

# **Notes paper - A data-driven and scalable approach for window operation detection in multi-family residential buildings**

Juliet Nwagwu Ume-Ezeoke, Kopal Nihar, Catehrine Gorle, Rishee Jain

June 13, 2023

## **1 Abstract**

Natural cooling, utilizing non-mechanical cooling, presents a low-carbon and low-cost way to provide thermal comfort in residential buildings. However, naturally cooled residential buildings are difficult to design without a clear understanding of the effect opening and closing windows will have on occupant thermal comfort. Predicting when and understanding why occupants open windows is a difficult task in itself. Although work has been done on creating predictive window opening and closing models, they traditionally rely on large datasets, specialized sensing, and/or require training data for each building. This makes them ineffective tools for scaling natural cooling to a large number of buildings. Here, we propose a novel unsupervised window opening and closing detection method that utilizes easy-to-deploy off-the-shelf temperature and humidity sensors. We assess our proposed models' effectiveness on an empirical dataset. We also compare the method's performance with a support vector machine (SVM) model that represents the state-of-the-art in window detection. We find that the proposed method outperformed the SVM on a number of key indicators, but had subpar performance on data in which the indoor and outdoor temperature had small differences. While the SVM's performance was sensitive to the characteristics of the time series features, our proposed method, which relies on indoor temperature alone, demonstrated robust performance in our pilot studies, making it a potential candidate for developing a highly scalable and generalizable model for window operation detection. Overall, this work aims to demonstrate the viability of using unsupervised data-driven methods for understanding window operations in residential buildings. In doing so, we aim to enable more accurate modeling of naturally cooled buildings and further catalyze widespread adoption of this low-cost and low-carbon technology.

## 2 Introduction

## 3 Literature Review

Several research studies have been performed to detect window opening/closing actions in the past. Environmental parameters such as indoor and outdoor temperature, relative humidity, wind speeds, solar radiation, etc. have been widely investigated, and numerous studies have found window operation to be affected primarily by these factors. Almost all of these studies have considered indoor and outdoor temperatures and deemed them to be significant. Data-driven machine learning models have been employed to predict these dynamic window actions taken by the users. These models predict the probability of either the window state (i.e. open or close) or a certain action being taken (opening or closing a window). Supervised ML models such as Logistic regression and Markov chain models, followed by artificial neural networks, have been heavily relied upon to study the probabilistic correlation between the window state and explanatory variables monitored. Other ML models such as support vector machines, and random forest classifiers have been explored in more recent studies to predict window-opening behavior based on both environmental and contextual features. However, these supervised models often rely on labeled long-term data collection, usually spanning over 6 months or multiple seasons. They suffer from a lack of generalizability as their results have to be specifically tuned for the context in which data was collected, either constrained by building type or climate. This makes them ineffective tools for scaling natural cooling to a large number of buildings.

Additionally, while environmental parameters, for eg, air temperature and relative humidity are commonly used, sensors to monitor window opening actions are relatively difficult to install and could be intrusive for building occupants. Therefore, we propose a novel unsupervised window opening and closing detection method that utilizes easy to deploy off-the-shelf temperature and humidity sensors.

As window operation detection is commonly a binary classification task in machine learning, several metrics have been used to assess the prediction performance of window status. These evaluation metrics, frequently referred to in literature are overall accuracy of the model, precision, recall, f-1 score, coefficient of determination ( $R^2$ ), mean average error (MAE), and root mean square error (RMSE).  $R^2$ , MAE, and RMSE are typically used when the output of model is numeric such as the probability of the action being taken. Otherwise, the ability to correctly estimate the window state (open or close) is normally assessed by accuracy, precision, recall, and f-1 scores. Some domain specific metrics for evaluating the prediction performance of the window open/close classification task have also been proposed in previous studies in order to investigate the overall accuracy of the model as well as the consistency of true estimates for window states. These include ratio of total window opening time over total monitoring period, number of actions taken, and median opening/closing durations. Recently, de Rautlin et al. proposed domain oriented metrics to improve coherence of results such as number of true and false openings, total time of true and false openings, and average true opening accuracy

score -> we will discuss in more detail further in the paper, as we will evaluate our model based on these metrics, in addition to, our developed custom metrics.

## 4 Methods

TODO brief overview here...

### 4.1 Data Collection

We collected data for short periods of time over three months in the summer. We placed HOBO sensors<sup>1</sup> temperature/ relative humidity sensors in two adjacent rooms in a multi-family residential building. The opening and closing schedules were not determined beforehand but were allowed to proceed naturally. This created three distinct scenarios for testing the efficacy of the model we developed.

Descriptions of the data are provided in **?@tbl-data-collected**. The experiments are labeled as experiments A, B, and C, corresponding to data recorded at the beginning on July 20, July 27, and September 8 respectively. Experiment B was much longer than the other experiments, about 14 days vs experiments A and C were collected over 3 and 4 days each. Experiment B also had a much shorter period of the window being open, with the window open around 5 percent of the time, vs experiments A and C which had the window open for at least 20 percent of the time.<sup>2 3</sup>

	A	B	C
Starting Day	2022-07-20 07:15:00	2022-07-27 09:00:00	2022-09-08 08:00:00
Data Length	4 days 00:00:00	14 days 10:45:00	3 days 00:00:00
Room	1.0	0.0	0.0
Opening Percentage	0.412987	0.051873	0.221453
Average Open Time	0 days 09:24:59.999999999	3 days 10:15:00	0 days 18:00:00

Figure 1: Data Collected

To develop an understanding of how temperature and relative humidity, the quantities that were measured, changed across the different experiments, we examined the distributions of the data in Figure 2. We observe that experiment B had the widest distribution of internal temperatures. Experiment C had internal temperatures that were higher and closer to the mean/median\*\* of the external temperatures. This was due to observed higher temperatures

<sup>1</sup>sensor name

<sup>2</sup>going to call this EWM from now on and fix when editing/ come up with a better name

<sup>3</sup>TODO include tables underneath the graphs or annotate the graphs showing the best performing metrics..

during this time period, which makes for a particularly interesting edge case. The internal relative humidity was higher than the ambient relative humidity across experiments.

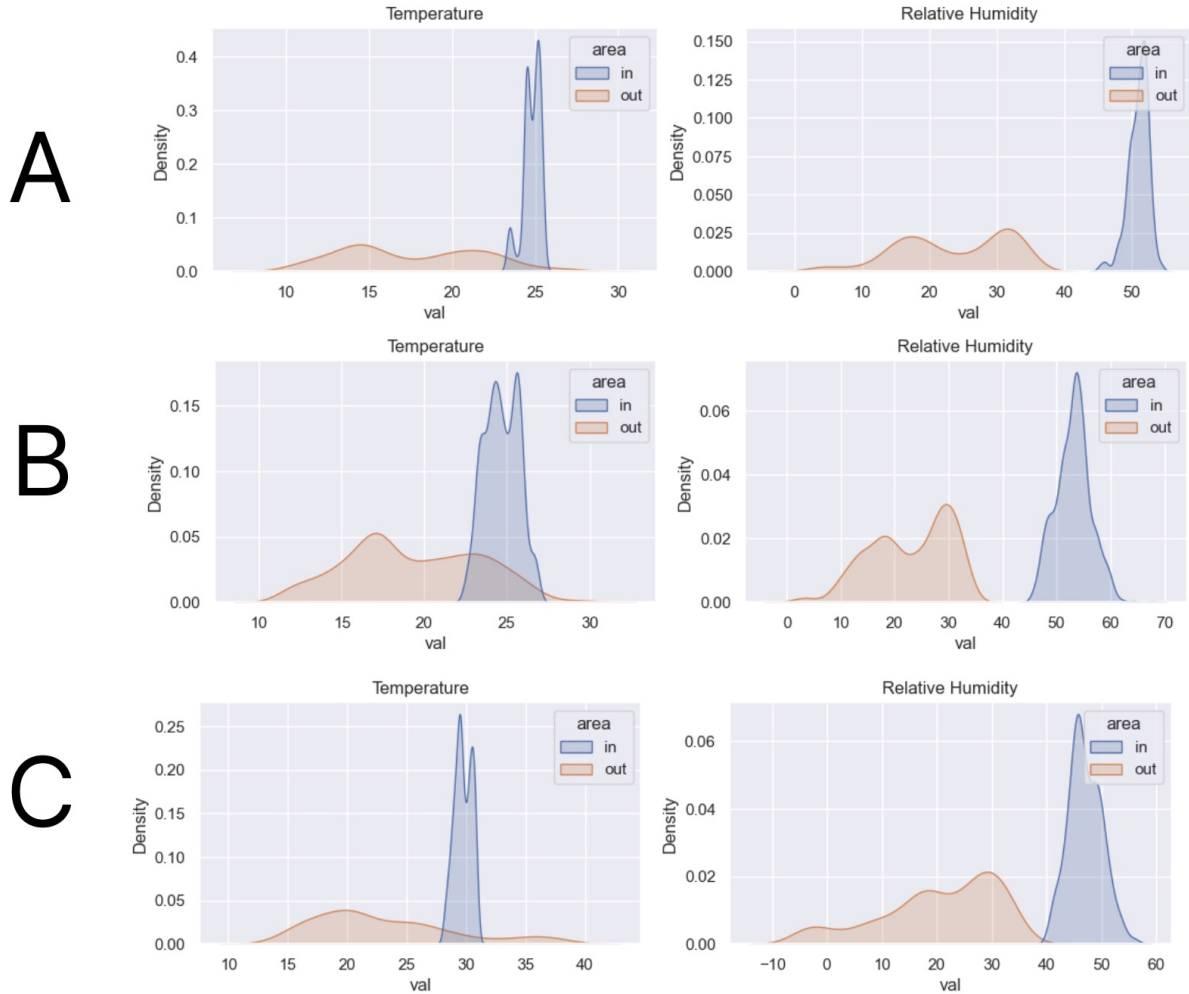


Figure 2: Distributions of Temperature and Relative Humidity Across Experiments

## 4.2 Detection Methods

some intro here <sup>4</sup>

### 4.2.1 Smoothing Technique

*Why we chose to develop this method*

---

<sup>4</sup>understanding?

Our approach relies on an intuitive understanding of window operation detection. We expect a typical time series recording of a quantity of interest, such as temperature or relative humidity, to contain information that reflects the seasonality of ambient quantities, noise due to occurrences within the space where measurements are being taken, and the desired signal of due to a change in the window state. The way in which these three components of a measurement are combined is unknown, and, the noise component in particular cannot be known based on the measurements we have collected. Therefore, our task is to attempt to extricate the seasonality that is due to ambient processes from the measured data, which will ideally reveal the window detection and a bit of noise.

To remove the seasonality from the measured data, we examined three methods of de-seasonalizing: using seasonal-trend-decomposition, using a sinusoidal fit, and using an exponentially weighted mean function. These all reflect a trade-off between capturing pure seasonality and capturing a bit of noise. The sinusoidal method, which is optimized for the indoor temperature time series, purely reflects seasonality whereas the exponentially weighted mean function is simply a moving average of the indoor temperature signal. Therefore, some noise is captured as well. The seasonal component of a seasonal trend decomposition is somewhere in the middle. We found that the exponentially weighted mean function performed the best in isolating window detection, which follows from our intuition that both the seasonal components and noise need to be isolated. <sup>5</sup>

#### *How this method works*

The smoothing technique proceeds as follows. The goal is to identify  $W(t)$ , window state as a function of time. This can take on two values: 0, representing window closed, or 1, representing window open. We have an observed variable  $T(t)$ , which represents the measurement of the indoor temperature. We apply an exponentially weighted mean function to the observed temperature, creating a smoothed time series,  $\overline{T(t)}$ , which in theory removes strong peaks that would reflect changes in window state, and isolates information concerning the seasonal response and additional noise <sup>6</sup>. Subtracting  $\overline{T(t)}$  from  $T(t)$  yields  $T'(t)$ , which is a time series that reflects changes in the window state and some additional noise. In order to more confidently identify the where the changes in window state occur, we examine the first and second derivatives of  $T'(t)$ ,  $\frac{dT'(t)}{dt}$  and  $\frac{d^2T'(t)}{dt^2}$ . Observing the second derivative in particular, is an effective way to identify change points. In order to predict where window changes occur, we apply a hypothesis test <sup>7</sup>. We assume that the time series  $\frac{d^2T'(t)}{dt^2}$  is normally distributed. Therefore, any value in  $\frac{d^2T'(t)}{dt^2}$  that is more than 2 standard deviations away from the mean is unlikely to occur, and could possibly indicate an instance of a change in window state. We will call these initial guesses  $G(t)$ . They take on positive or negative values depending on whether they are predicting a transition from window open to close, or window close to open. Therefore, we round the values of  $G(t)$  to 0 or 1 to reflect this. Finally, we interpolate

---

<sup>5</sup>drdr

<sup>6</sup>need to have subfigure labels, figures a - g...

<sup>7</sup>TODO: Some discussion here about poor performance on experiment C that can be attributed to the similarity between indoor and outdoor temperatures

between the rounded values of  $G(t)$ , so that we only predict a change in window state when  $G(t)$  transitions between 0 and 1. This prediction of the window state is called  $I(t)$ . We have included **?@fig-smoothing-tech** to show how the smoothing technique for window detection works.

[Smoothing Technique in Action ]{#fig-smoothing-tech}

#### 4.2.2 Machine Learning Method

*Why we chose to use this method*

We used a separate method in order to compare the efficacy of the smoothing technique that we developed. We chose to use an SVM, as it has been shown in literature (de Rautlin de Roy) to have robust performance across a wide range of features. SVMs are straightforward to implement with relatively few hyper-parameters, as compared to other methods that have been recently shown to have high performance on the window detection problem. <sup>8</sup>

*How this method works*

SVMs extend simple linear<sup>9</sup> regressions. In a similar way to a linear regression, SVMs approximate a line of best fit using the features and labels provided <sup>10</sup>. However, with the use of kernel functions that can map high dimensional feature spaces to ones of lower dimensions, SVMs also effectively create hyperplanes of best fit in order to classify datasets. The use of a kernel function also enables the introduction of mappings that may be a better fit for a given data set.

*How we implemented it*

Our focus was trying to get the best performance from the SVM using an optimal data set with optimal pre-processing functions applied<sup>11</sup>. We therefore came up with combinations of the various features we had access to, as well as their derivatives and differences from one another. The base features were: ambient temperature, measured temperature, the difference between ambient temperature and measured temperature, the difference between measured temperature and the derivative of measured temperature. We also had the same features for relative humidity. After creating combinations of these base features, we had a test set of 113<sup>12</sup> combinations as shown in **?@fig-feature-combos**.

[Combinations of Features for SVM ]{#fig-feature-combos}

---

<sup>8</sup>have been “tested” on a wide range of combinations of time series, and then sorted based on performance

<sup>9</sup>could also say logistic regression to reinforce the SVM logistic regression similarity

<sup>10</sup><https://developers.google.com/machine-learning/glossary#f>

<sup>11</sup>better way to say this?

<sup>12</sup>program this...

For each data set we collected, we created an SVM for each of the combinations in the test set, which resulted in 3x113 <sup>13</sup>. We chose not to perform hyperparameter tuning for the SVMs, and used the default parameters defined by sklearn <sup>14</sup> which consists of a radial basis function kernel. As we are interested in the developing an unsupervised detection method, we used the One-Class SVM implementation also from sklearn, which makes decisions by essentially using clustering to create the hyperplanes for classification. <sup>15</sup>

*How it compares to the previous method*

The SVM approach differs from the smoothing technique that we have developed, and is similar to other window techniques in that it is a highly generalizable method that is not developed with window detection in mind. A wide array of input features can be tested in order to get an acceptable prediction accuracy. However, this might not be acceptable in practice, as different scenarios/ window operation behaviors, will demand different sets of input features. In the event the true window detection pattern is unknown, it will be difficult to know which input features are giving a trustworthy result.

### 4.3 Evaluation Metrics

We used 3 sets of metrics to evaluate the methods we chose. The first two follow from a recent paper that compared the efficacy of different machine learning algorithms <sup>16</sup>, and the final set of metrics were developed to more closely examine specific window detection behavior.

*We chose to use this many diverse metrics because...*

#### 4.3.1 Standard Metrics

**Macro average F-1 score:** The F-1 score is a classic metric for evaluating the performance of classification algorithms. It condenses information about the precision of a model in predicting a certain class, as well as its ability to recall the available data. The macro averaged F-1 score averages the F-1 scores for all individual classes, but does not introduce weights in the averaging to reflect that the classes may be unbalanced. As shown in **tbl-data-collected**, the data sets we have vary from fairly balanced in experiment A, to highly unbalanced in experiment B. Therefore, using the macro average F-1 score provides a worst case performance. Using this F1-score also provides a basis of comparison to <sup>17</sup>. The implementation for computing this metric comes from sci kit learns classification report. <sup>18</sup>

---

<sup>13</sup>program this...

<sup>14</sup>need to cite this and include version..., and include the other default settings.

<sup>15</sup>need to double check this...

<sup>16</sup>dRdR

<sup>17</sup>dRdR

<sup>18</sup>include equation for F1 score

### 4.3.2 De Rautlin de Roy Metrics

The following metrics were introduced by <sup>19</sup>, and represent novel window detection specific algorithms.

**Opening Accuracy:** The opening accuracy gives a score to each opening instance predicted by the window detection algorithm. An opening instance is an interval on  $t$ , such as  $t_{nk} = t_n, t_{n+1}, \dots, t_k$ , where  $I(t_{nk}) = 1$ . An opening receives a score of 0 if it does not correspond to a true opening:  $I(t_{nk}) = 1, W(t_{nk}) = 0$ , and a score of 1 if it perfectly corresponds to a true opening:  $I(t_{nk}) = 1, W(t_{nk}) = 1$ . The score for a given predicted opening interval decreases by a penalty of 0.33 for each time step that does not align to a true opening. Therefore, a predicted opening interval  $I(t_{nk}) = 1$  that aligns a true opening interval can have a score of not less than 0 if it extends for more than two times <sup>20</sup> steps away from a true value. The value for each opening interval are averaged to get the opening accuracy for a given instance of a model. An unbounded metric is also computed when the penalty applied to a given interval's score does not stop at 0 but is allowed to further decrement. This provides an unbounded opening accuracy that gives a clearer indication of how "off" models are, given that many models we examining actually had a bounded opening accuracy of 0.

**True/False Opening Time:** This true opening time reflects the total amount of time when  $I(t) = W(t) = 1$ . The false opening time is when  $I(t) = 1, W(t) = 0$ . The metric is somewhat similar to the F-1 score in that rather than looking at specifically at change points, it provides information the entire time period.

### 4.3.3 Custom Metrics

We developed the final set of metrics based on the intuition that if a model is perfectly able to capture the times when a window state changes, then it is performing extremely well. We classify a "guess", as any value of  $I(t)$ , and an "action" as any value of  $W(t)$ . A "hit" occurs when  $I(t) = W(t)$  Like <sup>21</sup> we consider a prediction accurate if it is within at most 2 timestamps of the true occurrence. Therefore, a "near hit" occurs when  $I(t) = W(t \pm 1)$ . The metrics we examine are given below.

**Hits + Near Hits/ Guesses:** This metric accounts for the variability in guesses. A ratio of 1 indicates that all guesses taken by the model were accurate within two timestamps. A ratio close to 0 indicates that a lot of guesses were taken, but relatively few were close to where change occurred in the true window state.

**Guesses/Actions:** This metric implicitly accounts for the ability of the model to capture the pattern of the changes in window state. If the true window state changed only 10 times,

---

<sup>19</sup>dRdR

<sup>20</sup>check drdr math/logic here

<sup>21</sup>dRdR



but the model predicts changes on the order of 100, then the model is not performing well. A perfect score is 1.

## 5 Results

In Figure 3, the results comparing the smoothing technique<sup>22</sup> to the performance of the SVM across the metrics of interest are shown. While the smoothing technique represents one model, each of the SVM models shown in each of the subplots represents a model that performed unsupervised classification using different datasets<sup>23</sup>. The performance of the models is considered across the three empirical datasets that were collected.<sup>24</sup>

In Figure 3a<sup>25</sup>, We see that the EWM performs better than the best-performing SVMs on experiments A and B, but has considerably worse performance on experiment C. The range of f1-scores displayed in Figure 3a are comparable to those shown in<sup>26</sup>. The low values across all models for experiments B and C can be attributed to the rather unbalanced nature of the datasets, which a macro average f1-score penalizes.

The EWM has less accuracy in detecting when hits/guesses in experiments B and C compared to the best-performing SVMs, as shown in Figure 3b. However, as Figure 3c highlights, the SVMs take far more guesses per recorded action of window state change, while the EWM takes on the order of 1 guess per recorded action. This suggests that the SVMs' superior performance on the hits/guesses metric is due to probability, rather than an embedding<sup>27</sup> of the underlying physical dynamics.

Figure 3d and Figure 3e display the performance on the opening and unbounded opening accuracy metrics. We observe that the EWM has a far higher performance for experiment A, a slightly better performance for experiment C, and a score of 0 for experiment B. The unbounded opening accuracy allows us to compare just how "off" models are. We see that the EWM performs "less poorly" than the SVM across all experiments. This suggests that the EWM predictions of when window states change are, on average, closer in time to reality than the predictions of the SVM models.

The true time and false time metrics, shown in Figure 3f and Figure 3g operate as an analog to a weighted average f1-score, in that the the imbalance between window open and close is

---

<sup>22</sup>going to call this EWM from now on and fix when editing/ come up with a better name

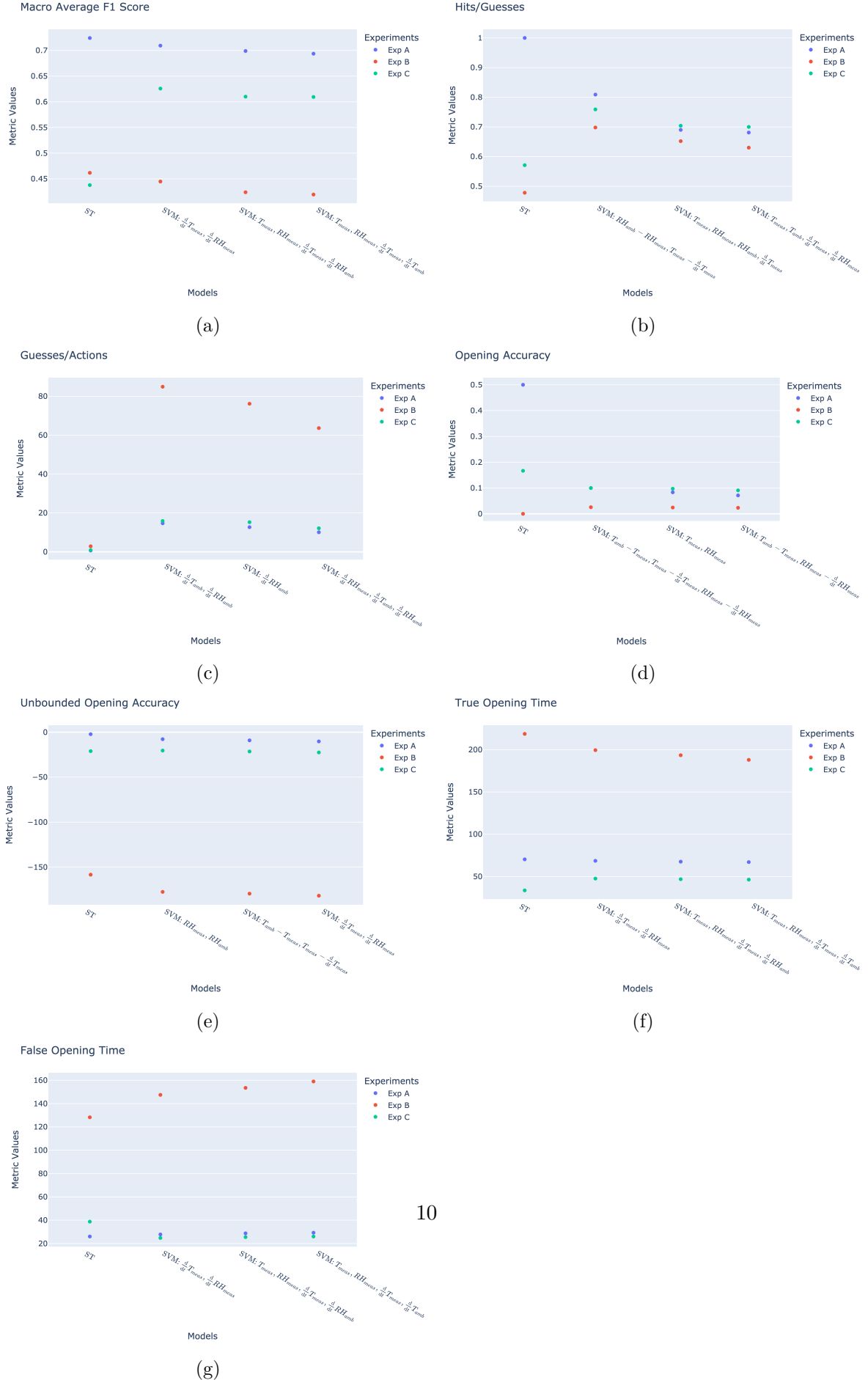
<sup>23</sup>TODO include tables underneath the graphs or annotate the graphs showing the best performing metrics..

<sup>24</sup>A is kind of like an favorable case, B is less ideal (less balanced window open/close), and C is even less ideal due to lack of strong differences in internal/external temperatures => include in methods

<sup>25</sup>need to have subfigure labels, figures a - g...

<sup>26</sup>drdr

<sup>27</sup>understanding?



inherently taken into account. The EWM is true more often and false less often for experiments A and B. The reverse is true for experiment C. <sup>28</sup>

Overall, EWM shows comparable performance to the SVMs across all metrics and experiments. This is significant, given that the SVMs shown in the graphs have made unsupervised classification decisions a wide range of combinations of time series, and then sorted based on performance to reveal the optimal time series. <sup>29</sup> The EWM method thus shows potential as a powerful technique with robust performance across both favorable and unfavorable datasets.

## 6 Conclusions

[Windows presentation](#)

---

<sup>28</sup>TODO: Some discussion here about poor performance on experiment C that can be attributed to the similarity between indoor and outdoor temperatures

<sup>29</sup>have been “tested” on a wide range of combinations of time series, and then sorted based on performance