# ENVIRONMENTAL RESEARCH
## INFRASTRUCTURE AND SUSTAINABILITY

# Is smart water meter temporal resolution a limiting factor to residential water end-use classification? A quantitative experimental analysis

1     # Is smart water meter temporal resolution a limiting factor to residential water
2     # end-use classification? A quantitative experimental analysis

3

4     Zahra Heydari[1], Andrea Cominola[2,3], and Ashlynn S. Stillwell[1]

5     [1]Civil and Environmental Engineering, University of Illinois Urbana-Champaign; [2]Chair of
6     Smart Water Networks, Technische Universität Berlin, Berlin, Germany; [3]Einstein Center
7     Digital Future, Berlin, Germany.

8

9     **Abstract**

10    Water monitoring in households provides occupants and utilities with key information to support
11    water conservation and efficiency in the residential sector. High costs, intrusiveness, and
12    practical complexity limit appliance-level monitoring via sub-meters on every water-consuming
13    end use in households. Non-intrusive machine learning methods have emerged as promising
14    techniques to analyze observed data collected by a single meter at the inlet of the house and
15    estimated the disaggregated contribution of each water end use. While fine temporal resolution
16    data allow for more accurate end-use disaggregation, there is an inevitable increase in the
17    amount of data that needs to be stored and analyzed. To explore this tradeoff and advance
18    previous studies based on synthetic data, we first collected 1-second resolution indoor water use
19    data from a residential single-point smart water metering system installed at a 4-person
20    household, as well as ground-truth end-use labels based on a water diary recorded over a 4-week
21    study period. Second, we trained a supervised machine learning model (random forest classifier)
22    to classify six water end use categories across different temporal resolutions and two different
23    model calibration scenarios. Finally, we evaluated the results based on three different
24    performance metrics (micro, weighted, and macro F1 scores). Our findings show that data
25    collected at 1- to 5-second intervals allow for better end-use classification (weighted F-score
26    higher than 0.85), particularly for toilet events; however, certain water end uses (e.g., shower and
27    washing machine events) can still be predicted with acceptable accuracy even at coarser
28    resolutions, up to 1 minute, provided that these end use categories are well represented in the
29    training dataset. Overall, our study provides insights for further water sustainability research and
30    widespread deployment of smart water meters.

31    Keywords: smart water meter, temporal resolution, residential water use, water sustainability,
32    supervised machine learning

33    **1. Introduction**

34    Strong emphasis on sustainability in water use has been increasingly brought to light by growing
35    population and urbanization (Cosgrove and Loucks 2015), coupled with climate change impacts
36    on water resources (Jabaloyes et al. 2018; Karamouz and Heydari 2020). With existing

37  limitations on water resource availability, new developments to increase water storage and
38  supply are often physically or economically constrained. Therefore, better management of
39  existing water resources has become an issue of paramount importance (Mazzoni et al. 2021).
40  Public utilities are now investing significant resources and efforts in the development and
41  implementation of water management strategies, both on the supply and the demand side, to
42  ensure future water security (Jain and Ormsbee 2002; Herrera et al. 2010). On the demand side,
43  these strategies include water saving technologies, new water policy regulations, rebate programs
44  for water-efficient devices, leakage management, and source substitution (e.g., replacing non-
45  potable end-uses with grey, recycled, or harvested rainwater (Dixon et al. 1999)) (Gleick et al.
46  2003; Inman and Jeffrey, 2006; Stewart et al. 2013; Cominola et al. 2015; Ntuli and Abu-
47  Mahfouz 2016).

48  Beside their direct effect on water resources, residential water conservation and efficiency
49  strategies can help save water-related energy required for water treatment, distribution, and
50  heating (Srinivasan et al. 2011). Residential end uses are responsible for more than 70% of all
51  water-related energy use (Escriva-Bou et al. 2018). However, the effectiveness of these measures
52  hinges on an accurate estimate of water demand from detailed understanding of how and when
53  water is used in the residential sector. Access to high resolution water consumption data can help
54  improve our knowledge of water demand, identify specific fixture/appliance end uses (e.g.,
55  toilet, shower, washing machine, outdoor irrigation), or detect anomalies, such as leaks (Luciani
56  et al. 2019). Smart water meters, which can provide the fine resolution data necessary to discern
57  end uses, have been proven essential in supporting water conservation and efficiency measures in
58  practice (Britton et al., 2008).

59  Conventional residential water meters typically collect coarse resolution data and require manual
60  readings, limiting the understanding of household-scale water use characteristics and its patterns
61  in time. Conversely, smart (or digital) water meters enable the collection and automated
62  reporting of fine resolution water use data, thereby allowing planners and utilities to better
63  understand demand patterns and enact management strategies. Smart metering can help the
64  development of accurate demand characterization and forecasts and, hence, improve the
65  operation and long-term planning of water supply and distribution systems (Stewart et al. 2018),
66  or promote durable conservation behaviors (Cominola et al., 2021). In addition, detailed
67  knowledge about water consumption at the household level can also translate into financial
68  savings for home occupants, especially when complemented with information about individual
69  end uses (e.g., Blokker et al. 2010).

70  Obtaining information on residential end uses is not a trivial problem. Information about
71  residential water demand at the end-use level could, in principle, be obtained through direct
72  measurements via intrusive monitoring, i.e., by installing sub-meters at all household end uses.
73  However, this approach is often practically or economically infeasible from a utility perspective
74  and would likely be rejected by home occupants due to its intrusive nature. Instead, water
75  utilities are increasingly installing residential smart water meters that collect fine resolution

76    water consumption data at the service line or entrance into the home, providing aggregate water

77    data, which are so far primarily used for billing purposes (Fogarty and Hudson 2006; Froehlich

78    et al. 2009). Similarly to previous experiences in the electricity sector, limits to directly

79    collecting water-use data at the residential end uses has motivated the development of several

80    non-intrusive disaggregation approaches, which have the advantage of allowing the

81    decomposition of a signal measured at the household level (i.e., aggregate water use) into the

82    individual contribution of each end use (Cominola et al. 2017; Di Mauro et al. 2020; Bethke et

83    al. 2021).

84    Several state-of-the-art disaggregation techniques require additional sensing on the premise

85    plumbing infrastructure and/or a manual characterization of each end use (Fogarty and Hudson

86    2006; Kim et al. 2008). These techniques can be intrusive, expensive, and time consuming, thus

87    they are not easy to develop or replicate at large scales (Froehlich et al. 2009, 2011; Srinivasan et

88    al. 2011; Ellert et al. 2015; Ntuli and Abu-Mahfouz). Other disaggregation techniques use only

89    flow (or volume) data collected at the household water inlet point. They can classify end uses in

90    a non-intrusive way, with the accuracy of results varying across different data sampling temporal

91    resolutions (e.g., 1-10 seconds vs. minutes; Clifford et al. 2018; Vitter and Webber 2018).

92    Understanding the tradeoff between the value of the information provided by fine-resolution data

93    and the economic and operational costs of the metering system is crucial to inform the design of

94    future metering networks and associated analytics to facilitate customer data interpretation.

95    The availability of fine-resolution smart meter data generates several opportunities for advancing

96    water demand management. Sub-minute sampling resolution is essential for most water end-use

97    disaggregation algorithms to provide a reliable categorization of household level water use into

98    different fixtures/appliances (e.g., shower, toilet, dishwasher, etc.) (Willis et al. 2010; Nguyen et

99    al. 2013; Abdallah and Rosenberg, 2014; Horsburgh et al. 2017; Cominola et al. 2018).

100   However, high resolution metering inevitably increases the amount of data the water utility must

101   collect, process, and manage. Sampling at 1-second resolution, for instance, implies replacing the

102   typical 12 monthly readings per user with over 31.5 million data readings. Large amounts of data

103   can compromise hardware and software performance due to issues with meter power sources,

104   battery life, telemetry network capacity, data gaps, and billing software, besides requiring

105   utilities to acquire new skill sets for their employees (Stewart et al., 2010; Suero et al. 2012).

106   Among the existing literature that has already explored the implication of data sampling

107   resolution on water end use disaggregation (e.g., Wonders et al. 2016), Cominola et al. (2018)

108   developed an analysis based on synthetic time series of water end use generated with STREaM,

109   the STochastic Residential water End-use Model. Their model relied on statistical distributions

110   of end-use characteristics derived from a large dataset of disaggregated water end-uses from over

111   300 single-family households in nine U.S. cities (DeOreo, 2011). STREaM generated synthetic

112   time series of water end uses with diverse sampling resolutions, which were analyzed with a

113   multi-resolution assessment framework to identify potentially critical thresholds in data

114   resolution for different aspects of information retrieval and demand management. While such

115   studies tend to make up for the shortness of (or limited access to) data through stochastic

3

116 modeling to generate synthetic disaggregated water use data, a data gap remains with limited
117 availability of ground-truth water end-use observations from real-world data (Di Mauro et al.
118 2020; Di Mauro et al. 2021).

119 Here, we address the challenge of testing if and how the theoretical results obtained in the
120 literature on synthetic data change when similar analysis is replicated directly on real-world data.
121 Compared to synthetic data, real-world data might be characterized by higher signal noise, data
122 gaps, and limited dataset size for model calibration. We build on the above modeling efforts
123 through collection and analysis of observed data from a monitored study home in the Midwest
124 United States, exploring the tradeoffs between data sampling resolution and performance in
125 water end-use classification. We examine different data sampling resolutions and explore water
126 end use disaggregation results by aggregating 1-second water consumption data from a 4-person
127 study household to coarser resolutions. We evaluate a set of performance metrics regarding water
128 end-use classification using supervised machine learning informed by ground-truth end-use
129 labels obtained from a water diary recorded over a 4-week study period. Findings from our
130 multi-resolution assessment can support further research and assist utilities in quantifying the
131 benefits associated with second-to-minute data sampling resolutions and the costs of adopting
132 and maintaining fine-resolution metering infrastructures.

133 The major contributions of this work include:

- 134 • Training and testing a water end-use classification model on real-world observation data
  135 obtained with a single-point smart meter for a 4-person household coupled with labels
  136 from a water diary.
- 137 • Quantifying the effects of temporal data sampling resolution on the performance of water
  138 end-use classification.
- 139 • Analyzing the tradeoff between end-use classification performance and data sampling
  140 resolution under two scenarios characterized by different model calibration strategies.

141 **2. Material and Methods**

142 **2.1. Metering setup, data collection, and temporal aggregation**

143 In this study, we used data from a single-point smart water metering system installed at a 4-
144 person, single-family, fully-detached residence in the Midwest United States, collecting 1-second
145 resolution flow rate data over a 4-week study period from September 3 to October 1, 2021.
146 Aggregate indoor household water use data were collected from a custom ally® electromagnetic
147 flow meter provided by Sensus, installed on the main water supply pipe into the residents' home.
148 In addition to measuring flow rate (gal/min), the meter also sensed temperature (K) and pressure
149 (psi) data at a 1-second resolution. Although these pressure and temperature data are useful to
150 water system operations, they are not as valuable to demand disaggregation due the large impact
151 the distribution system has on these variables. We validated this assumption through feature
152 analysis based on correlation and data visualization (see Figures S15-S18 in the Supporting
153 Information). Consequently, we focused our analysis on flow rate data. The water meter writes

4

154 data to a computer running a script that parses the raw data into a suitable format for further
155 analysis. A data acquisition system connected to the water meter parsed the raw data into a
156 timestamped comma separated value (csv) format for further analysis.

157 To examine the effects of data sampling temporal resolution on water end-use classification, we
158 aggregated the 1-second resolution time series to resolutions of 5 seconds, 10 seconds, 30
159 seconds, and 1 minute. The 1-minute resolution has been recognized as a critical threshold for
160 certain end-use data analytics in the electricity sector (Armel et al. 2013). Similarly, a previous
161 study based on analysis of synthetic data identified the same threshold as critical for end-use
162 disaggregation in the water sector (Cominola et al. 2018). Here, we test these findings with an
163 experimental study based on real-world data and aim to identify a similar critical data sampling
164 resolution threshold for water end-use classification in the residential sector. Meanwhile, since
165 the study is only based on a 4-person household, we preliminarily compare water consumption
166 patterns with a larger study to ensure the study home is representative of larger scale behavioral
167 patterns.

168 During the study period, the home occupants manually recorded a water diary of labeled end
169 uses. In this study, six types of indoor water end uses contributed to the total household water
170 demand: faucets, toilets, showers, refrigerator, dishwasher, and washing machine. We used a
171 written water diary over the 4-week study period to collect ground truth end use data for model
172 training and validation. The 4-week period was selected based on previous studies and
173 practicality (Beal et al. 2011; DeOreo et al. 2016; Horsburgh et al. 2017). The water diary
174 included end use labels, start time, and date that were completed by the household occupants.
175 More details about the diary are reported in the Supporting Information, including the water
176 diary template (Figure S19) and an example of completed recordings (Figure S20). This data
177 collection included only factual data such that this work was determined not to meet the
178 definition of human subjects research and, therefore, did not require Institutional Review Board
179 (IRB) approval. Documentation of this IRB decision is available upon request. Limitations that
180 naturally arose during the water diary process were as follows:

181 • Events that occupants would forget to fill in the diary could not be labeled after the
182 disaggregation of the data.
183 • Start times listed in the diary would sometimes correspond to events that occurred 1-2
184 minutes before the reported time, implying that occupants would sometimes fill in the
185 diary after the event.
186 • Specifically for faucet events, occupants mentioned occasionally leaving the faucet on to
187 avoid reporting multiple events, resulting in long faucet durations that can represent
188 atypical behavior in the model training process.
189 • The water diary was completed manually and was unreadable for some events.
190 • Some reported events did not match the information received from the meter.

191 In addition to these limitations, a power outage created a 2-day data gap in the smart water meter
192 dataset, where the water diary was completed but measured water flow was missing.

## 2.2. End-use disaggregation

193

194 The end-use disaggregation step separates concurrent water use events along with single events,
195 that, aggregated on the axis of time, would give the original time series collected at the single-
196 point residential water meter. While end-use disaggregation and end-use classification
197 sometimes coalesce into one concept in literature, in this study we consider disaggregation as the
198 first step of the end-use classification process (Nguyen et al. 2018). Single events are defined as
199 those that occur in isolation (e.g., dishwasher only), while combined or concurrent events have
200 simultaneous occurrences of water usage (e.g., a toilet flush during a shower). A single $i$-th water
201 use event $E_i$ can be quantitatively characterized by a vector of features $F_i$, which include values
202 of, e.g., start time, end time, average flow, and volume of that event. Separating and identifying
203 overlapping, or concurrent, water use events is a significant challenge in residential water
204 studies, and the accuracy of existing smart meter disaggregation models decreases significantly
205 when these types of events are encountered (Cominola et al. 2015). Concurrent events occur
206 often, especially during longer duration events such as showers or outdoor irrigation. Thus,
207 disaggregating concurrent events from one another by leveraging information on the
208 characteristics of individual fixtures or by learning the patterns of individual end uses is essential
209 for the purpose of creating a comprehensive water profile for the household.

210 In this analysis, we used the disaggregation model from Bethke et al. (2021), developed based on
211 Nguyen et al.'s (2013) method of separating concurrent events by calculating the vector
212 gradients of the flow rate data to identify start and end times of overlapping events. Once we
213 separated events with the above disaggregation approach, we manually labeled each
214 appliance/fixture water event based on the water diary and examined the events further with the
215 classification model described below. We repeated this process for every resampled resolution as
216 well as the original 1-second data. At coarser resolutions, the performance of the disaggregation
217 model deteriorated when detecting multiple short duration events happening simultaneously
218 (e.g., hand washing), or short duration events happening on top of a long duration event.
219 Therefore, in addition to naturally having fewer observations at coarser resolutions, the number
220 of events that we were able to match with the diary also decreased (Figure S21).

## 2.3. End-use classification

221

222 After disaggregating the original water use time series, we labeled each event by matching with
223 the water diary. We then trained a random forest (RF) classifier to perform appliance/fixture end-
224 uses classification, using the disaggregated water events resulting from the previous step of end-
225 use disaggregation. The classification algorithm allocated each data point (i.e., a $i$-th water use
226 event $E_i$) in the dataset to one of the labeled classes, after training on tuples of events and
227 associated features ($E_i, F_i$).

228 RF models have been presented by Breiman (2001) as classical ensemble learning algorithms and
229 have shown to be outstanding predictive models in classification tasks (Herrera et al. 2010;
230 James et al.2013). Random forests are built using the same fundamental principles as decision

6

231    trees and bagging (Bootstrap Aggregation). Bagging introduces randomness into the tree
232    building process by building many trees on random subsets of the training data with replacement;
233    this process is also known as bootstrapping. Bagging then aggregates the predictions across all
234    the trees, which reduces the variance of the overall procedure and improves predictive
235    performance (Géron 2019). However, bagging trees could result in tree correlation that limits the
236    effect of variance reduction. Random forests help reduce variance by injecting more randomness
237    into the training process (Hastie et al. 2009). The random forest algorithm is a bagging algorithm
238    that draws random bootstrap samples from the training set. However, while bagging provides
239    each tree with the full set of features, random forests have a random feature selection that makes
240    trees more independent of each other, which often results in better variance-bias tradeoffs (Table
241    S1) (Friedman et al. 2001; Probst et al. 2019). In this study, the two features of average flow and
242    duration were eventually selected to build the final models, based on the results of our feature
243    importance analysis (Figure S22). Therefore, the search for the split variable was limited to a
244    random subset of the two chosen features. Feature importance was performed based on
245    permutation-based feature importance (Breiman 2001) by evaluating which features contributed
246    the most to the generalization power of the model.

247    To understand the mechanism used by RF models, it is necessary to understand the construction
248    of classification decision trees. The goal of such a tree is to partition data into small and
249    homogeneous groups. When travelling down the tree, data are split into possible responses called
250    nodes that symbolize the branches of a tree. To perform each partitioning operation, a decision is
251    based on an index (e.g., the Gini index), which allows RF models to partition the nodes of each
252    tree into more homogenous groups that contain a larger proportion of one class in each
253    subsequent node (Kuhn & Johnson, 2013). The Gini index is calculated as in Eq. 1, where $C$ is
254    the total number of classes in the model and $p_{nk}$ is the probability of the occurrence of class $k$ at
255    node $n$. In this study, six different classes were evaluated based on typical household end uses:
256    faucets (f), toilets (t), showers (s), refrigerator (r), dishwasher (d), and washing machine (w). The
257    sum of all probabilities at a certain node is equal to one (see Eq. 2):

258
$$G = \sum_{k \epsilon C} p_{nk}(1 - p_{nk}) \qquad (Eq.1)$$

259
$$p_t + p_s + p_f + p_r + p_w + p_d = 1 \qquad (Eq.2)$$

260    Other metrics similar to the Gini index can be used to build decision trees, including cross
261    entropy and misclassification error. However, the Gini coefficient is the most commonly used
262    metric in the literature (James et al. 2013). Moreover, according to Raileanu and Stoffel (2004),
263    the frequency of disagreement of the Gini index and entropy is only 2% of all cases, yet entropy
264    is generally slower to compute because it requires a logarithmic function. For the above reasons,
265    we used the Gini coefficient in this study.

266    Besides Gini, the RF algorithm involves several other hyperparameters that can be tuned to
267    optimize model performance. While studies have shown that RF models are less sensitive

7

268 towards tuning than other algorithms such as support vector machines (Probst et al. 2019),
269 modest performance gains can still be valuable considering the limitations that naturally come
270 with a small dataset. Using grid search, we gave ranges to RF hyperparameters to exhaustively
271 try all possible combinations and select the best hyperparameter combination. Minimum sample
272 at each leaf (2- 5), minimum sample split (2, 5, 8, 12) number of sub-features (1, 2), maximum
273 depth (3-10), and the number of trees (10, 20, 50, 100, 200) were initially given to the grid for
274 hyperparameter tuning.

### 2.4. Model calibration and data sampling resolution scenarios

276 We considered two scenarios for calibration analysis of the classification model: the "1-second
277 only calibration" (Scenario 1), and "multi-resolution calibration" (Scenario 2).

278 The "1-second only calibration" (Scenario 1): In this scenario, the RF model was trained only on
279 the measured data at the 1-second resolution. Extended time series of 1-second resolution water
280 use data are not usually available from utility records, but they can be collected in small-scale
281 customized and experimental smart meter installations. With this scenario, we test whether
282 investing efforts and resources in gathering a small model calibration dataset at sub-minute
283 resolution is worth the potential gain of model disaggregation accuracy at coarser resolutions.
284 Our assumption behind this scenario is that the features of water use events can be more
285 accurately learned from data collected at higher resolutions. In the 1-second trained RF model
286 scenario, we split the labeled data into train (70% of the data) and validation (30% of the data)
287 datasets. The validation set was used to assess the model performance on the 1-second trained
288 data. Then, the entire resampled dataset from all other resolutions were separately used as test
289 sets to compare the performance of the model on coarser resolutions.

290 The "multi-resolution calibration" (Scenario 2): In this scenario, we trained different RF models
291 for each resolution (5 seconds, 10 seconds, 30 seconds, and 1 minute) on their own dataset and
292 compare their performances both with one another and with Scenario 1. In this scenario, we
293 examine the value of retraining the RF model specifically for different temporal resolutions to
294 quantify differences in performance between sampling resolution and, comparatively with
295 Scenario 1, across different model training strategies. To retain the value of limited data and
296 improve generalizability of the models, we implemented a k-fold cross-validation strategy
297 (Hawkins et al. 2003). We thus split the training set into k subsets, called folds, and then
298 iteratively fit the model k times, each time training the data on k-1 folds and evaluating on the
299 remaining single fold (representing the validation data). In this study, we fit the model with k =
300 10. At the end of training, we averaged the performance across all validation folds as the final
301 performance metric for the model.

### 2.5. Performance metrics

303 RF is a noise robust technique. However, when considering imbalanced problems,
304 canonical machine learning algorithms generally tend to be biased towards the majority group.

8

305 This behavior happens because such algorithms consider the number of objects in each group to
306 be roughly similar (Krawczyk, 2016, Ribeiro and Reynoso-Meza 2020). However, the minority
307 class is often the most important when dealing with skewed distributions, and a performance
308 metric should be chosen in a way to overcome such bias. While we do not directly balance the
309 dataset used in this study because of its limited size, in this analysis we evaluate and compare the
310 model performance using different formulations of the F1-score (FS). Specifically, we compare
311 (i) micro-FS, which is a global metric attributing equal importance to each sample, thus giving
312 emphasis on common labels, (ii) macro-FS, which attributes equal importance to each class, and
313 (iii) weighted FS, which computes the weighted average of the FS values obtained for individual
314 classes. While using these metrics does not solve class imbalance, we examine different F-score
315 formulations to see whether our classifier gets biased towards well represented classes or not.

316 Micro-FS (usually referred to as simply FS) is a global performance metric that puts more
317 emphasis on the most represented labels in the data set since it gives each sample the same
318 importance. Labels that are underrepresented in the dataset may not be intended to influence the
319 overall micro-FS heavily if the model is performing well on the other more common classes.
320 Micro-FS (Eq. 3) is defined as the harmonic mean of the precision (Eq. 4) and recall (Eq. 5):

321
$$Micro\text{-}FS = 2 \times \frac{(precision \times recall)}{(precision + recall)} \quad (Eq.3)$$

322
$$Precision = \frac{TP}{TP + FP} \quad (Eq.4)$$

323
$$Recall = \frac{TP}{TP + FN} \quad (Eq.5)$$

324

325 where true positives (TP) are the number of correctly classified positive instances, false positives
326 (FP) are the number of negative instances incorrectly classified as positive, and false negatives
327 (FN) are the number of positive instances incorrectly classified as negative.

328 Macro-FS (short for macro-averaged F1 score) is used to assess the quality of classification in
329 problems with multiple classes. The macro-FS gives the same importance to each class, with low
330 values for models that only perform well on the common classes while performing poorly on the
331 classes with less data. The macro-FS is defined as the mean of class-wise FS in Eq. 6:

332
$$Macro\text{-}FS = \frac{1}{N}\sum_{i=1}^{N} FS_i \quad (Eq.6)$$

333 where $i$ is the class index and $N$ is the number of classes/labels.

334 The weighted-average FS (Eq. 7) is calculated by taking the mean of all per-class F1
335 scores while considering the number of actual occurrences of each class in the dataset.

9

336 $$Weighted\text{-}FS = \frac{1}{H}\sum_{i=1}^{N}|i| \times FS_i \ (Eq.7)$$

337 where $i$ and $N$ are as above, and $H$ is the total number of aggregated elements across all classes
338 (Cominola et al. 2021).

339 The weighted-FS formulation modifies the macro-FS to account for class imbalance, while
340 imbalance is not considered in micro-FS and macro-FS.

341 **3. Results and discussion**

342 **3.1 Data characterization - time of the day visualization**

343 To make sure our study home could be a proper representative of a larger study scale, we
344 initially visualized the time-of-day and day-of-week distribution of three major classes of events
345 (shower, washing machine, and dishwasher) to find regular patterns of consumption similar to
346 those displayed in larger datasets. Much of the occupants' water consumption occurs during
347 typical weekday mornings and evenings. Figure 1(a) depicts shower end use distribution
348 throughout the week and time of the day in our study home. The results show that showers have
349 a more sporadic pattern of use on weekends while during weekdays most of them occur during
350 regular morning and evening peak hours. These behavioral patterns align with the time-of-day
351 and day-of-week distribution of showers reported in an analysis of water end use data gathered
352 for 762 U.S. households (Cominola et al. 2020), shown in Figure 1(b). The time-of-day and day-
353 of-week distribution figures for the washing machine (Figure S1) and dishwasher (Figure S2) are
354 also shown in the Supporting Information, with similar results. Washing machine events are
355 observed mostly during weekends with a wide distribution throughout time of the day, while
356 dishwashers are typically used throughout the week, either mornings or evenings. Comparison of
357 the results show similar patterns between our study home and the larger study of U.S. households
358 used in Cominola et al. (2020), demonstrating the potentially transferrable nature of our study
359 home results. Similar widespread end-use data would help water planners and managers
360 understand water consumption patterns, consumer behavior, and temporal variability. Decreasing
361 consumption during peak time on a widespread scale could contribute to lowering overall peak
362 demand for the local utility and reduce pressure on existing water infrastructure.
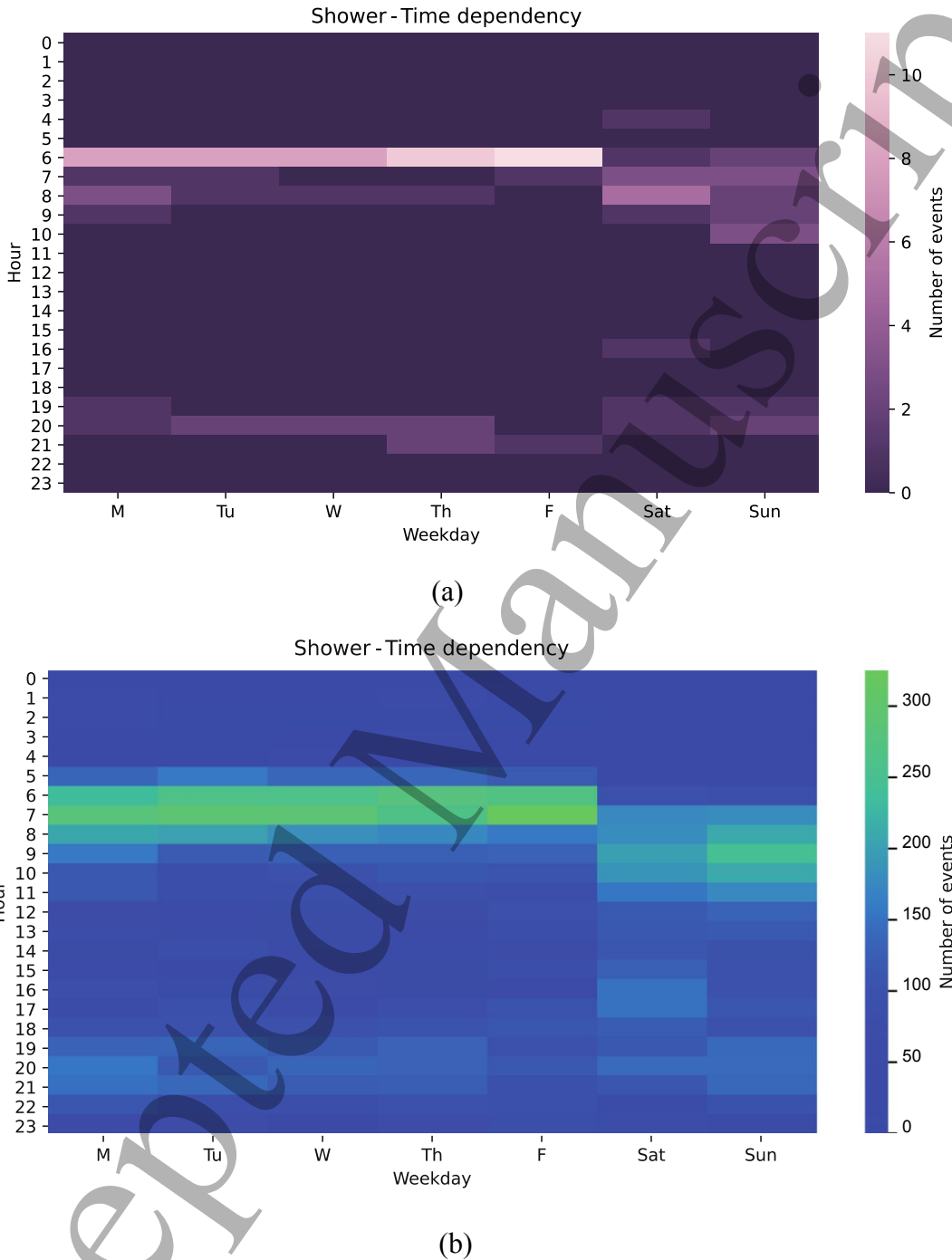
363

(a)



(b)

Figure 1. Time-of-day and day-of-week analysis: (a) Results from shower end use in this study home, 4-week study with 1-second resolution data; (b) Results adapted (with permission) from Cominola et al. (2020) from shower end-uses in 762 U.S. homes, 2-week study period with 10-second resolution data.

## 3.2. Comparative multi-resolution scenario analysis

The overall RF model performance across different resolutions in both calibration scenarios is presented in Figure 2. Grey lines represent Scenario 1 (1-second only calibration) and blue lines

11

371 represent Scenario 2 (multi-resolution calibration). The micro-FS, weighted-FS, and macro-FS
372 are represented with dashed, solid, and dotted lines, respectively. We observe that Scenario 2
373 gives higher performance across different temporal resolutions regardless of the performance
374 metric. For both 1-second and 5-second resolutions, the micro-FS and weighted-FS values are
375 similar: 0.91 and 0.89 for the micro- and weighted-FSs, respectively, at the 1-second resolution,
376 and 0.87 and 0.85 for the micro- and weighted-FSs, respectively, at the 5-second resolution. The
377 macro-FS generally shows the lowest values for all resolutions for both scenarios. We observe a
378 mild decrease in performance metrics with coarser temporal resolutions in Scenario 2, while
379 performance metrics decrease notably for resolutions coarser than 5 seconds in Scenario 1,
380 dropping as low as 0.2 for the 1-minute resolution.
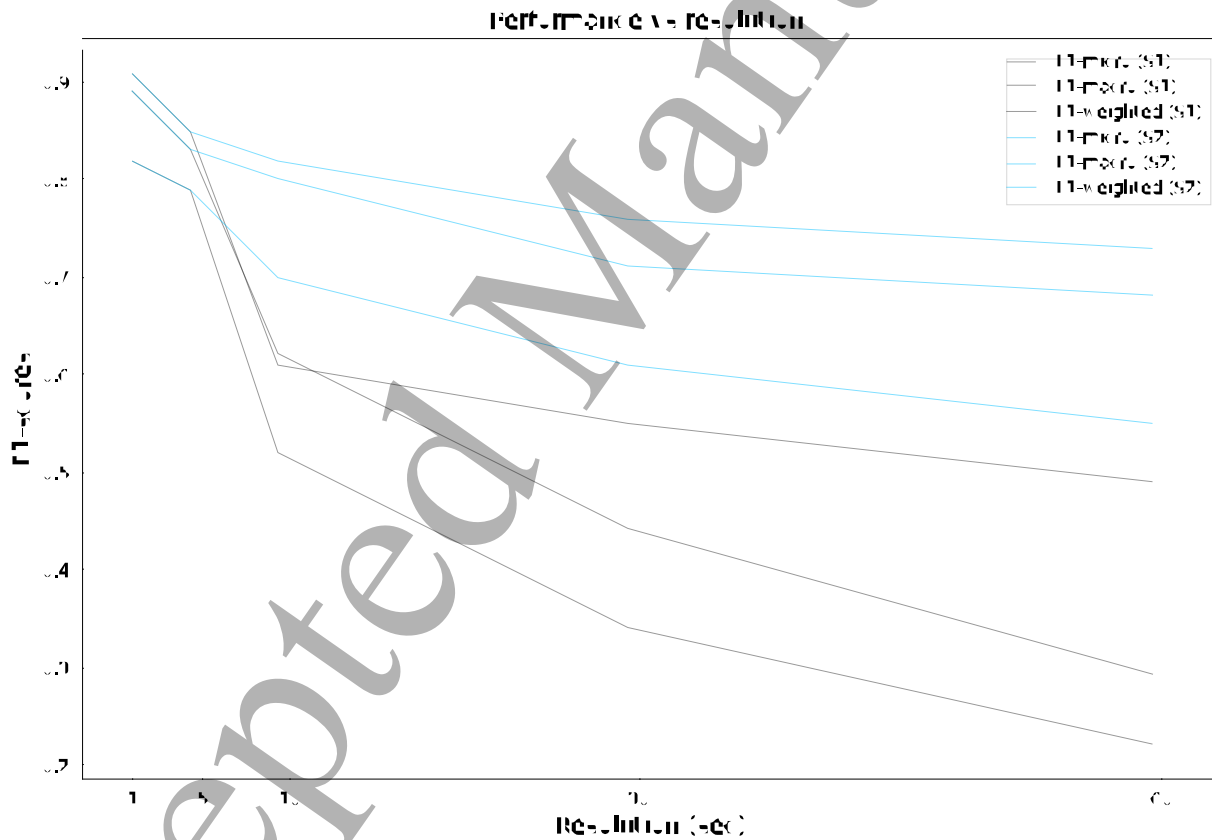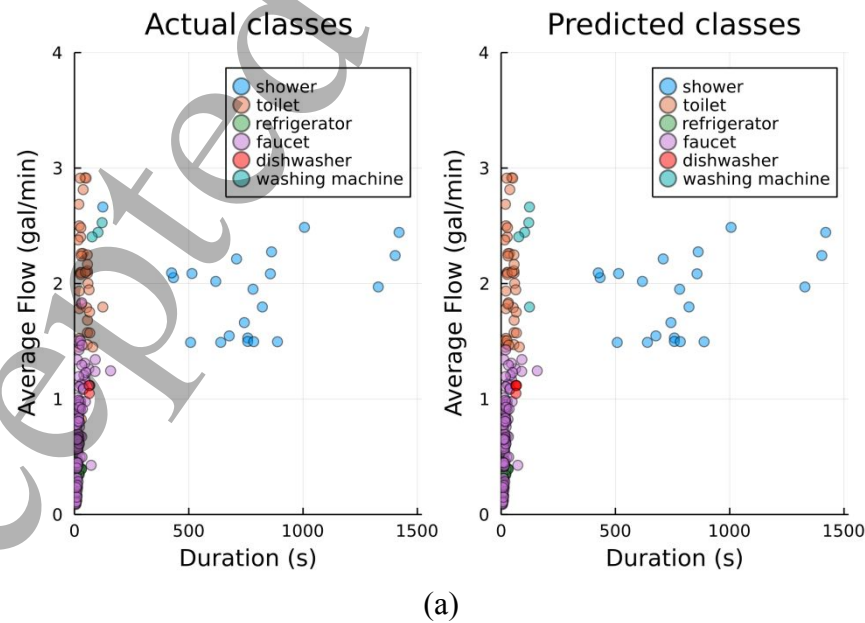
381

382



383

384 Figure 2. F1-score vs resolution curves for different F1-score formulations for Scenario 1 (grey lines) and
385 Scenario 2 (blue lines). The micro-FS (dashed lines), weighted-FS (solid lines), and macro-FS (dotted
386 lines) are represented.

387 Overall, our results indicate that the RF models learned end use event features better when
388 trained at the same data sampling resolution that they are tasked to use to classify unseen events,
389 provided that a training dataset with labelled events at that resolution is available. If

12

390  classification models are trained for application on data measured at the same resolution

391  (Scenario 2), those models can perform at an acceptable level of performance even at coarser

392  resolutions, depending on the relative importance of different end use classes. This observation

393  has important implications related to the tradeoffs between fine-resolution data collection and

394  increased data analytics needs. For instance, if a utility wants an estimate of water consumption

395  by the main indoor water uses in households (e.g., toilets and showers), the 1-minute resolution

396  model still provides an acceptable performance (weighted-FS equal to 0.73). This performance is

397  lower than the FS of 0.89 obtained for the 1-second resolution model, but this loss in model

398  accuracy is balanced by the benefit of gathering, storing, and analyzing fewer data observations

399  at the coarser temporal resolution. Conversely, if detailed information on all end uses is required,

400  only the 1-second and 5-second resolutions provide high performance predictions on all end use

401  classes; for less represented end uses, performance is compromised at coarser resolutions.

### 3.3. Detailed end-use classification results

403  Our detailed RF model validation results are presented in Figure 3, where the predicted classes

404  (right) are compared to the actual classes (left). Figure 3(a) represents the entire 1-second

405  resolution set of events, while Figure 3(b) zooms in on shorter duration events for clarity. The

406  average flow rate (gal/min) and duration (s) were used as identifying features for our model. Of

407  the total 654 events labeled, we used 196 events as a validation set. The model predicts the test

408  set with an accuracy of 92% and a weighted-FS of 0.89, which is noteworthy given the fact that

409  the training dataset had limited observations in some classes such as dishwasher and washing

410  machine. The model correctly predicts 179 events out of 196 total events of the test set.
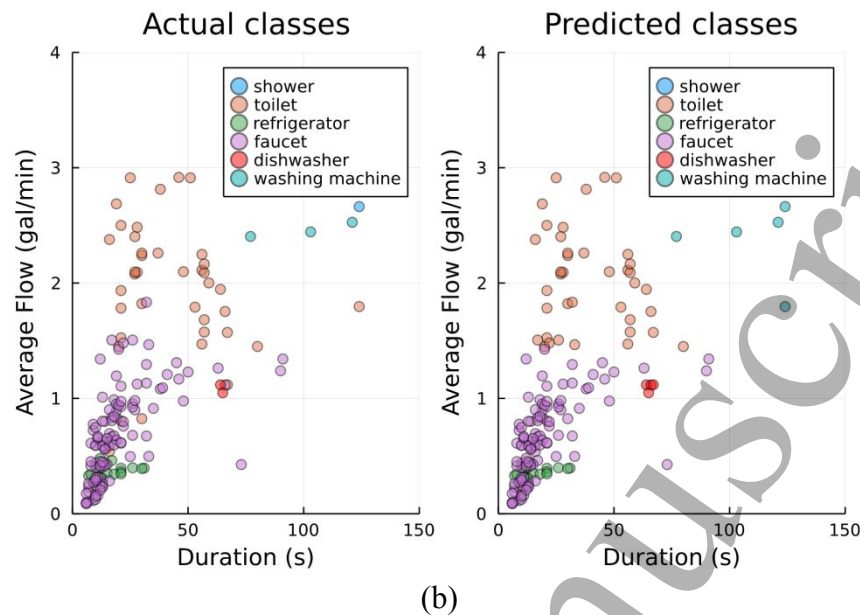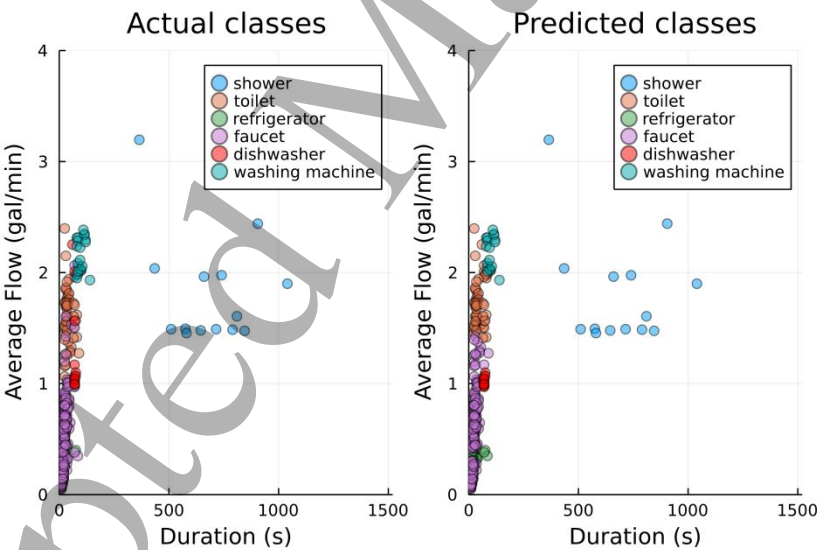
411



(a)

(b)

Figure 3. Actual and predicted water end-use classes. Predicted classes are obtained as results of the RF classifier on the 1-second resolution test set: (a) shows the entire dataset with durations ranging from 1-1500 seconds, and (b) shows the same results zoomed in on events within a duration range of 150 seconds (excluding 23 shower events) for clarity.

412   Yet, the high model performance in all classes might overrepresent the overall ability of our RF
413   models to classify unseen end use events. Our results might imply that, due to the fine temporal
414   resolution of the data, the model discerns the constant range of duration and average flow of
415   those end uses with automatic water consumption cycles (e.g., washing machine, dishwasher)
416   and detects them correctly. However, since our study represents a single household only, the
417   model might be overfitting on data from automatic appliances due to the invariance of duration
418   and flow in these specific automatic appliances, thus results on these specific end uses may not
419   be generalizable.

420   It is important to note that, while individual toilet uses are typically homogeneous in terms of
421   water consumption volume and duration, even considering dual-flush systems, the combination
422   of toilet and bathroom faucet uses are difficult to detect and disaggregate because such uses are
423   often almost simultaneous (e.g., use of toilet and consequent handwashing in a same minute).
424   Although temporal resolutions finer than 1 minute reduce disaggregation errors (Mazzoni et al.
425   2021), we were not able to disaggregate all toilet events followed by faucets. Rather, we labeled
426   the mentioned events as toilets since we attributed the subsequent faucet use due to the toilet use.
427   As a result, toilets have a wider range of flow and duration, as shown in Figure 3.

428   Figure 4 shows the classification results for Scenario 1 (1-second only calibration) applied to the
429   resampled 5-second (Figure 4(a)) and 1-minute resolutions (Figure 4(b)), respectively, selected
430   as examples at the two extremes of the considered spectrum of data resolutions. We report our
431   analysis results in both U.S. customary units (gal/min) and SI units (L/min). In comparing

14

432  different temporal resolutions, coarser resolutions tend to compress data points on the vertical

433  axis (i.e., decrease average event flow) and extend their range on the horizontal axis (i.e.,

434  increase event duration) due to temporal averaging. For example, toilet events that originally

435  ranged from 1.7-3 gal/min (6.4-11.4 L/min) average flow in the measured 1-second resolution

436  tend to shift to 1-2.5 gal/min (3.8-9.5 L/min) in the 5-second resolution and decrease further to

437  0.4-0.8 gal/min (1.5-3 L/min) in the 1-minute resolution. The duration of events increases with

438  coarser temporal resolution to an extent that the total volume of events is the same as to the

439  volume in the original 1-second resolution measurements. The mentioned shifts in values of end-

440  use features leads to decreased model performance with coarser temporal resolutions, up to a

441  point where, as shown in Figure 4(b), the model can no longer detect any toilet events. The

442  model still correctly predicts showers and a few washing machine events at the 1-minute

443  resolution; however, the model application to the 1-minute data predicts most other end uses as a

444  faucet under Scenario 1. Similar Scenario 1 classification results for the 10- and 30-second

445  resolutions are presented in the Supporting Information (Figures S3-S4) along with the zoomed

446  in figures of the 5-second and 1-minute resolutions for a detailed view (Figures S5-S8).
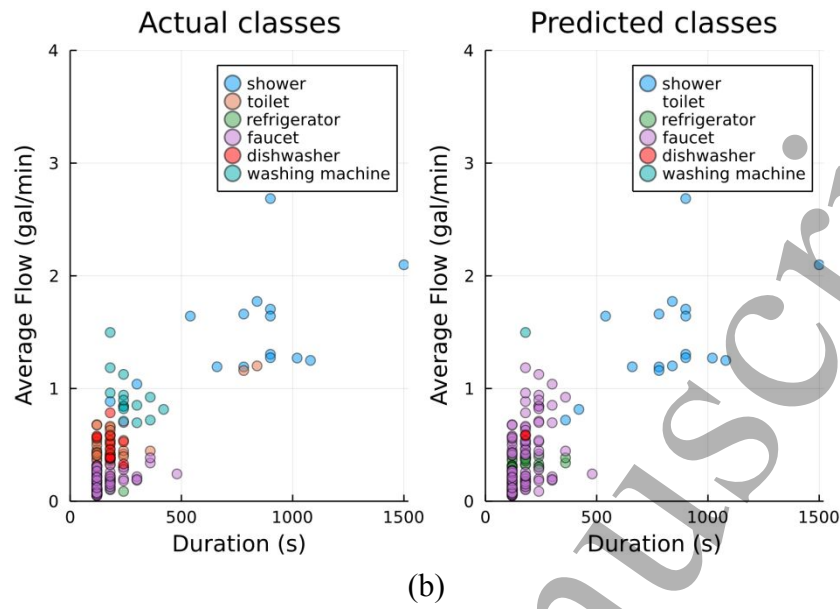
447

448



(a)

(b)

Figure 4. Actual and predicted water end-use classes. Predicted classes are obtained as results of the RF classifier trained on the 1-second data (Scenario 1) and applied to the (a) 5-second resolution test set, and (b) 1-minute resolution test set.

Figure 5 shows the confusion matrices of water end use classification across the events of our 4-person study household for Scenario 1. Faucets (f) account for the most frequent end uses, followed by toilets (t). The matrices show the total number of events labeled for each resolution, the actual classes, and the predicted classes by the model. The results for the 5-second resolution show that of 382 total events that we were able to match with the water diary, 324 events were classified correctly (Figure 5(a)). The main misclassifications were in predicting 14 actual toilet end uses as faucets and 4 actual faucet end uses as toilets. This misdetection mostly occurs for data that fall in the area with average flows of 1-1.5 gal/min (3.8-5.7 L/min) and durations of 25-50 seconds (see Figure S5 in the Supporting Information). For the 1-minute resolution (Figure 5(b)), only 187 events had corresponding end uses in the water diary due to disaggregation errors where the model was not able to separate concurrent events because of loss of information that naturally accompanies coarser resolutions. Out of these 187 events, 92 were classified correctly. The classification model predicts 135 events as faucets. While only 73 of these events are actually faucets, they still account for 40% of the prediction accuracy, motivating consideration of F1-score metrics due to the imbalanced dataset.

16

**Confusion Matrix (a)** — 5-second resolution

| Actual \ Predicted | d | f | r | s | t | w |
|---|---|---|---|---|---|---|
| d | 9 | 1 | 0 | 0 | 3 | 0 |
| f | 5 | 216 | 13 | 0 | 4 | 0 |
| r | 0 | 9 | 25 | 0 | 0 | 0 |
| s | 0 | 0 | 0 | 17 | 3 | 0 |
| t | 0 | 14 | 0 | 0 | 43 | 0 |
| w | 0 | 0 | 0 | 0 | 6 | 14 |

**Confusion Matrix (b)** — 1-minute resolution

| Actual \ Predicted | d | f | r | s | t | w |
|---|---|---|---|---|---|---|
| d | 1 | 6 | 4 | 0 | 0 | 0 |
| f | 0 | 73 | 9 | 0 | 0 | 0 |
| r | 0 | 15 | 0 | 0 | 0 | 0 |
| s | 0 | 1 | 0 | 17 | 0 | 0 |
| t | 2 | 31 | 11 | 2 | 0 | 0 |
| w | 2 | 9 | 0 | 3 | 0 | 1 |

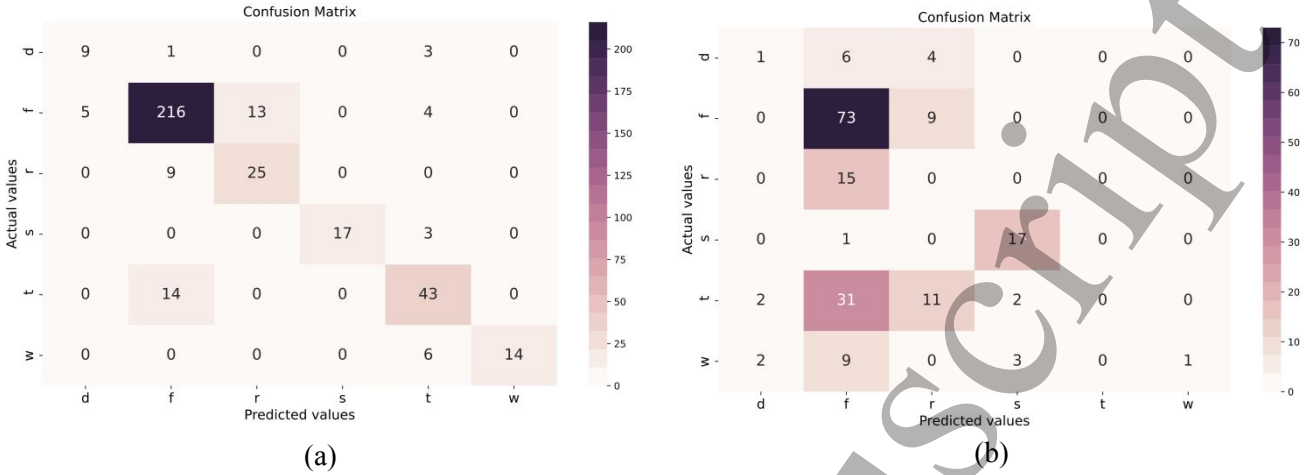(a)                                                      (b)

Figure 5. Confusion matrices for Scenario 1 (1-second trained random forest model): (a) 5-second resolution with 382 total events; (b) 1-minute resolution with 187 total events. Matrix rows show the actual classes and columns show the predicted classes for the following end uses: w (washing machine), s (shower), f (faucet), t (toilet), r (refrigerator), and d (dishwasher). Cell color is proportional to the number of events in that cell.

Figure 6 shows the confusion matrices of end use water consumption for Scenario 2, the multi-resolution calibration, for the 5-second (Figure 6(a)) and 1-minute (Figure 6(b)) resolutions. These results illustrate how the percentage of correct predictions changes from the 5-second resolution to the 1-minute resolution. Compared to Scenario 1 (Figure 5(a)), the 5-second resolution performance in Scenario 2 has either slightly improved or stayed the same, with the exception of refrigerator events (r). The 1-second resolution trained model in Scenario 1 had a better performance in predicting 5-second resolution refrigerator events. The prediction of toilets improved notably from 43 out of 57 to 51 out of 57 events. The main misclassifications were in predicting 5 actual toilet end uses as faucets and 6 actual faucet end uses as toilets. This misdetection mostly occurs for data that fall in the area with average flows of 1-1.5 gal/min (3.8-5.7 L/min) and durations of 25-50 seconds (see Figure S5 in the Supporting Information). Overall, the 5-second resolution has a high performance under both scenarios, with performance metrics slightly less than those of the 1-second resolution (as shown in Figure 2). In the 1-minute resolution, our model correctly predicts 139 of 187 labeled events, having the highest prediction accuracy in washing machine (100%), faucet (90%), and shower (88%) events. These results imply that if any of the aforementioned end uses are of importance, the 1-minute resolution can still be informative.

With further investigation of the diagonals of the confusion matrices, we see how Figure 6(b) has ameliorated in comparison to Figure 5(b), increasing correct predictions from 93 to 139. The 1-minute resolution model is still not able to discern refrigerator faucet events (r) from tap faucet events (f); however, this misclassification is not a critical issue since the refrigerator faucet is a faucet in nature. A noteworthy observation is that although the 1-minute resolution model under Scenario 2 incorrectly classifies one shower and one actual faucet event as a washing machine (i.e., false positive, $FP$ in Eq. 4), it does not label any other actual washing machine event as

496  other events (i.e., false negative, *FN* in Eq. 5), which leads to a higher Recall in this specific
497  class (100%) than in the 1-second resolution model (86%) (refer to Figure S12) and the 5-second
498  resolution model (95%), with the tradeoff of lower Precision (88% versus 95% in the 1-minute
499  and 5-second resolutions, respectively). Additional confusion matrices at other temporal
500  resolutions are available in Figures S9-S14 of the Supporting Information.

501



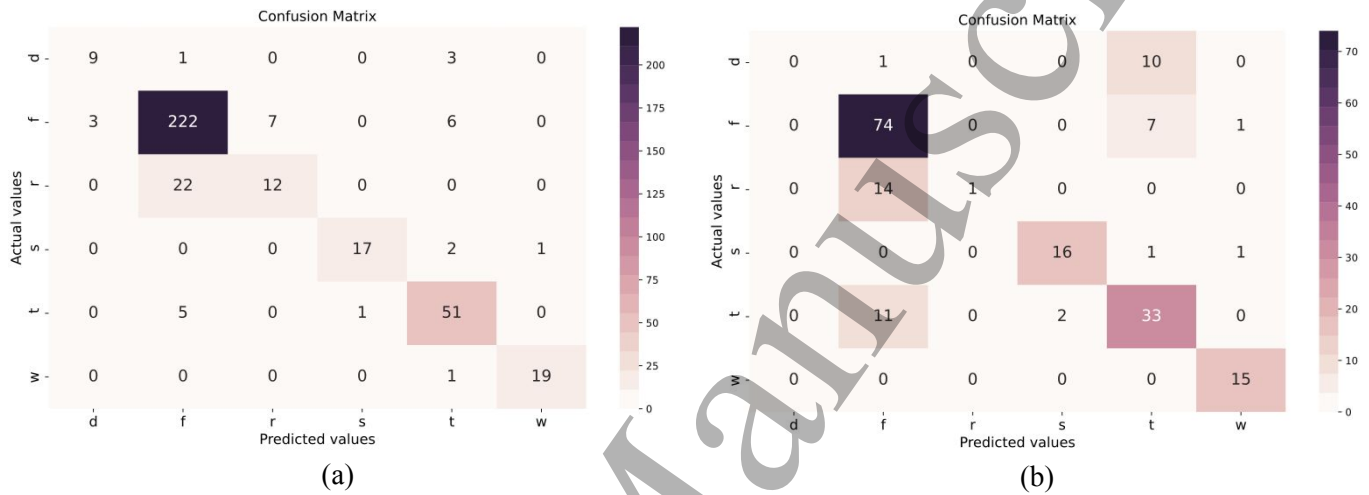(a)                                                                                     (b)

502   Figure 6. Confusion matrices for Scenario 2 (random forest model trained at the resampled temporal
503   resolutions): (a) 5-second resolution with 382 total events, and (b) 1-minute resolution with 187 total
504   events. Matrix rows show the actual classes and columns show the predicted classes for the following end
505   uses: w (washing machine), s (shower), f (faucet), t (toilet), r (refrigerator), and d (dishwasher). Cell color
506                          is proportional to the number of events in that cell.

507  In general, misclassifications do not cause significant degradation in predicting total water
508  consumption if they are infrequent and roughly symmetric across the diagonal (Srinivasan et al.
509  2011). For example, if toilet events are misclassified as faucet events while the same (or nearly
510  the same) number of faucet events are misclassified as toilet events, these misclassifications can
511  cancel out in terms of the accurate total number of events for those classes.

512  **4. Broader implications**

513  Overall, our study contributes to the literature showing that smart water meters provide water
514  utilities with more accurate and less labor-intensive information, enabling better knowledge on
515  changing water demands (Gurung et al. 2015; Stewart et al., 2018). High resolution temporal and
516  spatial water consumption data have undeniable social and technical benefits. Smart metering
517  contributes to more accurate water demand forecasting, demand management strategies, and
518  better informed utility operations and planning strategies (McDaniel and McLaughlin 2009;
519  Cominola et al. 2015; Salomons et al. 2020). Detailed water consumption patterns, which enable
520  researchers to investigate the relationships between human behaviors and the water cycle as part
521  of a broader socio-environmental scale, can be now obtained with advanced analytics, enabled
522  by fast paced computing power improvement and metering technology allowing data collection

18

523 with unprecedented temporal and spatial granularity (Flint et al. 2017; Zipper et al. 2019). While
524 these advances support greater understanding of water consumption patterns and water-related
525 human behaviors, we also acknowledge that there are potential privacy concerns regarding
526 individuals and communities that need to be addressed and appreciated. Water consumption
527 information transformed from the meter acts as an information side channel (McDaniel and
528 McLaughlin 2009), exposing household habits and behaviors. End-uses like showers and toilets
529 have detectable water consumption signatures, making end use classification information prone
530 to potential privacy abuse. Consequently, well established privacy policies would benefit utilities
531 in appropriate water demand management. Additionally, researchers have an ethical
532 responsibility to protect participant confidentiality.

533 Recent studies have addressed privacy issues in both the water and energy sectors and presented
534 solutions to overcome privacy related constraints to maximize the potential of granular data
535 (Khurana et al. 2010; Molina-Markham et al. 2010; Gurstein 2011; Amin 2012; Cole and
536 Stewart, 2013; Harter et al. 2013; Sankar et al. 2013; Helveston 2015; Park and Cominola 2020;
537 Salomons et al. 2020). For instance, smart meter data can be used without invading individual
538 privacy by aggregating data to coarser spatial or temporal scales as presented in our study.
539 Nevertheless, as shown in this study, aggregation limits the ability of end-use classification, or
540 any water consumption related research, to explore fine-scale behavioral dynamics for better
541 demand modeling. Therefore, any research intersecting with human behavior should prioritize
542 confidentiality (e.g., via anonymized data collected over a large sample of households) while
543 providing sufficient information to enable future improvements in that field. While the
544 formulation of privacy and security protection strategies is not within the scope of this study, we
545 acknowledge that privacy and security considerations must be addressed and proactively planned
546 for prior to collecting data throughout the research process so that modern metering technologies
547 could be leveraged to their full extent while securing customer privacy (Meyer 2018).

548 From the findings of this study, we can identify the following limitations and opportunities for
549 future research. First, future studies could focus on assessing how our results generalize when
550 data from a larger household sample or homes from different socio-demographic, geographical,
551 and climate contexts are available. Second, in this study we only considered six classes of indoor
552 water uses from a 4-person household. Further research could include outdoor water use and test
553 end-use disaggregation capabilities on houses with different sizes. Third, as highlighted in the
554 methods, end use datasets are often imbalanced, i.e., the number of events in each end use class
555 might vary substantially. While here we considered class imbalance *a posteriori*, by assessing
556 the disaggregation results with different formulations of the F-score, an alternative approach to
557 be tested when larger datasets are available is to balance the classes *a priori* (i.e., before
558 performing the classification), e.g., by oversampling/undersampling, which would solve the
559 problem of class imbalance. Finally, while here we only considered RF classifiers and a specific
560 approach for disaggregation, future studies could comparatively assess the performance of
561 different models, possibly accounting for multi-class events.

19

## 5. Conclusion

562

563 In this analysis, we presented a supervised approach to classify residential water consumption
564 end use events and tested it on data collected in a 4-person household through consideration of
565 multiple temporal resolutions by measuring water use data with a 1-second resolution smart
566 water metering system and labeling events based on a water diary for a 4-week study period. We
567 investigated two different scenarios of model calibration in evaluating the effect of temporal
568 resolution on end use classification performance. The first scenario consisted of training a
569 random forest classifier on the original 1-second resolution data only and testing it on other
570 labeled temporal resolution datasets (i.e., 5 seconds, 10 seconds, 30 seconds, 1 minute). In this
571 scenario, our model exhibited high overall performance on the 1-second and 5-second resolution
572 water use events and classified certain classes of end uses with fairly good accuracy for the 10-
573 second resolution. The performance decreased notably for the 30-second and 1-minute
574 resolutions.

575 The second scenario consisted of training separate models for each temporal resolution using k-
576 fold cross-validation. We saw that coarser temporal resolutions ameliorated in this second
577 scenario, with F1-score performance metrics as high as 0.89 for certain end use classes at the
578 finer resolutions. A weighted F1-score above 0.85 was obtained in this scenario for
579 disaggregation tasks performed at 1- and 5-second resolutions.

580 Our results reveal detailed information that can help utilities and residents make informed water
581 conservation and efficiency decisions based on detailed knowledge on water demands. The
582 analysis of classification model performance versus temporal resolution considering different F1-
583 scoore formulations provides insight for future water management regarding the selection of an
584 efficient monitoring resolution based on priorities and data management capabilities.

585 In addition, our approach performed end use disaggregation of data aggregated at different
586 temporal resolutions that are closer to the resolutions of commercial smart water meters (i.e., 1
587 minute). Thus, while making use of data collected at a finer resolution (e.g., 1 second) might not
588 be available to water utilities due to data management and analysis tradeoffs, we demonstrate
589 possible model extensions to broader and further contexts in the field of residential water
590 demand monitoring.

591 Ultimately, disaggregating and classifying water events obtained from residential smart water
592 meter data reveals detailed information about how water is consumed within households.
593 Understanding the overall water consumption profile and performance of different resolutions
594 presents opportunities for improved residential water conservation and efficiency and long-term
595 water resource sustainability (Attari 2014; Inskeep and Attari 2014; Horsburgh et al. 2017;
596 Goulas et al. 2022). Our study presents an experimental example of how using smart water meter
597 data can provide end use information to pinpoint opportunities for improved efficiency within
598 residential buildings.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

614 **References**

615

616 Abdallah, A. M., & Rosenberg, D. E. (2014). Heterogeneous residential water and energy
617 linkages and implications for conservation and management. *Journal of Water Resources*
618 *Planning and Management*, *140*(3), 288-297.

619 Amin, S. M. (2012, July). Smart grid security, privacy, and resilient architectures: Opportunities
620 and challenges. In *2012 IEEE Power and Energy Society General Meeting* (pp. 1-2). IEEE.

621 Armel, K. C., Gupta, A., Shrimali, G., & Albert, A. (2013). Is disaggregation the holy grail of
622 energy efficiency? The case of electricity. *Energy policy*, *52*, 213-234.

623 Attari, S. Z. (2014). Perceptions of water use. *Proceedings of the National Academy of*
624 *Sciences*, *111*(14), 5129-5134.

625 Beal, C. D., Gurung, T. R., & Stewart, R. A. (2016). Demand-side management for supply-side
626 efficiency: Modeling tailored strategies for reducing peak residential water demand. *Sustainable*
627 *Production and Consumption*, *6*, 1-11.

628 Beal, C., Stewart, R., Huang, T., & Rey, E. (2011). *South East Queensland residential end use*
629 *study*. Brisbane, Australia: Urban Water Security Research Alliance.

630 Bethke, G. M., Cohen, A. R., & Stillwell, A. S. (2021). Emerging investigator series:
631 disaggregating residential sector high-resolution smart water meter data into appliance end-uses
632 with unsupervised machine learning. *Environmental Science: Water Research &*
633 *Technology*, *7*(3), 487-503.

634 Blokker, E. J. M., Vreeburg, J. H. G., & Van Dijk, J. C. (2010). Simulating residential water
635 demand with a stochastic end-use model. *Journal of Water Resources Planning and*
636 *Management*, *136*(1), 19-26.

637 Boyle, T., Giurco, D., Mukheibir, P., Liu, A., Moy, C., White, S., & Stewart, R. (2013).
638 Intelligent metering for urban water: A review. *Water*, *5*(3), 1052-1081.

639 Breiman, L. (2001). Random forests. *Machine learning*, *45*(1), 5-32.

640 Clifford, E., Mulligan, S., Comer, J., & Hannon, L. (2018). Flow-signature analysis of water
641 consumption in nonresidential building water networks using high-resolution and
642 medium-resolution smart meter data: Two case studies. *Water Resources Research*, *54*(1), 88-
643 106.

644 Cole, G., & Stewart, R. A. (2013). Smart meter enabled disaggregation of urban peak water
645 demand: precursor to effective urban water planning. *Urban Water Journal*, *10*(3), 174-194.

646 Cominola, A., Becker, M.-P., and Taormina, R.: Benchmarking machine learning algorithms for
647 Non-Intrusive Water Monitoring, EGU General Assembly 2020, Online, 4–8 May 2020,
648 EGU2020-16119, https://doi.org/10.5194/egusphere-egu2020-16119, 2020

649 Cominola, A., Giuliani, M., Castelletti, A., Fraternali, P., Gonzalez, S. L. H., Herrero, J. C. G., ...
650 & Rizzoli, A. E. (2021). Long-term water conservation is fostered by smart meter-based
651 feedback and digital user engagement. *npj Clean Water*, *4*(1), 1-10.

652 Cominola, A., Giuliani, M., Castelletti, A., Rosenberg, D. E., & Abdallah, A. M. (2018).
653 Implications of data sampling resolution on water use simulation, end-use disaggregation, and
654 demand management. *Environmental Modelling & Software*, *102*, 199-212.

655 Cominola, A., Giuliani, M., Piga, D., Castelletti, A., & Rizzoli, A. E. (2015). Benefits and
656 challenges of using smart meters for advancing residential water demand modeling and
657 management: A review. *Environmental Modelling & Software*, *72*, 198-214.

658 Cominola, A., Giuliani, M., Piga, D., Castelletti, A., & Rizzoli, A. E. (2017). A hybrid signature-
659 based iterative disaggregation algorithm for non-intrusive load monitoring. *Applied energy*, *185*,
660 331-344.

661 DeOreo, W. B. (2011). Analysis of water use in new single family homes. *By Aquacraft. For Salt
662 Lake City Corporation and US EPA*.

663 DeOreo, W. B., Mayer, P. W., Dziegielewski, B., & Kiefer, J. (2016). *Residential end uses of
664 water, version 2*. Water Research Foundation.

665 Di Mauro, A., Cominola, A., Castelletti, A., & Di Nardo, A. (2021). Urban Water Consumption
666 at Multiple Spatial and Temporal Scales. A Review of Existing Datasets. *Water*, *13*(1), 36.

667 Di Mauro, A., Di Nardo, A., Francesco Santonastaso, G., & Venticinque, S. (2020).
668 Development of an IoT System for the Generation of a Database of Residential Water End-Use
669 Consumption Time Series. *Environmental Science Proceedings*, *2*(1), 20.

670 Escriva-Bou, A., Lund, J. R., & Pulido-Velazquez, M. (2018). Saving energy from urban water
671 demand management. *Water Resources Research*, *54*(7), 4265-4276.

672 Figueiredo, M., Ribeiro, B., & de Almeida, A. (2013). Electrical signal source separation via
673 nonnegative tensor factorization using on site measurements in a smart home. *IEEE Transactions
674 on Instrumentation and Measurement*, *63*(2), 364-373.

675 Friedman, J., Hastie, T., & Tibshirani, R. (2001). 4.3: Linear Discriminant Analysis. *The
676 elements of statistical learning*, *1*, 106-119.

677 Froehlich, J., Larson, E., Saba, E., Campbell, T., Atlas, L., Fogarty, J., & Patel, S. (2011, June).
678 A longitudinal study of pressure sensing to infer real-world water usage events in the home.
679 In *International conference on pervasive computing* (pp. 50-69). Springer, Berlin, Heidelberg.

680 Géron, A. (2019). *Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow:*
681 *Concepts, tools, and techniques to build intelligent systems*. " O'Reilly Media, Inc.".

682 Goulas, A., Goodwin, D., Shannon, C., Jeffrey, P., & Smith, H. M. Public perceptions of
683 household IoT smart water 'event'meters in the UK–implications for urban water governance.

684 Gurstein, M. B. (2011). Open data: Empowering the empowered or effective data use for
685 everyone?. *First Monday*.

686 Gurung, T. R., Stewart, R. A., Beal, C. D., & Sharma, A. K. (2015). Smart meter enabled water
687 end-use demand data: platform for the enhanced infrastructure planning of contemporary urban
688 water supply networks. *Journal of Cleaner Production*, *87*, 642-654.

689 Hartter, J., Ryan, S. J., MacKenzie, C. A., Parker, J. N., & Strasser, C. A. (2013). Spatially
690 explicit data: stewardship and ethical challenges in science. *PLoS Biology*, *11*(9), e1001634.

691 Hastie, T., Tibshirani, R., & Friedman, J. (2009). Random forests. In *The elements of statistical*
692 *learning* (pp. 587-604). Springer, New York, NY.

693 Hastie, T., Tibshirani, R., Botstein, D., & Brown, P. (2001). Supervised harvesting of expression
694 trees. *Genome Biology*, *2*(1), 1-12.

695 Hawkins, D. M., Basak, S. C., & Mills, D. (2003). Assessing model fit by cross-
696 validation. *Journal of chemical information and computer sciences*, *43*(2), 579-586.

697 Helveston, M. N. (2015). Consumer protection in the age of big data. *Wash. UL Rev.*, *93*, 859.

698 Herrera, M., Torgo, L., Izquierdo, J., & Pérez-García, R. (2010). Predictive models for
699 forecasting hourly urban water demand. *Journal of hydrology*, *387*(1-2), 141-150.

700 Horsburgh, J. S., Leonardo, M. E., Abdallah, A. M., & Rosenberg, D. E. (2017). Measuring
701 water use, conservation, and differences by gender using an inexpensive, high frequency
702 metering system. *Environmental modelling & software*, *96*, 83-94.

703 Inskeep, B. D., & Attari, S. Z. (2014). The water short list: The most effective actions US
704 households can take to curb water use. *Environment: Science and policy for sustainable*
705 *development*, *56*(4), 4-15.

706 J. Fogarty, C. Au, and S. Hudson. Sensing from the basement: a feasibility study of unobtrusive
707 and low-cost home activity recognition. In Proceedings of UIST, 2006.

708 J. Froehlich, E. Larson, T. Campbell, C. Haggerty, J. Fogarty, and S. Patel. Hydrosense:
709 infrastructure-mediated singlepoint sensing of whole-home water activity. In In UbiComp, 2009.

710 Jain, A., & Ormsbee, L. E. (2002). Short-term water demand forecast modeling techniques—
711 CONVENTIONAL METHODS VERSUS AI. *Journal-American Water Works*
712 *Association*, *94*(7), 64-72.

713    James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical*
714    *learning* (Vol. 112, p. 18). New York: springer.

715    Karamouz, M., & Heydari, Z. (2020). Conceptual design framework for coastal flood best
716    management practices. *Journal of Water Resources Planning and Management*, *146*(6),
717    04020041.

718    Khurana, H., Hadley, M., Lu, N., & Frincke, D. A. (2010). Smart-grid security issues. *IEEE*
719    *Security & Privacy*, *8*(1), 81-85.

720    Kowalski, M., & Marshallsay, D. (2003, April). A system for improved assessment of domestic
721    water use components. In *II International Conference Efficient Use and Management of Urban*
722    *Water Supply*.

723    Makonin, S., Popowich, F., Bajić, I. V., Gill, B., & Bartram, L. (2015). Exploiting HMM
724    sparsity to perform online real-time nonintrusive load monitoring. *IEEE Transactions on smart*
725    *grid*, *7*(6), 2575-2585.

726    Mantovani, R. G., Horváth, T., Cerri, R., Junior, S. B., Vanschoren, J., & de Carvalho, A. C. P.
727    D. L. F. (2018). An empirical study on hyperparameter tuning of decision trees. *arXiv preprint*
728    *arXiv:1812.02207*.

729    Mayer, P. W., DeOreo, W. B., Opitz, E. M., Kiefer, J. C., Davis, W. Y., Dziegielewski, B., &
730    Nelson, J. O. (1999). Residential end uses of water.

731    Mazzoni, F., Alvisi, S., Franchini, M., Ferraris, M., & Kapelan, Z. (2021). Automated Household
732    Water End-Use Disaggregation through Rule-Based Methodology. *Journal of Water Resources*
733    *Planning and Management*, *147*(6), 04021024.

734    McDaniel, P., & McLaughlin, S. (2009). Security and privacy challenges in the smart grid. *IEEE*
735    *security & privacy*, *7*(3), 75-77.

736    Meyer, M. N. (2018). Practical tips for ethical data sharing. *Advances in methods and practices*
737    *in psychological science*, *1*(1), 131-144.

738    Molina-Markham, A., Shenoy, P., Fu, K., Cecchet, E., & Irwin, D. (2010, November). Private
739    memoirs of a smart meter. In *Proceedings of the 2nd ACM workshop on embedded sensing*
740    *systems for energy-efficiency in building* (pp. 61-66).

741    Nguyen, K. A., Zhang, H., & Stewart, R. A. (2013). Development of an intelligent model to
742    categorise residential water end use events. *Journal of hydro-environment research*, *7*(3), 182-
743    201.

744    Ntuli, N., & Abu-Mahfouz, A. (2016). A simple security architecture for smart water
745    management system. *Procedia Computer Science*, *83*, 1164-1169.

746 Park, S., & Cominola, A. When Privacy Protection Meets Non-Intrusive Load Monitoring:
747 Trade-off Analysis and Privacy Schemes via Residential Energy Storage.

748 Pastor-Jabaloyes, L., Arregui, F. J., & Cobacho, R. (2018). Water end use disaggregation based
749 on soft computing techniques. *Water*, *10*(1), 46.

750 Piga, D., Cominola, A., Giuliani, M., Castelletti, A., & Rizzoli, A. E. (2015). Sparse optimization
751 for automated energy end use disaggregation. *IEEE Transactions on Control Systems*
752 *Technology*, *24*(3), 1044-1051.

753 Probst, P., Wright, M. N., & Boulesteix, A. L. (2019). Hyperparameters and tuning strategies for
754 random forest. *Wiley Interdisciplinary Reviews: data mining and knowledge discovery*, *9*(3),
755 e1301.

756 Probst, P., Wright, M. N., & Boulesteix, A. L. (2019). Hyperparameters and tuning strategies for
757 random forest. *Wiley Interdisciplinary Reviews: data mining and knowledge discovery*, *9*(3),
758 e1301.

759 Probst, P., Wright, M. N., & Boulesteix, A. L. (2019). Hyperparameters and tuning strategies for
760 random forest. *Wiley Interdisciplinary Reviews: data mining and knowledge discovery*, *9*(3),
761 e1301.

762 Rahimpour, A., Qi, H., Fugate, D., & Kuruganti, T. (2017). Non-intrusive energy disaggregation
763 using non-negative matrix factorization with sum-to-k constraint. *IEEE Transactions on Power*
764 *Systems*, *32*(6), 4430-4441.

765 Raileanu, L. E., & Stoffel, K. (2004). Theoretical comparison between the gini index and
766 information gain criteria. *Annals of Mathematics and Artificial Intelligence*, *41*(1), 77-93.

767 Ribeiro, V. H. A., & Reynoso-Meza, G. (2020). Ensemble learning by means of a multi-
768 objective optimization design approach for dealing with imbalanced data sets. *Expert Systems*
769 *with Applications*, *147*, 113232.

770 Salomons, E., Sela, L., & Housh, M. (2020). Hedging for privacy in smart water meters. *Water*
771 *Resources Research*, *56*(9), e2020WR027917.

772 Sankar, L., Rajagopalan, S. R., Mohajer, S., & Poor, H. V. (2012). Smart meter privacy: A
773 theoretical framework. *IEEE Transactions on Smart Grid*, *4*(2), 837-846.

774 Sønderlund, A. L., Smith, J. R., Hutton, C. J., Kapelan, Z., & Savic, D. (2016). Effectiveness of
775 smart meter-based consumption feedback in curbing household water use: Knowns and
776 unknowns. *Journal of Water Resources Planning and Management*, *142*(12), 04016060.

777 Srinivasan, V., Stankovic, J., & Whitehouse, K. (2011, November). Watersense: Water flow
778 disaggregation using motion sensors. In *Proceedings of the Third ACM Workshop on Embedded*
779 *Sensing Systems for Energy-Efficiency in Buildings* (pp. 19-24).

780 Stewart, R. A., Nguyen, K., Beal, C., Zhang, H., Sahin, O., Bertone, E., ... & Kossieris, P.
781 (2018). Integrated intelligent water-energy metering systems and informatics: Visioning a digital
782 multi-utility service provider. *Environmental Modelling & Software*, *105*, 94-117.

783 Stewart, R. A., Willis, R., Giurco, D., Panuwatwanich, K., & Capati, G. (2010). Web-based
784 knowledge management system: linking smart metering to the future of urban water
785 planning. *Australian Planner*, *47*(2), 66-74.

786 Suero, F. J., Mayer, P. W., & Rosenberg, D. E. (2012). Estimating and verifying United States
787 households' potential to conserve water. *Journal of Water Resources Planning and*
788 *Management*, *138*(3), 299-306.

789 Vitter, J. S., & Webber, M. E. (2018). A non-intrusive approach for classifying residential water
790 events using coincident electricity data. *Environmental modelling & software*, *100*, 302-313.

791 Willis, R. M., Stewart, R. A., Panuwatwanich, K., Jones, S., & Kyriakides, A. (2010). Alarming
792 visual display monitors affecting shower end use water and energy conservation in Australian
793 residential households. *Resources, Conservation and Recycling*, *54*(12), 1117-1127.

794 Wonders, M., Ghassemlooy, Z., & Hossain, M. A. (2016). Training with synthesised data for
795 disaggregated event classification at the water meter. *Expert Systems with Applications*, *43*, 15-
796 22.

797 Y. Kim, T. Schmid, Z. Charbiwala, J. Friedman, and M. Srivastava. Nawms: nonintrusive
798 autonomous water monitoring system. In Sensys, 2008.

799 Zipper, S. C., Stack Whitney, K., Deines, J. M., Befus, K. M., Bhatia, U., Albers, S. J., ... &
800 Schlager, E. (2019). Balancing open science and data privacy in the water sciences. *Water*
801 *Resources Research*, *55*(7), 5202-5211.

802