) the error mean converges towards zero and the standard deviation reduces; (c) the expected cost becomes a more accurate prediction of the ex-post cost while diminishing.



**Figure 7:** Time-of-use tariff with an appreciation of electricity from 6am to 7pm.

validation losses tend to decrease until the 34<sup>th</sup> epoch before leveling off. At epoch 34, the training and validation losses stand at 252 and 248, respectively.

The center plot (figure 6.b) shows the evolution of the error, assuming a Gaussian distribution. After about five epochs, the error mean converges around 0 with a standard deviation at about 2.5 kW. Then, the error mean ripples slightly until the 25<sup>th</sup> epoch

where it stabilizes at 0.15 with a standard deviation of 2.35. The training and validation sets exhibit very similar trends.

The last plot displays the evolution of the expected and ex-post electricity costs for the training and validation sets. The expected and ex-post costs are very volatile over the first half of the training. Over the second half, the expected and ex-post costs level off at about 450 € and 455 €, respectively. The bill prediction of the day-ahead planning is therefore accurate. It is worth noting that the training converges at the minimum ex-post cost over all epochs. In conclusion, when surveying simultaneously the three metrics, the set of parameters obtained at the end of the 34<sup>th</sup> epoch offers the highest performance.

In comparison, the ITO model performs extremely poorly. The RC parameters of the ITO model were found by minimization of the MSE. Nevertheless, the MSE reported in Table 1 indicate clearly that the DFL model at the 34<sup>th</sup> epoch is better for each data set: training, validation, and test. The two models show very consistent metrics across the three data sets.

The hierarchical loss of the ITO test set stands at 652, more than two and a half times the value of the DFL model (253). However, the MAE is slightly lower for the ITO model, rising from 2.66 for



the ITO model to 2.94 for the DFL model. This can be explained by the hierarchical loss that does not aim at minimizing the prediction error at the zonal level. The test set average error of the ITO model stands at -1.81 kW with a standard deviation of 2.39 kW. Since, on average, the expected power is much lower than the ex-post power, the electricity cost is underestimated. The expected cost is 85 € for an ex-post cost at 474 €. It represents a prediction error of 389 €. In contrast, the DFL model achieve a prediction error of 16 €. Very interestingly, the ex-post cost of the DFL model (468 €) is lower than for the ITO model (474 €). The foremost hindrance to DFL is the heavy computational burden associated with training. On a personal MacBook Pro, equipped with an Apple M1 Pro chip and 32 GB of memory, the DFL training took 7 hours compared to 16 minutes for the conventional supervised learning.

	ITO			DFL		
	Train.	Val.	Test	Train.	Val.	Test
Hierarchical loss	655	646	652	252	248	253
MAE (kW)	2.68	2.62	2.66	2.93	2.87	2.94
MSE (kW <sup>2</sup> )	16.8	16.6	16.7	15.4	14.6	15.6
Error mean (kW)	-1.85	-1.86	-1.81	0.15	0.12	0.08
Error std (kW)	2.34	2.28	2.39	2.35	2.27	2.38
Expected cost (€)	82	75	85	454	446	452
Ex-post cost (€)	478	474	474	456	455	468
Cost error (€)	395	398	389	2	8	16
Training Time	16 min			7h02min		

**Table 1: Performances of the Identify-Then-Optimize (ITO) model versus the Decision-Focused Learning (DFL) model at epoch 34.**

The scheduling of the building lower floor on the coldest day is shown in Figure 8. The upper part displays the temperature profiles within the five zones. The blue shades correspond to a given penalty per hour for deviating from the 21°C target. The discomfort penalty is the lowest during the night, from midnight to 7am, and the highest during the conventional working hours, reflecting the higher occupancy.

Overall, and aligned with the discomfort penalty, the temperature profiles stay close to the target of 21°C during the day, with greater deviations occurring in the evening and night. Still the deviations remain within admissible ranges. The electric power scheduling of the HVAC is shown in the lower two-thirds figure. The zonal electric power scheduling in subfigure 8.b is quite inaccurate. This can be explained by the hierarchical loss, which focuses primarily on having a precise day-ahead scheduling of the whole building. This is depicted by subfigure 8.c where the power trends are much better for the floor level and the entire building. At the building level, the attempt to shift the HVAC consumption towards cheap electricity stands out. Indeed, the HVAC consumption is high during the night, from 2am to 8am, and peaks again towards the end of the day, when the price of electricity goes down. The RC model is a linear approximation of the dynamics considering only the previous time step. Moreover, the outdoor temperature is the only weather input. Consequently, the limited modeling capability and input of

the RC model, along with the complexity of the building’s thermodynamics, are responsible for the remaining inaccuracies. The real-time adjustment to the model inaccuracies and the changing conditions should be handled by a real-time controller.

### 3.6 Ambient Temperature Shift

We investigate the robustness of the proposed framework to the input distribution shift. In particular, we assess the model response to a hot year. We build the hot year by selecting the hottest sample from each cluster. For the sample associated with the hottest cluster, we add two Celsius degrees on top of the hottest day, assuming record-breaking temperatures. Figure 9 compares the ambient temperature distribution of the four data sets. The average and maximum ambient temperatures in the hot year set are higher than in any other set. The average is 2°C higher for the hot-year set than for the training set, thus reflecting a distribution shift.

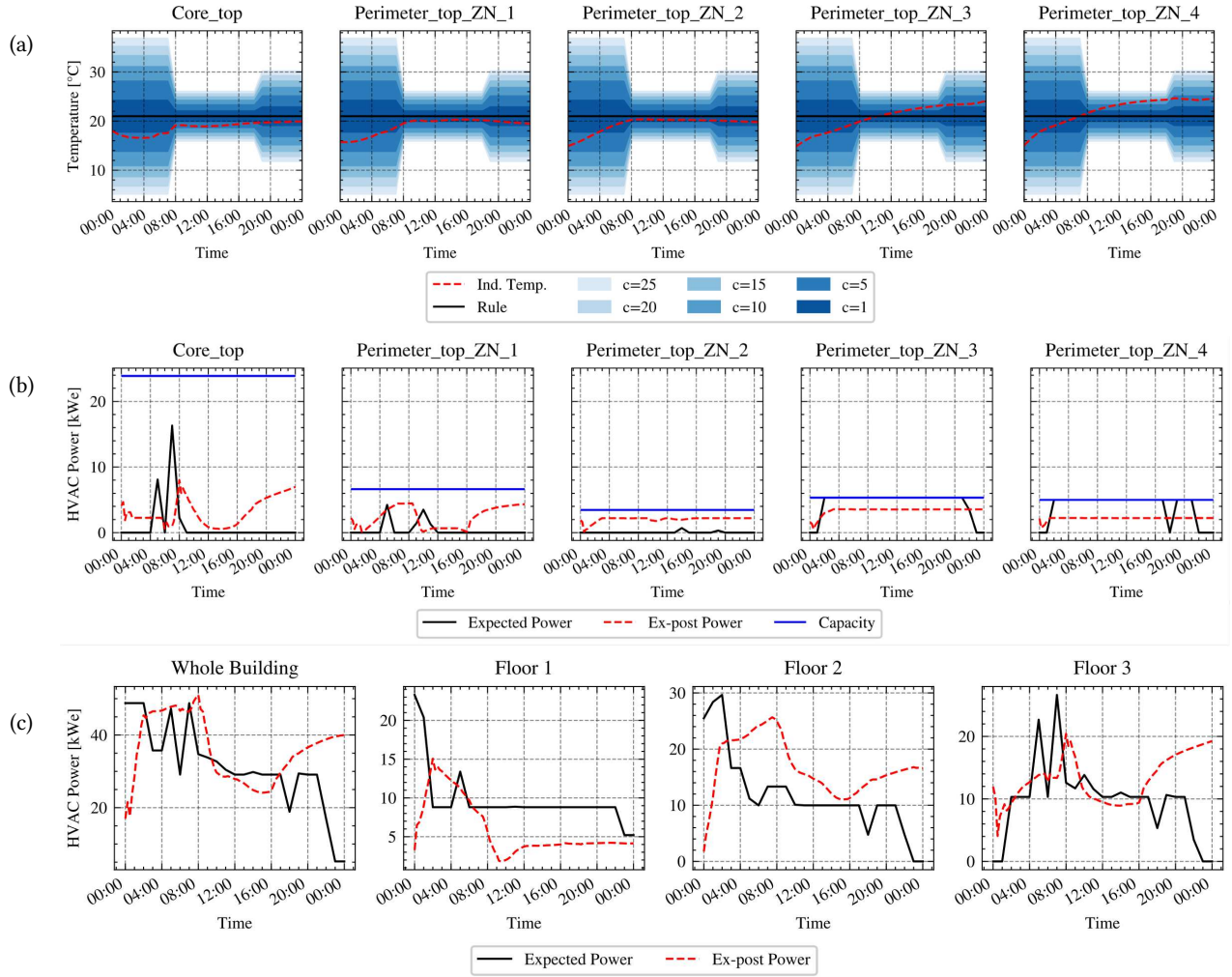
Table 2 shows the metrics of the ITO and DFL models over the hot year data set. The results show a slight performance degradation for both models compared to the randomly sampled test set. Performance degradation is expected because the hot year data distribution is purposely different from the previous distribution. The hierarchical loss stands at 685 (+5.1% compared to the test set) and 271 (+7.1%) for the ITO and DFL models, respectively. The error standard deviation remains virtually unchanged. The mean, however, worsens for the ITO model but improves for the DFL model. This results in a greater cost error for the ITO but not for the DFL. Overall, DFL outperforms ITO and is less affected by input distribution shift than the ITO model.

	ITO	DFL
Hierarchical loss	685	271
MAE (kW)	2.76	2.95
MSE (kW <sup>2</sup> )	18.3	15.8
Error mean (kW)	-1.96	0.1
Error std (kW)	2.39	2.39
Expected cost (€)	79	467
Ex-post cost (€)	500	482
Cost error (€)	420	15

**Table 2: Metrics of the Identify-Then-Optimize (ITO) model and the Decision-Focused Learning (DFL) model over a dataset representing a distribution shift towards higher ambient temperature (i.e., hot year).**

## 4 Discussion

The proposed DFL strategy outperforms the ITO approach. The naive ITO framework that consists in MSE minimization over historical data leads to extremely poor dynamic model. In particular, the RC model is unable to appropriately fit the operating areas relevant for the control policy. This leads to potentially severe electric power scheduling underestimation resulting in unexpectedly high ex-post bills. In the case study, the ex-post cost was on average 6 times higher than the expected cost for the ITO model. This raises important questions about the use of RC model for the



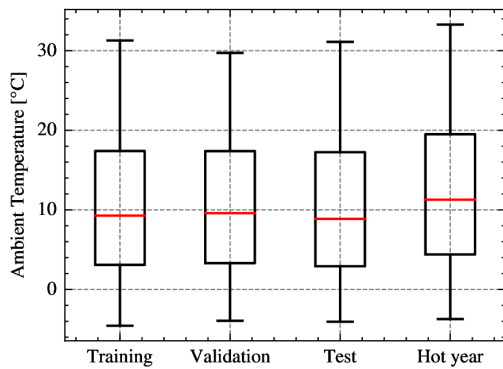
**Figure 8: Scheduling for the coldest day obtained by the DFL RC model at epoch 34. Subfigures are (a) the zonal indoor temperature of the top floor; (b) the zonal HVAC electric power of the top floor, and (c) the aggregated HVAC power for the whole building and each floor. The blue shades indicate the temperature deviations causing an objective value penalty lower than a given threshold. Temperature control is more loose during the night and evening because of the worker absence. The mismatch between the expected and ex-post HVAC powers in (b) can be explained by the intrinsic RC model limitations and the task-aware loss which gives little importance to accurate zonal predictions.**

day-ahead HVAC scheduling, at least within an ITO framework. In contrast, the proposed DFL approach significantly improves both the accuracy of the bill prediction and the cost reduction. All error metrics show improvement, except for the MAE. However, the DFL approach is much more computationally demanding. The DFL duration before triggering the early stopping is about 7 hours compared to 16 minutes for the ITO approach. This is due to the DFL need to solve, for each sample, the day-ahead HVAC scheduling optimization problem and simulating the results. Furthermore, software for DFL are not as mature as the ones for supervised learning. Future software development for DFL will likely harvest the power of parallel computing and new hardware leading to significant time reduction. In this specific case, due to the linear formulation of the

optimization, the computational burden mainly lies with the simulation. We recommend considering using a reduced-order physics model of the simulator in the initial training steps. Finally, we showed that the model is robust to an input data shift towards higher temperature.

## 5 Conclusion

This paper presents a new framework that enables simultaneous identification and control (or planning) of a complex system. This is of particular interest for black-box systems in which dynamics change radically once the control policy is applied, thus making the historical observation little informative. The method exploits recent advances in decision-focused learning, and more specifically,



**Figure 9: Ambient temperature distribution of each data set. Whiskers represent minimum and maximum temperatures.**

in the calculation of the gradient of cone programs. The gradient of the convex optimization control policy output (i.e., the command) with respect to the policy parameters (e.g., the system dynamics model parameters) can be computed. Because a key constraint of the control policy is being learned, the feasibility domain evolves at each gradient descent step. We handle the potential infeasibility by relaxing the constraints on the system states and pre-learning the uncertain parameters on historical data. Furthermore, we propose a new type of loss function that bypasses the need for the system to be differentiable. Not only is this necessary for black-box systems, but also for actual physical systems that cannot be easily modeled.

We apply the proposed approach to the day-ahead HVAC scheduling of a 15-zone office located in Denver, CO. This case study showcases the efficient learning provided by our framework, and the ability to design complex task-aware loss functions to reflect the impact of the parameter misestimation on the objective value. The results show unambiguously the added value of simultaneous system identification and planning by displaying a lower ex-post cost along with a more accurate price forecast and a reduced error on the day-ahead HVAC scheduling.

Future work might consider building a surrogate model to bypass the simulation in the early DFL steps. More detailed models for building dynamics should also be investigated. Such models will likely turn the optimization formulation into a mixed integer linear program that will bring new challenges in the computation of the gradient. In addition, the appropriate number of representative scenarios necessary for DFL training should be explored. A trade-off must be found between a good generalization of the DFL model to the whole initial dataset and the substantial computational time that can be saved. Lastly, future work should quantify the impact of a change in building use that results in new internal heat load, occupancy pattern, or temperature setpoints, and thus, HVAC operating points.

## 6 Acknowledgements

Pietro Favaro is an FNRS-F.R.S. Research Fellow (grant number FC 49537) and Research Fellow of the Belgian American Educational Foundation (B.A.E.F.).

## References

- [1] Abdul Afram and Farrokh Janabi-Sharifi. 2014. Review of modeling methods for HVAC systems. *Applied Thermal Engineering* 67, 1 (June 2014), 507–519. <https://doi.org/10.1016/j.applthermaleng.2014.03.055>
- [2] Zakia Afroz, GM Shafiullah, Tania Urmee, and Gary Higgins. 2018. Modeling techniques used in building HVAC control systems: A review. *Renewable and Sustainable Energy Reviews* 83 (March 2018), 64–84. <https://doi.org/10.1016/j.rser.2017.10.044>
- [3] Akshay Agrawal, Shane Barratt, Stephen Boyd, Enzo Busseti, and Walaa M. Moursi. 2019. Differentiating Through a Cone Program. *Journal of Applied and Numerical Optimization* 1 (2019), 107–115. <https://jano.biemdas.com/archives/931>
- [4] Akshay Agrawal, Shane Barratt, Stephen Boyd, and Bartolomeo Stellato. 2020. Learning Convex Optimization Control Policies. In *Proceedings of the 2nd Conference on Learning for Dynamics and Control*. PMLR, Berkeley, CA, 361–373. <https://proceedings.mlr.press/v120/agrawal20a.html> ISSN: 2640-3498.
- [5] Ercan Atam and Lieve Helsen. 2015. A convex approach to a class of non-convex building HVAC control problems: Illustration by two case studies. *Energy and Buildings* 93 (April 2015), 269–281. <https://doi.org/10.1016/j.enbuild.2015.02.026>
- [6] Mokhtar S. Bazaraa, John J. Jarvis, and Hanif D. Sherali. 2009. *The Simplex Method*. John Wiley & Sons, Ltd, Hoboken, NJ, US, Chapter 3, 91–149. <https://doi.org/10.1002/9780471703778.ch3> arXiv:<https://onlinelibrary.wiley.com/doi/pdf/10.1002/9780471703778.ch3>
- [7] Filip Belić, Željko Hocenski, and Dražen Slišković. 2016. Thermal modeling of buildings with RC method and parameter estimation. In *2016 International Conference on Smart Systems and Technologies (SST)*. IEEE, Osijek, Croatia, 19–25. <https://doi.org/10.1109/SST.2016.7765626>
- [8] Enzo Busseti, Walaa M. Moursi, and Stephen Boyd. 2019. Solution refinement at regular points of conic problems. *Computational Optimization and Applications* 74, 3 (Dec. 2019), 627–643. <https://doi.org/10.1007/s10589-019-00122-9>
- [9] Marie Chau and Michael C. Fu. 2015. *An Overview of Stochastic Approximation*. Springer, New York, NY, 149–178 pages. [https://doi.org/10.1007/978-1-4939-1384-8\\_6](https://doi.org/10.1007/978-1-4939-1384-8_6)
- [10] Drury Crawley, Linda Lawrie, Frederick Winkelmann, W.F. Buhl, Y. Joe Huang, Curtis Pedersen, Richard Strand, Richard Liesen, Daniel Fisher, Michael Witte, and Jason Glazer. 2001. EnergyPlus: Creating a New-Generation Building Energy Simulation Program. *Energy and Buildings* 33 (April 2001), 319–331. [https://doi.org/10.1016/S0378-7788\(00\)00114-6](https://doi.org/10.1016/S0378-7788(00)00114-6)
- [11] Priya L. Donti, Brandon Amos, and J. Zico Kolter. 2017. Task-based end-to-end model learning in stochastic optimization. In *Proceedings of the 31st International Conference on Neural Information Processing Systems (NIPS'17)*. Curran Associates Inc., Red Hook, NY, USA, 5490–5500.
- [12] Ján Dragoňa, Aaron R. Tuor, Vikas Chandan, and Draguna L. Vrabie. 2021. Physics-constrained deep learning of multi-zone building thermal dynamics. *Energy and Buildings* 243 (July 2021), 110992. <https://doi.org/10.1016/j.enbuild.2021.110992>
- [13] Diana D’Agostino, Roberto Landolfi, Maurizio Nicoletta, and Francesco Minichiello. 2022. Experimental Study on the Performance Decay of Thermal Insulation and Related Influence on Heating Energy Consumption in Buildings. *Sustainability* 14, 5 (Jan. 2022), 2947. <https://doi.org/10.3390/su14052947> Number: 5 Publisher: Multidisciplinary Digital Publishing Institute.
- [14] Adam N. Elmachtoub and Paul Grigas. 2022. Smart “Predict, then Optimize”. *Manage. Sci.* 68, 1 (Jan. 2022), 9–26. <https://doi.org/10.1287/mnsc.2020.3922>
- [15] M. González-Torres, L. Pérez-Lombard, Juan F. Coronel, Ismael R. Maestre, and Da Yan. 2022. A review on buildings energy information: Trends, end-uses, fuels and drivers. *Energy Reports* 8 (Nov. 2022), 626–637. <https://doi.org/10.1016/j.egy.2021.11.280>
- [16] Raad Z. Homod. 2013. Review on the HVAC System Modeling Types and the Shortcomings of Their Application. *Journal of Energy* 2013, 1 (2013), 768632. <https://doi.org/10.1155/2013/768632> \_eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1155/2013/768632>.
- [17] Hao Huang, Lei Chen, and Eric Hu. 2014. Model predictive control for energy-efficient buildings: An airport terminal building study. In *11th IEEE International Conference on Control & Automation (ICCA)*. IEEE, Piscataway, NJ, USA, 1025–1030. <https://doi.org/10.1109/ICCA.2014.6871061> ISSN: 1948-3457.
- [18] Doseok Jang, Larry Yan, Lucas Spangher, and Costas J. Spanos. 2024. Active Reinforcement Learning for Robust Building Control. *Proceedings of the AAAI Conference on Artificial Intelligence* 38, 20 (Mar. 2024), 22150–22158. <https://doi.org/10.1609/aaai.v38i20.30219>
- [19] Scott Jeon, Alessandro Abate, and Jonathan M. Cullen. 2023. Low emission building control with zero-shot reinforcement learning. In *Proceedings of the Thirty-Seventh AAAI Conference on Artificial Intelligence and Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence and Thirteenth Symposium on Educational Advances in Artificial Intelligence (AAAI’23/IAAI’23/EAAI’23, Vol. 37)*. AAAI Press, Washington, DC, 14259–14267. <https://doi.org/10.1609/aaai.v37i12.26668>

- [20] Yi-hao Kao, Benjamin Roy, and Xiang Yan. 2009. Directed Regression. In *Advances in Neural Information Processing Systems*, Y. Bengio, D. Schuurmans, J. Lafferty, C. Williams, and A. Culotta (Eds.), Vol. 22. Curran Associates, Inc., Vancouver, B.C., Canada. [https://proceedings.neurips.cc/paper\\_files/paper/2009/file/0c74b7f78409a4022a2c4c5a54401e1e1-1-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2009/file/0c74b7f78409a4022a2c4c5a54401e1e1-1-Paper-Conference.pdf)
- [21] Diederik P. Kingma and Jimmy Ba. 2017. Adam: A Method for Stochastic Optimization. <https://doi.org/10.48550/arXiv.1412.6980> arXiv:1412.6980 [cs].
- [22] Kevin J. Kircher and K. Max Zhang. 2015. On the lumped capacitance approximation accuracy in RC network building models. *Energy and Buildings* 108 (Dec. 2015), 454–462. <https://doi.org/10.1016/j.enbuild.2015.09.053>
- [23] Jayanta Mandi, James Kotary, Senne Berden, Maxime Mulamba, Victor Bucarey, Tias Guns, and Ferdinando Fioretto. 2023. Decision-Focused Learning: Foundations, State of the Art, Benchmark and Future Opportunities. <https://doi.org/10.48550/arXiv.2307.13565> arXiv:2307.13565 [[cs, math]]
- [24] Arkadi Nemirovski. 2007. Advances in Convex Optimization: Conic Programming. In *International Congress of Mathematicians*, Vol. 1. EMS Press, Zürich, Switzerland, 413–444.
- [25] Lukas Ortmann, Fabian Böhm, Florian Klein-Helmkamp, Andreas Ullbig, Save-rio Bolognani, and Florian Dörfler. 2024. Tuning and Testing an Online Feedback Optimization Controller to Provide Curative Distribution Grid Flexibility. <http://arxiv.org/abs/2403.01782> arXiv:2403.01782 [cs, eess].
- [26] Brendan O'Donoghue, Eric Chu, Neal Parikh, and Stephen Boyd. 2016. Conic Optimization via Operator Splitting and Homogeneous Self-Dual Embedding. *Journal of Optimization Theory and Applications* 169, 3 (June 2016), 1042–1068. <https://doi.org/10.1007/s10957-016-0892-3>
- [27] Lolla Phani Raghav, Rangu Seshu Kumar, Dhenuvakonda Koteswara Raju, and Arvind R. Singh. 2022. Optimal day ahead energy consumption management in grid-connected microgrids. *International Journal of Energy Research* 46, 2 (Feb. 2022), 1864–1881. <https://doi.org/10.1002/er.7303> Publisher: John Wiley & Sons, Ltd.
- [28] S. M. Hosseini, R. Carli, and M. Dotoli. 2019. Robust Day-Ahead Energy Scheduling of a Smart Residential User Under Uncertainty. In *2019 18th European Control Conference (ECC)*. IEEE, Naples, Italy, 935–940. <https://doi.org/10.23919/ECC.2019.8796182> Journal Abbreviation: 2019 18th European Control Conference (ECC).
- [29] Sanket Shah, Kai Wang, Bryan Wilder, Andrew Perrault, and Milind Tambe. 2022. Decision-Focused Learning without Decision-Making: Learning Locally Optimized Decision Losses. In *Advances in Neural Information Processing Systems*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (Eds.), Vol. 35. Curran Associates, Inc., New Orleans, LA, 1320–1332. [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/0904c7edde20d7134a77fc7f9cd86ea2-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/0904c7edde20d7134a77fc7f9cd86ea2-Paper-Conference.pdf)
- [30] Susannah Shoemaker. 2023. NREL Researchers Reveal How Buildings Across United States Do—and Could—Use Energy. <https://www.nrel.gov/news/features/2023/nrel-researchers-reveal-how-buildings-across-the-united-states-do-and-could-use-energy.html>. Accessed: 2024-05-25.
- [31] Mattia Silvestri, Senne Berden, Jayanta Mandi, Ali Irfan Mahmutogullari, Brandon Amos, Tias Guns, and Michele Lombardi. 2024. Score Function Gradient Estimation to Widen the Applicability of Decision-Focused Learning. <http://arxiv.org/abs/2307.05213> arXiv:2307.05213 [cs].
- [32] Jianwen Sun, Yan Zheng, Jianye Hao, Zhaopeng Meng, and Yang Liu. 2020. Continuous Multiagent Control Using Collective Behavior Entropy for Large-Scale Home Energy Management. *Proceedings of the AAAI Conference on Artificial Intelligence* 34, 01 (April 2020), 922–929. <https://doi.org/10.1609/aaai.v34i01.5439> Section: AAAI Technical Track: Applications.
- [33] U.S. Department of Energy. 2019. [https://www.energycodes.gov/sites/default/files/2023-10/ASHRAE901\\_OfficeMedium\\_STD2019.zip](https://www.energycodes.gov/sites/default/files/2023-10/ASHRAE901_OfficeMedium_STD2019.zip). Accessed: 2024-06-24.
- [34] Dariush Wahdany, Carlo Schmitt, and Jochen L. Cremer. 2023. More than accuracy: end-to-end wind power forecasting that optimises the energy system. *Electric Power Systems Research* 221 (Aug. 2023), 109384. <https://doi.org/10.1016/j.epsr.2023.109384>
- [35] Zhenbo Wang. 2024. A survey on convex optimization for guidance and control of vehicular systems. *Annual Reviews in Control* 57 (Jan. 2024), 100957. <https://doi.org/10.1016/j.arcontrol.2024.100957>
- [36] Bryan Wilder, Bistra Dilikina, and Milind Tambe. 2019. Melding the Data-Decisions Pipeline: Decision-Focused Learning for Combinatorial Optimization. *Proceedings of the AAAI Conference on Artificial Intelligence* 33, 01 (July 2019), 1658–1665. <https://doi.org/10.1609/aaai.v33i01.33011658> Section: AAAI Technical Track: Constraint Satisfaction and Optimization.
- [37] Yinyu Ye, Michael J. Todd, and Shinji Mizuno. 1994. An  $O(\sqrt{nL})$ -Iteration Homogeneous and Self-Dual Linear Programming Algorithm. *Mathematics of Operations Research* 19, 1 (Feb. 1994), 53–67. <https://doi.org/10.1287/moor.19.1.53> Publisher: INFORMS.
- [38] Liang Yu, Shuqi Qin, Meng Zhang, Chao Shen, Tao Jiang, and Xiaohong Guan. 2021. A Review of Deep Reinforcement Learning for Smart Building Energy Management. *IEEE Internet of Things Journal* 8, 15 (Aug. 2021), 12046–12063. <https://doi.org/10.1109/JIOT.2021.3078462> Conference Name: IEEE Internet of Things Journal.
- [39] Yaohui Zeng, Zijun Zhang, and Andrew Kusiak. 2015. Predictive modeling and optimization of a multi-zone HVAC system with data mining and firefly algorithms. *Energy* 86 (June 2015), 393–402. <https://doi.org/10.1016/j.energy.2015.04.045>
- [40] S. E. Zou, A. A. Shah, P. K. Leung, X. Zhu, and Q. Liao. 2023. A comprehensive review of the applications of machine learning for HVAC. *DeCarbon* 2 (Sept. 2023), 100023. <https://doi.org/10.1016/j.decarb.2023.100023>