

Data requirements and performance evaluation of model predictive control in buildings: A modeling perspective

Sicheng Zhan, Adrian Chong*

Department of Building, School of Design and Environment, National University of Singapore, 4 Architecture Drive, Singapore, 117566, Singapore



ARTICLE INFO

Keywords:

Model predictive control
Control-oriented model
Data requirements
Level of detail
Performance evaluation
Model identification

ABSTRACT

Model predictive control (MPC) has shown great potential in improving building performance and saving energy. However, after over 20 years of research, it is yet to be adopted by the industry. The difficulty of obtaining a sufficient control-oriented model is one major factor that hinders the application. In particular, what data is required to build the model and what control performance can be expected with a certain model remain unclear. This study attempts to uncover the underlying reasons and guide future research to tackle the challenges. It starts by clarifying a finer categorization of past studies with respect to both modeling methods and modeling purposes. An extended Level of Detail (LoD) framework is proposed to quantify the data usage in each study. Accordingly, meta-analyses are conducted to compare the data requirements of different modeling categories. The criteria and approaches for model performance evaluation are summarized and classified into validation and verification methods, followed by a discussion about the relationship between the model and control performance. The critical review provides new perspectives on the data requirements and performance evaluation of control-oriented models. Ultimately, the paper concludes with five directions for future research to bridge the gaps between data requirements, model performance, and control performance.

1. Introduction

1.1. MPC and control-oriented models

Buildings take up 30–40% of global greenhouse gas emission and energy consumption [1], among which up to 85% is consumed in the operation phase [2]. Building system control is a challenging task because of the varying system dynamics and disturbances. At present, PID control are mostly used in practice, yielding the unsatisfactory performance if not well-tuned and the absence of multi-objective supervisory control [3]. These suggest the great energy-saving potential of implementing advanced optimal control schemes.

Model Predictive Control (MPC) was first applied for industrial process control [4] and has been tested in buildings since the 1990s [5, 6]. It is capable of adapting different system dynamics and disturbances, improving the thermal comfort conditions and energy performance simultaneously. The benefits are more prominent when the control task goes beyond setpoint tracking, such as occupancy-based control [7] or Demand Response (DR) applications [8]. However, not many actual implementation cases are spotted over the years, which can be

attributed to the relatively high requirements on modeling, expertise, data, hardware, usability, and costs [9,10]. Reducing the modeling effort and enhancing the model reliability are still essential problems to tackle.

Fig. 1 displays the typical framework of MPC. Three main processes are involved in obtaining the control decision: disturbance forecast, control-oriented model, and optimization. While all the three processes are essential, the control-oriented model has been acknowledged as the cornerstone of MPC [11,12]. Disturbance forecast provides the boundary condition for the control-oriented model over the prediction horizon, such as ambient conditions [13], occupant presence [14], and energy prices [15]. Under the boundary condition, the control-oriented model predicts the building's thermal response and energy performance with different control decisions. Based on the model, optimization is applied to identify the optimal control decision. According to the modeling purposes and methods, building metadata and/or time series training data may be needed for model identification. The optimization problem is defined by the objective function and the constraints, which are recently categorized in Ref. [16]. Different optimization algorithms may be selected depending on the problem formulation and the model

* Corresponding author.

E-mail address: adrian.chong@nus.edu.sg (A. Chong).

derivability [17]. To facilitate the optimization, a desirable model is expected to have a simplified structure and high accuracy, requiring less calibration and computational cost and maintaining certain physical significance [18].

1.2. Past reviews and research gaps

To procure a satisfactory model is one of the main barriers to implementing MPC in buildings. As reported in Henze [19], building and calibrating the models account for 70% of the total effort. In fact, it is not just modeling that is hard, but assessing the difficulty in advance as well [11]. Due to the importance and difficulty of this process, extensive research has been trying to tackle the challenge. Many review studies in the past few years have discussed relevant issues. Fundamentally, the modeling methods are usually categorized into physical-based (white-box), data-driven (black-box) and hybrid (gray-box) models [9,11,20]. In addition to that, Li & Wen [12] covered the mechanism of different building thermal response models, as well as models of energy storage and generation systems. Afram & Janabi-Sharifi [21] summarized the specific modeling techniques used in each of the three categories and introduced the general process from model creation to evaluation. Mirakhori & Dong [7] outlined the optimization techniques used corresponding to different modeling methods. Hilliard et al. [17] categorized the modeling methods by different spatial scales and listed the inputs and outputs of representative studies. Atam & Helsen [22] reviewed and compared different modeling methods, specifically for Ground Source Heat Pumps (GSHP). Afram et al. [23] reviewed the data collection and handling issues specifically for the Artificial Neural Networks (ANN) models. Rockett & Hathaway [24] talked about the effect of model update, occupancy uncertainty, and data handling on the model and control performance from a practical point of view. Afroz et al. [20] defined the physical processes in the physical models of different sub-systems, classified the sub-categories of black-box models (also covered in Ref. [8]), and elaborated the pros and cons of the three fundamental categories. The comparison is conducted with respect to prediction accuracy, generalization capability, training data requirement, and complexity. Serale et al. [25] separated the white-box models into detailed simulation models and reduced-order models, and also differentiated the models as building, HVAC systems, and building with HVAC systems. Fontenot & Dong [16] specified the main challenges in modeling as the high complexity of thermal models, and the uncertainties in disturbances. Pallonetto et al. [26] also distinguished the detailed and simplified white box, and thoroughly discussed the features and calibration issues of detailed simulation models.

There are three major gaps in the existing review studies. First, despite the model's well-known crucial role, a holistic review on model performance evaluation and its relationship with the control performance is missing. Model performance refers to the authenticity of control-oriented models, which is usually evaluated by the prediction

accuracy. In the meantime, it is also deemed to be necessary that the models represent the building dynamics for better extrapolation capabilities [27]. Control performance is reflected in the control results, such as energy consumption, thermal comfort, and the like. Model and control performance were shown closely related [28], but no quantitative relationship has been established.

Besides, the comparison of modeling methods is conducted mainly across the three fundamental categories. However, variations exist within the same category in terms of, for instance, data usage and prediction performance. For example, while black-box models are generally considered to be more accurate, different modeling techniques could result in up to 100% difference in prediction accuracy [29]. Also, gray-box models are declared to require less building metadata than white-box models and less training data than black-box models. Yet, depending on different modeling purposes and model fidelity, gray-box models can use extensive metadata [30] or training data [31].

Moreover, current discussions among different modeling methods are conceptual and qualitative. The lack of quantitative investigation leaves the modeling challenges unresolved. Hence, this study aims to promote the application of MPC in buildings by shedding some light on the modeling-related issues, specifically, two research questions:

- **What is the necessary and desirable data to build a satisfactory control-oriented model?** With the increasing deployment of building information modeling (BIM) and building management systems (BMS), much more data is being generated over the building life cycle, exerting a big challenge on data management and utilization [32]. For MPC, among the numerous potentially useful data, what is really needed to build a satisfactory model remains an open question. Therefore, a framework to quantify the data usage in different studies is needed to enable future studies and improve the generalization capability.
- **What is the minimum performance requirement of a control-oriented model?** Root Mean Square Error (RMSE) is a typical metric that quantifies models' prediction error, which has been shown insufficient to inform the control performance [33]. Meanwhile, there are other approaches and metrics for model evaluation. On the other hand, only few studies have focused on how the model performance would affect the control performance [33,34]. Thus, future research towards this direction could benefit from a systematic review and a deeper understanding of the model and controller performance evaluation.

1.3. Scope and structure

This paper gives a critical review on studies related to the control-oriented models used for MPC in buildings, mainly from the perspective of data requirements and performance evaluation. Section 2 presents a new model categorization regarding modeling purposes and

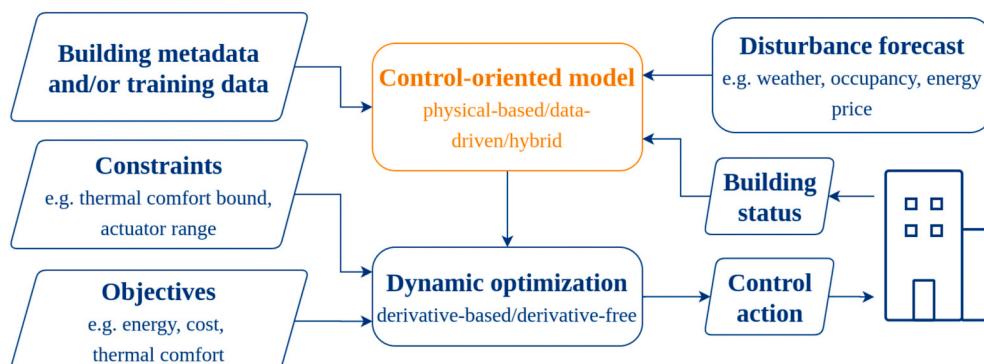
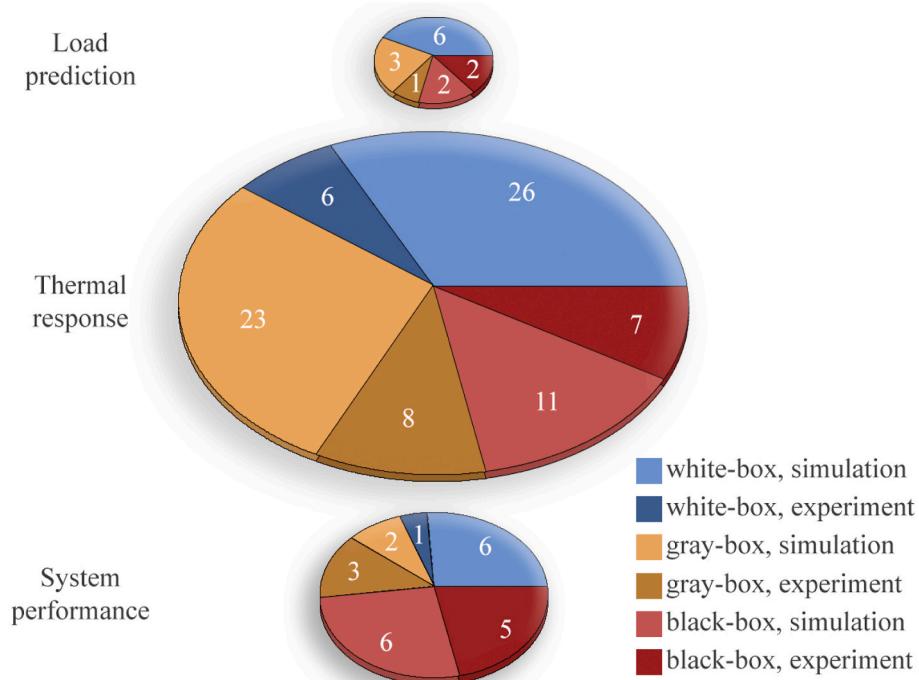


Fig. 1. Typical MPC framework, control-oriented model as the cornerstone.

Table 3

Categorization of 113 papers according to model methods and model purposes.

Modeling purposes	Modeling methods		
	White-box	Gray-box	Black-box
Load prediction	[53,87,91–93,104]	[63,82,88], [55] ^a	[86,89], [90,105] ^a
Thermal response	, [13–15,34,38–40,42,44,46,47,49,50,52,96, 106–116], [94,117–121] ^a	, , , , , [27,28,30,33,43,56–58,64,65,67,83,85,122–131], [10,59,61,62,95,97,132,133] ^a	[68–71,77,79,81,84,134–138], [72,74, 78,80,139–141] ^a
System performance	[99,142–146], [98] ^a	[31,147], [66,69] ^a	, , , [29,76,148–151], [55,73,97,132, 152] ^a

^a Demonstrated with actual experiments.**Fig. 2.** Proportion of each category in the reviewed studies.

details about grid-level or system-level models are referred to these two review papers [16,100].

- The geographical distribution of these studies is displayed in Fig. 3, where the color overlays represent the main climate zones [101], triangles are simulation studies, and circles are experimental studies. Most (83.9%) studies are located in the temperate, especially cool temperate zones. Compared with the large population in North America and Europe, there might be great potential to exploit in the large area with a similar climate in Asia. One desired building characteristic for MPC, especially for load shifting applications, is higher thermal mass [102]. Therefore, the smaller number of studies in the tropical area is possibly due to the relatively lighter envelop. However, performance improvement is still achievable [30], and integrating renewable energy may bring better opportunities [103].
- Regarding the modeling and control scales, the proportions of the number of controlled zones are displayed in pie charts (Fig. 4). In total, 34.7% (41) studies are demonstrated in single zones, while only 16.9% (20) are full-scale applications. The lack of large-scale cases, regardless of simulation/experiment, agrees with Rockett & Hathaway [24]. However, the ratio of full-scale demonstration increases to 35.5% considering only experiments, whereas 70.8% simulation studies used less than five zones.

3. Data requirements

Depending on the modeling methods and purposes, as well as building systems, different data or information is used to build the model. Data availability and resolution are critical for model calibration, for either white-box methods [153] or data-driven methods [154]. Meanwhile, data management and utilization have become a challenge with the increasingly available data over the building life cycles [32, 155]. Hence, a framework based on extended Level of Detail (LoD) is proposed to quantify and categorize data (including building metadata and time series data) used for modeling in past studies.

3.1. Definition of extended Level of detail

In the context of Building Information Models (BIM), Level of Detail (LoD), or Level of Development, defines and illustrates inputs and information requirements of the different levels for building elements. This clear articulation allows model authors to define what their models can be relied on for, and allows downstream users to clearly understand the usability and the limitations of models they are receiving [156]. The original theme of LoD aligns with the need for control-oriented models to clarify the required data and to further imply the performances. However, as illustrated in the left part of Fig. 5, the original LoD definition and some extension studies [157,158] all focused on the design and construction phases, considering the static characteristics of

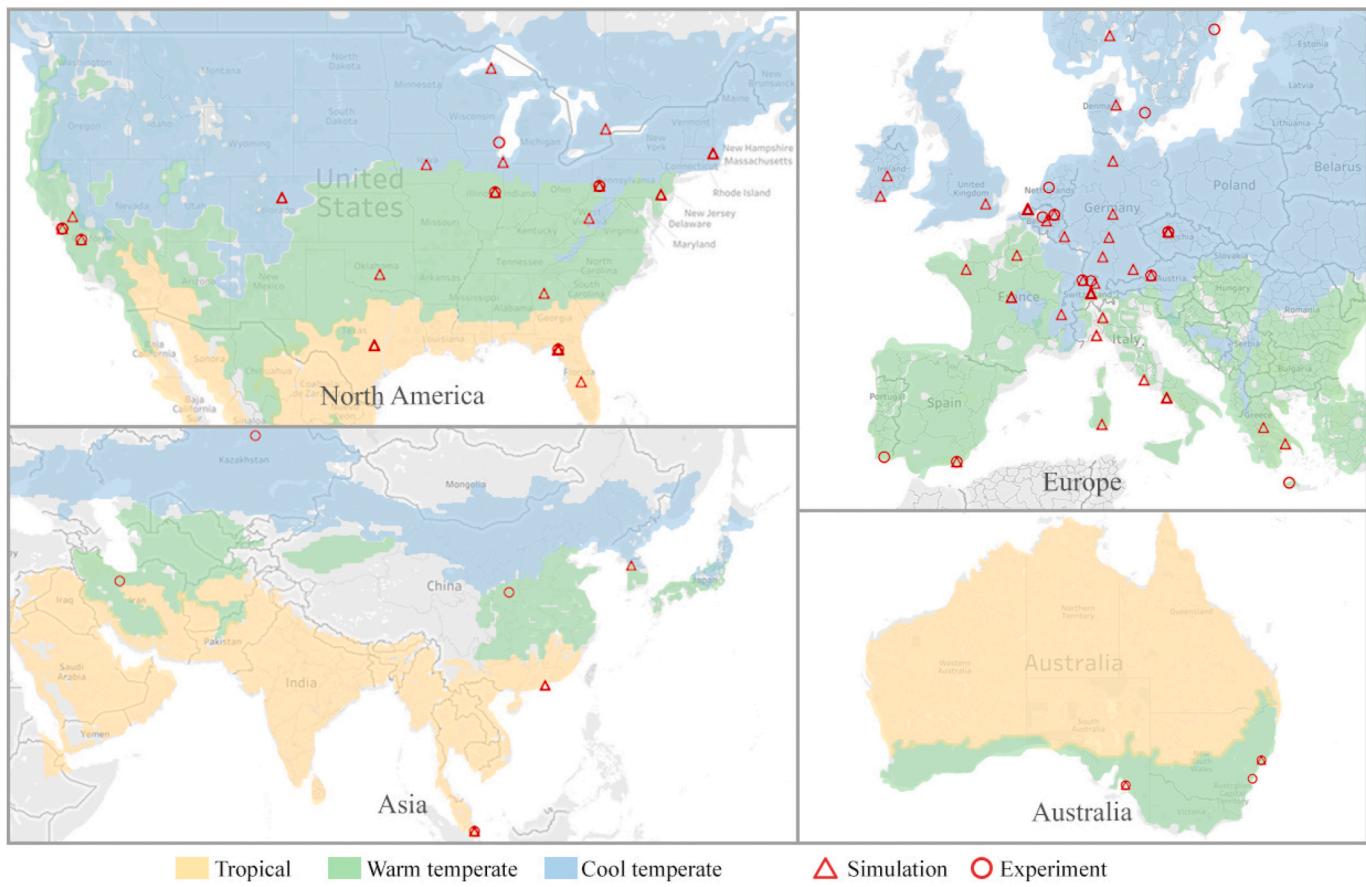


Fig. 3. Geographical distribution of the reviewed simulation and experiment studies. Color overlays reflect three main climate zones.

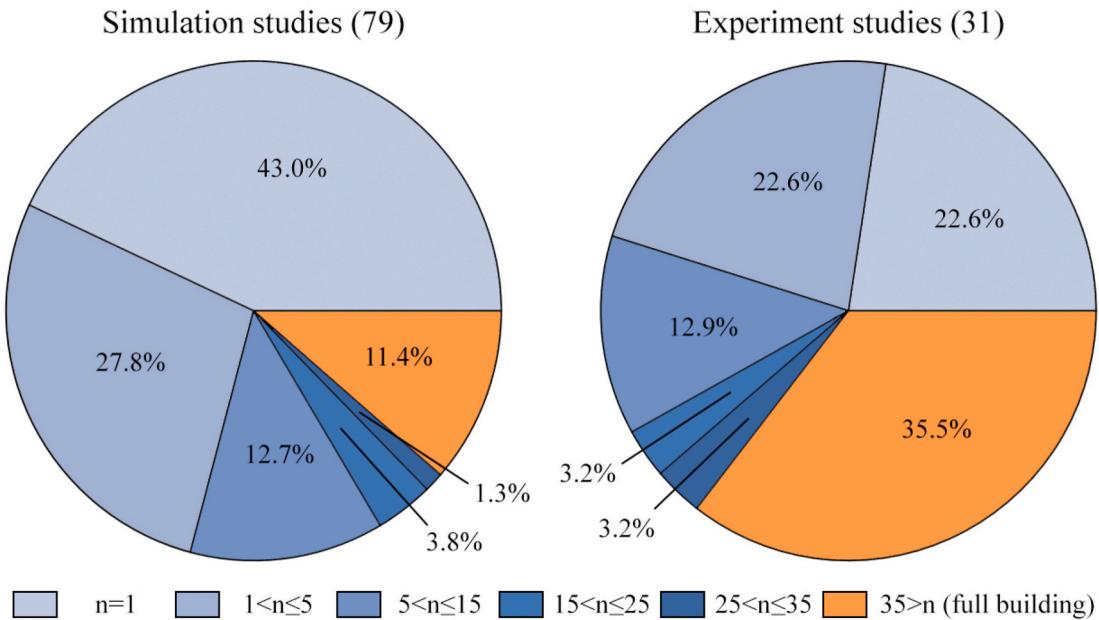


Fig. 4. Proportion of the number of controlled zones (n) in the reviewed studies. Full building applications are counted as >35.

building elements. Becerik-Gerber et al. [159] noted the increasing need for identifying non-geometric data requirements to apply BIM for facility management but overlooked the considerable variation in time series data.

This paper extends the original definition to take in time series sensor

data used for MPC. The extended LoD inherits the original form of using three-digit numbers to describe certain levels, yet endowing every digit with an actual interpretation: time validity, measurement granularity, and temporal resolution. Fig. 6 displays the definitions and interpretations of each digit. The sankey flows denote possible

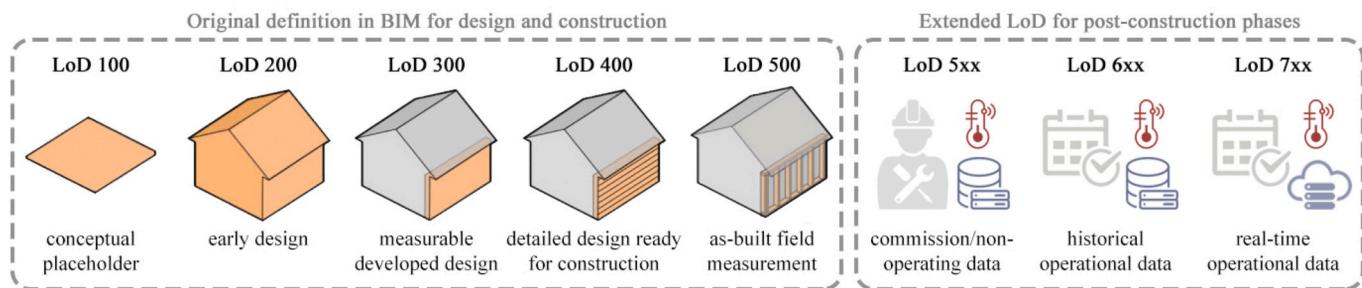


Fig. 5. The original LoD definition in BIM and the extended definition for time series sensor data.

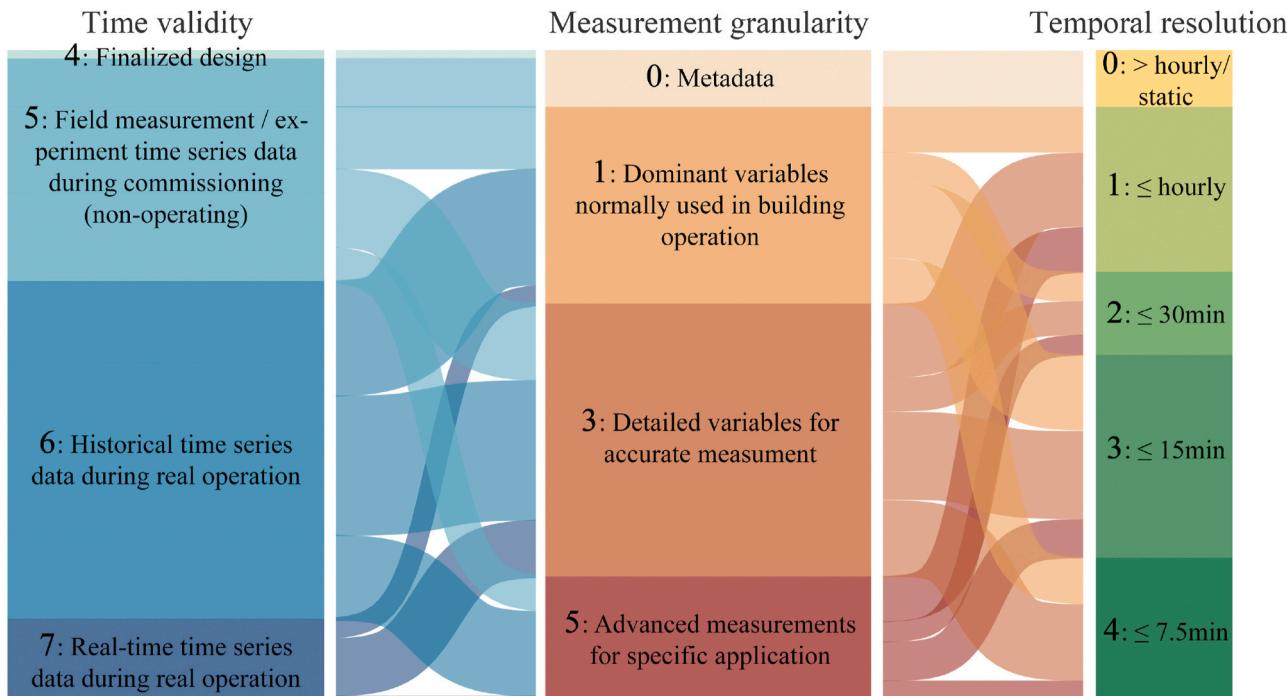


Fig. 6. Definition of the extended Level of Detail (LoD). The sankey flows represent the LoD levels appearing in the literature. The widths of Sankey flows and the heights of bars reflect the usage frequency of each possible level.

combinations of the three digits. Meanwhile, the usage frequencies of each LoD in the reviewed papers are reflected in the widths of sankey flows and the heights of bars. In general, larger LoD implies a higher cost of data acquisition. Apart from the sensor and operating costs, data storage, exchange, processing, and computing all introduce additional costs [160].

3.1.1. Time validity

The first digit represents the time validity of data, i.e., how up-to-date the data is. LoD 400 is simply adopted from BIM for building metadata as detailed and accurate design information, ready for construction. As an example, Kwak et al. [107] used design drawings to build an Energyplus model. Level 5 involves field measurement, including as-built metadata [14] and time series data from designed experiments [139]. The measurement is typically conducted during building commissioning or other non-operating periods to avoid intruding on occupants [73]. Since verification of metadata is usually carried out, the white-box models are considered using level 5 unless explicitly pointed out. Although LOD 400 is rarely mentioned, evolving from that improves the compatibility with BIM definition. Level 6 refers to historical data collected during real operation, which is commonly used for gray-box [33] and black-box [152] model identification. Level 7 stands for real-time operation data, requiring data exchange modules in

the system architecture [30]. For clarification, while real-time data is normally used when implementing control, level 7 is meant to distinguish models that are regularly updated [95,140]. Compared with level 5, time series data at level 6 and 7 are fully exposed to uncertainties and closer to real operation.

3.1.2. Measurement granularity

The second digit means the measurement granularity. Again, level 0 is stuck to the BIM definition for building metadata. For time series data, level 1 includes the principal variables that are usually measured for building operation, level 3 contains the detailed measurements to describe the object more accurately, and level 5 is the advanced measurements taken for specific purposes, subject to customized change. As illustrated in Table 4, The four basic levels of LoD refer to different specific variables for the six data categories: energy consumption, indoor condition, internal disturbance, external disturbance, system condition, and envelop condition. The categories are inspired by the one proposed by Mahdavi & Taheri [161] and modified to fit the data requirements of control-oriented models.

- **Energy consumption (EN)** of the entire building is usually measured for billing purpose, which forms level 1. Different energy sources, such as electricity and gas, are separated if applicable [74].

Table 4

Detailed definition of measurement granularity for the 6 data categories.

Level	0	1	3	5 ^a
Energy consumption	N/A	Total consumption by energy sources	Separated consumption by usage type	Separated consumption by sub-components
Indoor condition	N/A	Indoor air temperature	Variables affecting thermal comfort	Thermal comfort/sensation feedback
System condition	System specifications	On/off operating mode, thermostat setpoints	Temperature and flow rate variables	Static pressures
Envelop condition	Geometric and thermal properties	N/A	Surface or core temperature	N/A
Internal disturbance	Assumed operating schedules/profiles	N/A	Estimated operating profile	Additional occupant sensors
External disturbance	N/A	Weather data of the city/region	On-site weather station/sensors	Solar heat gain on different orientations

^a Level 5 here is illustrated with typical examples. Actual variables might be subject to customized change in specific studies.

For level 3 detailed measurements, the energy consumption is splitted into different end use, such as heating [84], cooling [82], lighting [14], and plug [107]. As advanced measurements, the energy consumption is further disaggregated by sub-components, such as boilers [97], heat pumps [61], pumps [68], and fans [149].

- **Indoor condition (IC)** means indoor thermal comfort conditions in most cases and is usually represented by indoor air temperature [33, 57,58]. For detailed measurements, other factors affecting thermal comfort are involved, including mean radiant temperature [117], humidity [80], and operative temperature [77]. At level 5, special measurements are taken to investigate specific problems. For example, occupant thermal comfort feedback is collected for an occupant-oriented MPC [144], and illuminance is considered for optimal control of blind position [52].
- **Internal disturbances (ID)** are the sources of internal heat gain: occupants, equipment, and lights. The metadata in this category refers to the assumed operating on/off schedules [55] or ratio based profiles [59]. The information can be based on standard or expert knowledge. There is no level 1 because they are not necessary for normal building operations. As level 3, the profiles are estimated based on electricity [97] or temperature [73] trends. CO₂ concentration is the most used advanced measuring method [122,132]. Passive infrared (PIR) sensors and people counters are also used [152]. Internal radiative and convective heat gain are used in many simulation studies [33,58], but is hardly measurable in real operation unless using load emulators [30,76].
- **External disturbances (ED)** are the climate conditions that cause external heat gain. While the dry-bulb temperature and solar irradiance are used much more frequently than other variables like wind speed [77] and ground temperature [10], they are not differentiated in different levels given their similar availability with the existence of weather stations. Level 1 refers to the publicly available weather data of the city or the region [74,94]. On-site weather stations are used to accurately measure the buildings' ambient condition [69,97, 107], therefore defined as level 3. The Typical Meteorological Year (TMY) used in most simulation studies [33,58] are also considered accurate measurements since the buildings are assumed to be under these typical conditions. At level 5, solar heat gain on different orientation serves as an example [43,106].
- **System conditions (SC)** describe how the HVAC systems are operated. Level 0 as static information requires information such as capacity and COP [39]. Note that sometimes COP is assumed to be constant just to estimate energy consumption [65], which is not considered to require level 0 SC information. Level 1 is usually available in building operation, including the on/off operating mode [140] and thermostat setpoints [55]. Detailed measurements cover the flow rates and temperatures on the water side [62] and the air side [73]. These points are often used to estimate thermal loads when the power meters are not in place [30,124]. Valve [81], damper [150] and blind [129] positions are also categorized into level 3 as

they imply the heat flow. As an advanced measurement example, supply air static pressure was taken for system performance estimation [76].

- **Envelop condition (EC)** only has level 0 static characteristics and level 3 detailed measurements. Level 0 may involve geometric properties like areas and volumes [30], and/or thermal properties like U-value and Solar Heat Gain Coefficient (SHGC) [50]. Information like the number of rooms can be easily observed and therefore is not explicitly accounted for. For level 3, surface and/or core temperatures of the envelops are measured. These variables are usually found in buildings with radiant systems such as TABS or Concrete Core Activation (CCA).

3.1.3. Temporal resolution

Larger numbers as the third digit indicate higher temporal resolution of the time series data. Level 0 includes the static building characteristics and time series data with the interval larger than an hour. Time interval of less than or equal to an hour but larger than 30 min falls into level 1. Similarly, less than or equal to 30 min but larger than 15 min belongs to level 2, and so forth. Thereby, level n corresponds to 2^{n-1} to 2^n data points per hour.

3.2. Data requirements of different models

The data usage of the 118 models is categorized and quantified according to the extended LoD framework. The average LoD of the six categories is also calculated for each study to enable the quantitative comparison between different model types. Since the resulting data is unpaired and non-Gaussian distributed, the Mann-Whitney U tests are applied. Representative studies are selected for presentation in Tables 5–7. The selection is done by stratified sampling from each of the nine model types.

3.2.1. Comparing modeling purposes

It is expected to see the most data used in system performance models (Table 7), followed by thermal response models (Table 6), and then load prediction models (Table 5). It is obvious that load prediction models require the least. While the difference between thermal response and system performance models may not be visually detected from Tables 6 and 7, the Mann-Whitney U test on the average LoD gets 0.016 p-value, indicating a significant difference. The medians are 626 for system performance models and 528 for thermal response models. The general data requirements of the three modeling methods are respectively summarized:

- With the assumption on indoor condition, load prediction models requires no IC data. They predict the total or thermal loads (level 1 or 3 measurement granularity) based on past values [89,92,93] or the disturbances [55,86,88]. Most studies assumed constant room

models to estimate internal heat gains (p-value 0.047). This could be a side benefit of deploying more sophisticated sensing systems.

3.2.2. Comparing modeling methods

It can be seen that white-box models mostly use just building metadata of LoD 400 and 500. One exception is Zhao et al. [144] embedded historical occupant thermal comfort feedback to predict thermal comfort. Most gray-box models need information such as the building layout, which is not counted as using metadata. Some used metadata to provide initial guesses [43] or value bounds [83] for parameter identification. The usage of metadata was shown crucial for some identification algorithms, especially in cases like MRI, when the optimization problem is non-convex [59]. It is worth noting that metadata, particularly the operation profile, was also used in black-box models [79,84]. As another example, Li et al. [151] used the system specification to obtain the heat pump performance and other temperature measurements to model other sub-systems.

Gray and black-box models use a similar amount of time series data. The statistical test gives 0.232 p-value, showing no significant difference. In fact, gray-box models have slightly higher median average LoD (617 over 606). This observation disagrees with the conclusion in Ref. [20] that gray-box models require less data. It is still arguable that the difference lies in the training data length. The length is not quantified in this review because many studies did not report. However, training data length varying from one day to one year is found in both modeling methods.

3.2.3. Other comparisons

Although the modeling methods greatly impact the usage of metadata, when it comes to time series data, the usage of level 5, 6, and 7 is almost evenly distributed in different model types. Faster system dynamics requires higher temporal resolution. For example, the average resolution level of Fan Coil Unit (FCU) models is 3.75, while the average of TABS models is 2.11. To summarize, among the three dimensions of LoD, time validity is partially influenced by the modeling methods, measurement granularity is affected by the modeling purpose, and temporal resolution is typically decided by the system dynamics.

4. Performance evaluation

A control-oriented model with acceptable prediction capability is the prerequisite to achieve good control. Potential model mismatch could lead to control performance degradation. Therefore, it is crucial to obtain a structured perception of model performance evaluation and its impact on control. With the intention to quantify the relationship, two gaps are noticed: a) the absence of a standard or comparable approach for model performance evaluation, and b) the paucity of research remarking the relationship. This section gives an overview on these two issues.

4.1. Model evaluation

The approaches to evaluate the credibility of a computerized model can be categorized into validation and verification [162]. Validation substantiates that a model, within its domain of applicability, possesses a satisfactory range of accuracy. Verification substantiates that a model represents the conceptual model within specified limits of accuracy.

4.1.1. Validation metrics

Most studies evaluate the model through validation. A number of them adopted a qualitative approach to plot the model outputs with the test data and show a good tracking [62,72,73]. To quantify the error over a period, Mean Bias Error (MBE, Equation (1)) is a basic metric. However, it is rarely used because positive and negative errors cancel each other when summing up and may distort the results. Therefore, Mean Absolute Error (MAE, Equation (2)), Mean Squared Error (MSE,

Equation (3)), and Root Mean Squared Error (RMSE, Equation (4)) are normally used [50,86,123]. Concerning about the variation of error, Maximum Absolute Error (MaxAE, Equation (5)) and Standard Deviation of Absolute Error (StdAE, Equation (6)) were also used [29,124]. These metrics are useful to avoid the potential thermal comfort violation caused by the model mismatch [97,139]. There is no consensus on the acceptable error, but many studies tried to contain the errors within $\pm 1^{\circ}\text{C}$.

$$\text{MBE} = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i) \quad (1)$$

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |(\hat{y}_i - y_i)| \quad (2)$$

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 \quad (3)$$

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2} \quad (4)$$

$$\text{MaxAE} = \max_{i \in [1,n]} |(\hat{y}_i - y_i)| \quad (5)$$

$$\text{StdAE} = \sqrt{\frac{1}{n-1} \sum_{i=1}^{i=1} \left(|(\hat{y}_i - y_i)| - \text{MAE} \right)^2} \quad (6)$$

To diminish the effect of absolute value scales when comparing different models, the error metrics are normalized to obtain percentages. Mean Absolute Percentage Error (MAPE, Equation (7)) and Coefficient of Variation (RMSE) (CV(RMSE), Equation (8), also known as Normalized RMSE) are commonly used [76,132]. Replacing the average of the measured value in the denominator of equation (8) with the range of predicted values yields standardized RMSE [43]. However, comparing models with these normalized percentages still requires some caution. For instance, consider model A that predicts room temperature around 26°C with 1°C RMSE and model B that gives similar RMSE around 20°C , the lower CV(RMSE) of model A does not make it more accurate. Alternatively, R squared (R², Equation (9), sometimes referred to as the goodness of fit), estimating the ratio of explained variance in the prediction, is frequently used as well [77,95].

$$\text{MAPE} = \frac{1}{n} \sum_{i=1}^n \left| \frac{(\hat{y}_i - y_i)}{y_i} \right| \quad (7)$$

$$\text{CV(RMSE)} = \frac{1}{\bar{y}} \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2} \quad (8)$$

$$R^2 = 1 - \frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (9)$$

Usually, closed-loop, or one-step ahead, prediction is validated. As a stricter approach, open-loop prediction was validated on occasion to examine the model performance over the horizon [44]. For example, multi-step RMSE (MS-RMSE, Equation (10)) was applied [27]. In equations (1)–(10), \hat{y} is the model output, y is the test data, \bar{y} is the average of test data, p is the prediction horizon, n is the length of test data, i and k refer to the time step.

$$\text{MS-RMSE} = \sqrt{\frac{1}{p(n-p)} \sum_{i=1}^{n-p} \sum_{k=1}^p (\hat{y}_{i+k|i} - y_{i+k})^2} \quad (10)$$

One consideration is the dataset used for validation. The mostly applied in practice is the historical data over a period of time. It can be obtained either from the real operation or using a high-fidelity simulation model. Otherwise, using real-time data, Finck et al. [90] validated the model outputs against measurements during the control experiment. Several datasets were designed for the purpose of better examining the models' prediction capability. Kim et al. [130] suggested validating the model in a cross-validation manner if significant disturbances present in operations. Li et al. [82] tested the extensibility (extrapolating capability) by designing scenarios when the boundary conditions exceed the range of the training dataset. Several studies made datasets with step signals to check impulse responses between input-output pairs [27,163]. These designed datasets are usually generated using a high-fidelity simulation model. It was noted in the field of hydrology that the uncertainty in measured data should be carefully contained to effectively validate the models [164], which was rarely concerned by building control studies.

4.1.2. Verification and identifiability

It has been recognized that a desired control-oriented model should not only predict with small error but also represent the actual building systems [27]. In line with this, a desired [58] or minimum [34] model complexity is needed by different building systems. Verification can also be done by physically interpreting the model outputs or parameters. In the frequency domain, model responses to different input stimuli were compared. The zone thermal reaction to heat input at different frequencies was examined [34], and "error with respect to input" was defined to quantify the performance [163]. Privara et al. [52] examined the whiteness of residuals using a cumulative periodogram, confirming that the system dynamics was properly modeled, and the residuals were caused by noises. For parameter inspection of gray box models, the identified RC values were compared with the physical meaning [58, 125]. Significance index and correlation index were introduced to evaluate how the parameters affect the model performance [66]. A model is expected to be less sensitive to parameter perturbation [37] and have less correlated parameters [66].

As a model becomes more complex and has more degree of freedom, identifiability comes to be an issue [126,153]. If a model is over-parameterized, i.e., has too many parameters to identify, parameter estimation, model evaluation, and further application could be more difficult [56,165,166]. It has been shown that the parameter values can be varied a lot without significantly changing the validation results [83,151]. The parameter identifiability can be decomposed into structural and output identifiability [123]. To correctly identify the models instead of overfitting the training data, an appropriate model structure is the most important factor [167]. Therefore, model selection methods were applied for both white-box [168], gray-box [43,127] and ANN models [80]. Further, Privara et al. [57] applied different model selection criteria for probabilistic and deterministic identification algorithms.

4.2. Relating to control performance

4.2.1. Control performance evaluation

MPC in buildings aims at energy saving, thermal comfort improvement, peak load reduction, system efficiency improvement, etc. The effectiveness can be demonstrated through simulation or experiment. A realistic simulation-based demonstration framework requires a control-oriented model and a high-fidelity simulation model. The simulation model serves as a virtual testbed, to which the control action based on the controller model is applied [25]. Tools and desired features of simulation models are reviewed in Ref. [169]. Typical tools include Building Controls Virtual Test Bed (BCVTB) [170], TRNSYS (type 15,17, 56,155) [171], and Modelica [172]. However, a number of simulation-based studies deployed an idealistic framework by using the same model for optimization and simulation. 72.0% (85) of the

categorized papers are demonstrated by simulation, out of which 41.2% (35) are idealized. Assuming that the controller model is perfectly representing the building overlooks the influence of the model performance on the control performance. Studies using white-box methods, such as model reduction, typically tend to be idealized. Yet, without comparing to experimental data or higher fidelity models, the effectiveness of model reduction methods was claimed to be questionable [44].

Unlike simulation-based demonstration, where different control strategies can be compared under the same boundary conditions, the comparison is not as straightforward in actual experiments. The most used approach is normalizing the control results, usually energy consumption, by degree days [62,97] or outdoor temperatures [173]. Alternatively, some studies showed that different strategies are applied under similar averages [59] or profiles [98] of outdoor temperature.

Baseline selection is another concern when evaluating the control performance. The most convenient way is to compare with the default control in the building. However, it was argued that the default settings in BMS are possibly poor-tuned, disputing the improvement of control performance [24]. It was noted that the saving potential brought by MPC could also be achieved by fine-tuning the rule-based controller (RBC) [152]. Accordingly, the RBC was pre-tuned to consolidate the control comparison [95,106]. Moreover, MPC with simpler configurations were used to show the superiority of robust MPC [123], non-linear MPC [117], and system performance MPC [99]. The upper performance bound of MPC was quantified by using the perfect model and disturbance prediction [13]. Additionally, different combinations of HVAC systems and control algorithms were considered as integrated baselines [149].

4.2.2. Affecting factors

The performance of MPC varied among different situations from worse than baselines to over 100% better. Apart from the model performance, the wide range is also affected by factors including building characteristics [102], ambient conditions [94], operation constraints [39], disturbances [122], and etc. These factors function in a combined and complicated way. For example, the impact of internal disturbance was moderated during the heating season, as compared with the cooling season [31]. Consequently, extracting the relationship between model and control performance involves explicitly designed experiments, which is rare in the past years.

Several studies showed that model mismatch could result in more energy consumption [119,132] and/or discomfort [97,139]. To quantify, 10% error led to 5% more energy cost and 100% more comfort violation [28]. On the other hand, a more accurate model, in terms of RMSE, did not necessarily lead to better control performance [33]. The prediction horizon matters as well. Zong et al. [10] found the energy cost decreased and then increased as the prediction horizon increased, which might be relevant to the open-loop prediction accuracy. Regarding the model structure and characteristics, a multi-zone model achieved better thermal comfort than a simplified single-zone model with a similar amount of energy [83], a certain number of states was found necessary to capture the thermal dynamics [34], and a non-convex model caused multiple local optima in the optimization, some of which deteriorated the control [174].

5. Discussion

Through the categorizations and discussions in the last three sections, several research gaps are spotted. This section summarizes the review with five directions for future study, three of which regarding the data requirements and the other two about the performance evaluation. These topics interrelate and should be studied together systematically.

5.1. What are the minimum data requirements to build a control-oriented model?

The answer to this question is subject to modeling purposes, modeling methods, and building systems (section 3.2). While the variation between different modeling purposes or building systems is clear, the border, in terms of data requirements, between modeling methods is not. Due to the lack of description, the usage of building metadata is mostly referred to as LoD 500. However, LoD 400 can be useful to build either white-box [175] or gray-box [43] models. On the other hand, the usage of time series data is typically omitted in white-box studies, although it is involved in both manual and automated calibration [176]. Better describing the data usage would help justify the modeling effort and correspondingly the scalability of the proposed methods. Therefore, future research would benefit from explicitly quantifying the data requirements.

Another important question to answer is whether extra excitation or regular update is necessary for model identification. As shown in Fig. 6, LoD 6xx is mostly needed, but both 5xx and 7xx are also used in many studies. It was argued that the operation data generated by normal operations led to poor identification, so the data must fulfill specific requirements [163]. There are pretty mature methodologies to generate excitation signals [177], but these experiments are usually not compatible with normal operations. On the other hand, unoccupied experiments are not exposed to the uncertainties brought by occupants and other internal disturbances [8]. As a potential solution, the excitation could be partially conducted within the comfort range during daily operations [147]. As for LoD 7xx, regularly-updated models are found to be desired in several cases [22,56]. The problem is the increased cost for data exchange and computation. Also, it is practically impossible to provide full excitation and regular update simultaneously.

The proposed LoD covers the variations in time validity, measurement granularity, and temporal resolution. Besides, there are other factors found in the literature. The length of training data is not counted in the current framework since no pattern is observed. However, it may change the cost of data acquisition, especially for LoD 5xx and 7xx. Hence, it will be added as an attribute in the future plan of reforming the LoD in an object-oriented way. The spatial resolution of indoor conditions also makes a difference when modeling large scale buildings. Reference room temperature [78] and average room temperature [97] were used in different studies. These choices are made mainly to reduce model complexity. Considering the room temperature is usually available with the thermostats, the spatial resolution is not specified in the framework. The data quality or sensor accuracy is another factor that might affect the modeling procedure. Yet, it is hard to estimate and rarely declared. Thus, unbiased measurements are assumed.

5.2. How to balance the trade-off between model complexity and data requirements?

When building a control-oriented model, fewer assumptions and better prediction capability are desired [20], calling for more states and higher complexity. For example, adding a state for wall temperature improved the performance [10]. On the other hand, higher-order models require more data to identify [58]. Insufficient training data could result in issues of identifiability (section 4.1.2). Inadequately informative data could also deteriorate model performance [178]. Therefore, it is essential to find a balanced point between model complexity and data requirements. The prediction error of a numerical model can be generally attributed to aleatoric uncertainty and epistemic uncertainty [179]. The epistemic uncertainty can be further decomposed into structural and parameter uncertainties, which paves the way to balance the trade-off [180]. illustrated an abstracted relationship between model uncertainty and complexity for HVAC simulation models. Inspired by the idea, the qualitative relationship between model complexity, potential prediction error, and data availability is depicted

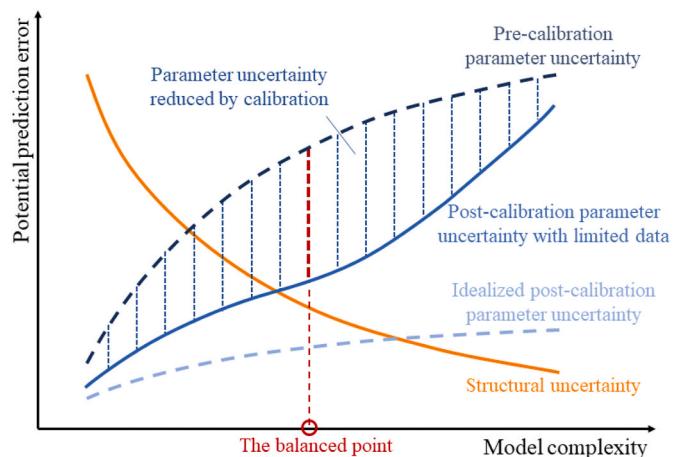


Fig. 7. Qualitative relationship between model complexity, potential prediction error, and data availability.

in Fig. 7.

As the model complexity increases, the model describes the physical process more precisely, reducing the structural uncertainty (orange line). On the contrary, the dark blue dash line represents the parameter uncertainty, which increases with the number of parameters in the model. With the building metadata and time series data, calibration (including calibration for white-box models and parameter identification for data-driven models) is conducted to bring down the parameter uncertainty. The cyan dash line refers to an idealized situation, where all data requirements can be fulfilled, and the parameter uncertainty can be mostly removed. The remaining part could be caused by aleatoric uncertainties. More complex models require more data for calibration and benefit more from the calibration. In such an idealized situation, the most complex model may be favored.

In reality, however, building systems are hardly data-rich [181]. The limited amount of available data results in the solid blue line in between. Given a specific level of data availability, if the model complexity is too high or too low, the data would be overfitted or underfitted [33]. Consequently, the parameter uncertainty reduced by the calibration, marked by the vertical dash lines, will increase and then decrease as the complexity increases. Thereby, the balanced point is where the improvement brought by the calibration reaches its maximum, highlighted in red. This point corresponds to the best performance that can be achieved with the data. As more data is available, a more complex model could be supported, and the balanced point could be shifted to the right, yielding better prediction capability and smaller potential error. More research is needed to quantify this relationship.

Pragmatically, transferability is a barrier to the commercial application of MPC in buildings [182]. The configuration effort remains high every time when it comes to a new building. With the LoD framework and the quantified relationship, the upper bound of model and control performance of a building can be estimated based on its data availability. Along this line, an automated modeling framework could promote the scalability and potential of MPC.

5.3. What is the significance or potential benefits of higher measurement granularity for the role of occupants?

Occupant behavior is one of the major uncertainty sources in building energy assessment [183]. Therefore, more accurate measurements and estimations are desired. Compared with other data categories, great variations exist in the advanced measurements of internal disturbances. CO₂ and PIR sensors were used in past MPC studies [122, 152]. However, the reliability of both methods was questioned [184, 185]. Meanwhile, there are other advanced techniques available such as

- [171] S. Klein, W. Beckman, J. Mitchell, J. Duffie, N. Duffie, T. Freeman, J. Mitchell, J. Braun, B. Evans, J. Kummer, et al., Trnsys 17: a transient system simulation program. solar energy laboratory, Madison, Madison, USA: University of Wisconsin.
- [172] Wetter M, Zuo W, Nouidui TS, Pang X. Modelica buildings library. *J Build Perform Simulat* 2014;7(4):253–70. <https://doi.org/10.1080/19401493.2013.765506>.
- [173] Chen B, Cai Z, Bergés M. Gnu-rl: a precocial reinforcement learning solution for building hvac control using a differentiable mpc policy. In: Proceedings of the 6th ACM international conference on systems for energy-efficient buildings, cities, and transportation; 2019. p. 316–25. <https://doi.org/10.1145/3360322.3360849>.
- [174] Kelman A, Ma Y, Borrelli F. Analysis of local optima in predictive control for energy efficient buildings. *J Build Perform Simulat* 2013;6(3):236–55. <https://doi.org/10.1080/19401493.2012.671959>.
- [175] Chong A, Xu W, Chao S, Ngo N-T. Continuous-time bayesian calibration of energy models using bim and energy data. *Energy Build* 2019;194:177–90. <https://doi.org/10.1016/j.enbuild.2019.04.017>.
- [176] Coakley D, Raftery P, Keane M. A review of methods to match building energy simulation models to measured data. *Renew Sustain Energy Rev* 2014;37:123–41. <https://doi.org/10.1016/j.rser.2014.05.007>.
- [177] Zhu Y. Multivariable system identification for process control. Elsevier; 2001.
- [178] Chong A, Augenbroe G, Yan D. Occupancy data at different spatial resolutions: building energy performance and model calibration. *Appl Energy* 2021;286:116492. <https://doi.org/10.1016/j.apenergy.2021.116492>.
- [179] Der Kiureghian A, Ditlevsen O. Aleatory or epistemic? does it matter? *Struct Saf* 2009;31(2):105–12. <https://doi.org/10.1016/j.strusafe.2008.06.020>.
- [180] Trčka M, Hensen JL. Overview of hvac system simulation. *Autom ConStruct* 2010;19(2):93–9. <https://doi.org/10.1016/j.autcon.2009.11.019>.
- [181] T. Dixon, S. Bright, P. Mallaburn, K. B. Janda, C. Bottrill, R. Layberry, Learning from the “data poor”: energy management in understudied organizations, *J Property Invest Finance*:10.1108/JPIF-03-2014-0018.
- [182] Schmidt M, Åhlund C. Smart buildings as cyber-physical systems: data-driven predictive control strategies for energy efficiency. *Renew Sustain Energy Rev* 2018;90:742–56. <https://doi.org/10.1016/j.rser.2018.04.013>.
- [183] Tian W, Heo Y, De Wilde P, Li Z, Yan D, Park CS, Feng X, Augenbroe G. A review of uncertainty analysis in building energy assessment. *Renew Sustain Energy Rev* 2018;93:285–301. <https://doi.org/10.1016/j.rser.2018.05.029>.
- [184] Fisk WJ. A pilot study of the accuracy of co2 sensors in commercial buildings. Tech. rep. Lawrence Berkeley National Laboratory; 2008.
- [185] Jin Y, Yan D, Sun H. Lighting system control in office building using occupancy prediction based on historical occupied ratio. *EES (Ecotoxicol Environ Saf)* 2019;238(1):012009. <https://doi.org/10.1088/1755-1315/238/1/012009>.
- [186] Chen Z, Jiang C, Xie L. Building occupancy estimation and detection: a review. *Energy Build* 2018;169:260–70. <https://doi.org/10.1016/j.enbuild.2018.03.084>.
- [187] Melfi R, Rosenblum B, Nordman B, Christensen K. Measuring building occupancy using existing network infrastructure. In: 2011 international green computing conference and workshops. IEEE; 2011. p. 1–8. <https://doi.org/10.1109/IGCC.2011.6008560>.
- [188] Hu S, Yan D, Azar E, Guo F. A systematic review of occupant behavior in building energy policy. *Build Environ* 2020;106807doi. <https://doi.org/10.1016/j.buildenv.2020.106807>.
- [189] Park JY, Nagy Z. Comprehensive analysis of the relationship between thermal comfort and building control research-a data-driven literature review. *Renew Sustain Energy Rev* 2018;82:2664–79. <https://doi.org/10.1016/j.rser.2017.09.102>.
- [190] Jung W, Jazizadeh F. Energy saving potentials of integrating personal thermal comfort models for control of building systems: comprehensive quantification through combinatorial consideration of influential parameters. *Appl Energy* 2020;268:114882. <https://doi.org/10.1016/j.apenergy.2020.114882>.
- [191] Arnold JG, Moriasi DN, Gassman PW, Abbaspour KC, White MJ, Srinivasan R, Santhi C, Harmel R, Van Griensven A, Van Liew MW, et al. Swat: model use, calibration, and validation. *Trans ASABE* 2012;55(4):1491–508. <https://doi.org/10.13031/2013.42256>.
- [192] Drgoňa J, Arroyo J, Cupeiro Figueroa I, Blum D, Arendt K, Kim D, Ollé EP, Oravec J, Wetter M, Vrabić DL, Helsen L. All you need to know about model predictive control for buildings. *Annu Rev Contr* 2020;50(1367-5788):190–232. <https://doi.org/10.1016/j.arcontrol.2020.09.001>.
- [193] Blum D, Jorissen F, Huang S, Arroyo J, Benne K, Li Y, Gavan V, Rivalin L, Helsen L, Vrabić D, et al. Prototyping the boptest framework for simulation-based testing of advanced control strategies in buildings. In: Proceedings of the international building performance simulation association, international building performance association (IBPSA); 2019.