

Towards a self-tuned data analytics-based process for an automatic context-aware detection and diagnosis of anomalies in building energy consumption timeseries

Original

Towards a self-tuned data analytics-based process for an automatic context-aware detection and diagnosis of anomalies in building energy consumption timeseries / Chiosa, R.; Piscitelli, M. S.; Fan, C.; Capozzoli, A.. - In: ENERGY AND BUILDINGS. - ISSN 0378-7788. - 270:(2022), p. 112302. [10.1016/j.enbuild.2022.112302]

Availability:

This version is available at: 11583/2970367 since: 2022-09-28T16:09:18Z

Publisher:

Elsevier

Published

DOI:10.1016/j.enbuild.2022.112302

Terms of use:

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

Publisher copyright

Elsevier postprint/Author's Accepted Manuscript

© 2022. This manuscript version is made available under the CC-BY-NC-ND 4.0 license
<http://creativecommons.org/licenses/by-nc-nd/4.0/>. The final authenticated version is available online at:
<http://dx.doi.org/10.1016/j.enbuild.2022.112302>

(Article begins on next page)

Towards a self-tuned data analytics-based process for an automatic context-aware detection and diagnosis of anomalies in building energy consumption timeseries

Roberto Chiosa^a, Marco Savino Piscitelli^{a*}, Cheng Fan^b, Alfonso Capozzoli^a

^a Department of Energy "Galileo Ferraris", TEBE Research Group, BAEDA Lab, Politecnico di Torino, Corso Duca degli Abruzzi 24, 10129 Turin, Italy

^b Department of Construction Management and Real Estate, College of Civil and Transportation Engineering, Shenzhen University, Shenzhen, China

Abstract

Recently, the spread of IoT technologies has led to an unprecedented acquisition of energy-related data providing accessible knowledge on the actual performance of buildings during their operation. A proper analysis of such data supports energy and facility managers in spotting valuable energy saving opportunities. In this context, anomaly detection and diagnosis (ADD) tools allow a prompt and automatic recognition of abnormal and non-optimal energy performance patterns enabling a better decision-making to reduce energy wastes and system inefficiencies. To this aim, this paper introduces a novel meter-level ADD process capable to identify energy consumption anomalies at meter-level and perform diagnosis by exploiting information at sub-load level. The process leverages supervised and unsupervised analytics techniques coupled with the distance-based contextual matrix profile (CMP) algorithm to discover infrequent subsequences in energy consumption timeseries considering specific boundary conditions. The proposed process has self-tuning capabilities and can rank anomalies at both meter and sub-load level by means of robust severity score. The methodology is tested on one-year energy consumption timeseries of a medium/low voltage transformation cabin of the university campus of Politecnico di Torino leading to the detection of 55 anomalous subsequences that are diagnosed by analysing a group of 8 different sub-loads.

Keywords: building energy consumption, anomaly detection and diagnosis, contextual matrix profile, timeseries analytics

1. Introduction

In the last few years, the increasing widespread use of IoT sensors in buildings for the pervasive monitoring of energy-related data, has led to an unprecedented acquisition of reliable and accessible knowledge related to the actual performance of buildings during their operation. Considering that in Europe the building sector accounts for 40% of final energy use [1] and almost 90% of the total energy consumed during the life cycle of a building depends on its operation [2], supporting building owners and energy managers to extract valuable information from monitoring data is of paramount importance to (i) reduce energy consumption, (ii) increase system efficiency, (iii) prevent energy wastes and (iv) operate their buildings more efficiently.

Although a great deal of research has been done, the increasing data volume still overwhelms end-users [3], making it hard to clearly spot energy reduction opportunities, find the root cause of anomalies or simply be aware of energy usage of buildings and systems. Building-related data are heterogeneous and reflects the complex interactions occurring between occupants, energy systems, building envelope, and forcing conditions [4]. If properly managed, ingested and analysed, such data provide the opportunity to gain insight on the building operational behaviour discovering valuable and ready-to-implement energy conservation measures [5].

A robust coupling of IoT sensors data, artificial intelligence (AI) and energy domain knowledge proved to be effective in achieving relevant energy saving by exploiting a variety of energy management solutions [6]. The tools providing such capabilities are the so-called energy management and information systems (EMIS) which are employed to monitor,

* Corresponding author

Email address: marco.piscitelli@polito.it (Marco Savino Piscitelli)

URL: www.baeda.polito.it (Marco Savino Piscitelli)

41 analyse and control energy systems in buildings leveraging advanced data analytics technologies for supporting facility
42 staff to enhance energy performance and efficiency [6]. Depending on the level of detail of the measured data, EMIS
43 solutions can be classified as meter-level or system-level: the first includes analysis of high-level measurements (e.g.,
44 data related to the whole-building electrical load or to the main sub-loads) while the second focuses on more detailed data
45 pertaining to the operation of specific energy systems or components (e.g., component operation of air handling units in
46 HVAC systems).

47 A subgroup of EMIS, conceived for the collection and analysis of meter-level data, are the so-called energy information
48 systems (EIS). EIS typically focus on data not usually collected through building automation systems (BAS) providing
49 visual and analytical insights also enabling predictive energy management strategies. According to Granderson et al. [9]
50 EIS tools, such as energy consumption forecasting [7–10], anomaly detection and diagnosis (ADD) [11–15], advanced
51 energy benchmarking [16–19], load profiling [20–22], and schedule optimization of building energy systems [23,24]
52 proved to lead to systematic energy saving and, as reported in [6], to an average two-year simple payback period, thus
53 demonstrating their high competitiveness as profitable investment option in the building sector.

54 Among these solutions, ADD has been the most underdeveloped for application on meter-level data [11,12]. ADD
55 tools allow a prompt and automatic recognition of abnormal and non-optimal performance patterns of energy systems
56 providing information for the identification of energy waste and for the prioritization of corrective interventions. While
57 fault detection and diagnosis (FDD) tools analyse system/component-level data to detect faults and anomalies, ADD tools
58 generally rely on aggregated meter-level data to automatically detect anomalous energy trends at whole-building scale.
59 Although performing a meter-level analysis poses several challenges related to the influence of several factors (e.g.,
60 building features, climatic conditions, occupant behaviour, comfort requirements, operating schedules and management)
61 it is of considerable value in real world case studies where the available measured variables are commonly related to the
62 meter scale [15]. In this context, the main objectives for an ADD process are: (i) the recognition of typical patterns in the
63 whole-building energy consumption timeseries, (ii) the detection of infrequent/anomalous patterns and (iii) the diagnosis
64 of the detected anomalies by inferring the occurrence of anomalous patterns at sub-load level.

65 According to these objectives, this paper introduces the conceptualization and development of a novel EIS tool capable
66 of robustly performing a meter-level ADD on building electrical energy consumption timeseries. The proposed approach
67 employs a pattern recognition technique derived from the matrix profile (MP) [25] algorithm, called contextual matrix
68 profile (CMP) [26] for the automatic detection of energy consumption anomalies at the whole-building level and their
69 diagnosis by analysing the group of associated sub-loads. The following sections review the main works related to ADD
70 in the energy and buildings field and to the application of MP-based algorithms as anomaly detection methods in
71 timeseries.

72 1.1. Related work

73 Generally speaking an anomaly is a region of data with significantly different behaviour from other data and that do
74 not conform to expected values [27]. It can be referred as discord, deviation or exception and its definition is significantly
75 different depending on the field of application and the analysis performed. In the energy and buildings field, which mainly
76 involves univariate timeseries data (e.g., electrical energy consumption), the definition of anomaly is very domain-
77 specific and may include abnormal behaviour of occupants, faulty operations of appliances, incorrect management of
78 energy systems, anomalous sub-load energy consumption and technical and non-technical energy losses [28]. Thus, the
79 nature of energy timeseries data in buildings requires to carefully address the definition of anomaly, which can be
80 classified as point, collective, or contextual [29]. A *point anomaly* is one individual instance or observation that can be

81 considered anomalous when compared to the remaining data. A *collective anomaly* is a collection of anomalous instances
82 with respect to the entire dataset. Eventually, a *context anomaly* is so defined only if considered in a certain context (i.e.,
83 boundary conditions) and may not be considered an anomaly in a different context [30].

84 ADD has been traditionally addressed by means of statistical analysis, however the increasing spread of advanced
85 machine learning and data analytics techniques coupled with the large volume of available data have opened new
86 possibility to develop more sophisticated and reliable data-driven processes [3].

87 In ADD, the detection phase is usually accomplished by estimating a reference baseline representing normal behaviour
88 according to specific boundary conditions, and labelling each observation that diverges from it as anomalous [29]. The
89 discrimination between normal and abnormal behaviour is essential, so both the robust development of a reference model
90 and the proper selection of the features, used to define anomalies, are of paramount importance.

91 Supervised methods have been used to train machine learning algorithms using labelled dataset (i.e., datasets with
92 verified ground truth about the presence of an anomaly) to create reference models able to distinguish anomalous from
93 normal energy consumption. Support vector machines (SVM) and multi-layer perceptron are largely used to perform
94 model-based anomaly detection [31]. Zhao et al. [32] developed an SVM-based process that proved to be effective in
95 detecting anomalies with different severities in chiller operation. Regression and decision tree classifiers are other widely
96 used supervised techniques to develop energy consumption reference models for discovering abnormal patterns [33]. An
97 evolutionary tree model was employed as a detection method in [12] to effectively discover frequent and infrequent
98 patterns in meter-level electrical loads. A hybrid neural net ARIMA model was employed in [34] to predict the energy
99 consumption and then identify anomalies comparing the actual and predicted energy consumption using the two-sigma
100 rule. In [11] a regression model, for detecting the anomalous trends in electrical energy consumption, was developed by
101 coupling artificial neural network (ANN) and regression tree (RT).

102 Although supervised approach can achieve very accurate results, its adoption in real-world ADD applications is still
103 limited compared to unsupervised methods, mainly due to the absence of reliable annotated datasets [28,35] and the
104 difficulty in obtaining high quality training data [23] for machine learning algorithms.

105 *Unsupervised* anomaly detection is more promising for practical ADD applications since it makes it possible to detect
106 rare and unknown anomalous patterns, without any a-priori knowledge and does not require pre-labelled datasets [36].
107 Beside statistical unsupervised methods, such as principal component analysis (PCA) [37] or generalized extreme
108 studentized deviate (GESD) [8], data mining methods, such as association rule mining (ARM) [27] and clustering analysis
109 [38], gained popularity thanks to their capability to automatically extract significant relations between complex and
110 massive data [5,39,40]. In this context, timeseries analytics has been stimulating a great deal of interest in the scientific
111 literature in recent years, since building-related data are often stored in datasets in form of timeseries on which the
112 anomaly detection is performed also considering the time domain. A typical timeseries data mining approach for anomaly
113 detection involves the extraction of subsequences related to discordant observations (i.e., discords) that diverges from the
114 rest of the dataset [28,41,42]. Li et al. [43] employed GESD to identify anomalous observations in electricity usage
115 timeseries referred to 40 buildings. Fan et al. [8] identified typical daily load profiles through entropy-weighted k-means
116 (EWKM) clustering and then abnormal daily energy consumption profiles were identified through GESD. An anomaly
117 detection framework applied on smart meter data stream was presented in [44], where ARM and categorical clustering
118 were used to detect anomalies.

119 In some cases, a reduction and transformation of timeseries data, could enhance the performance and reduce the
120 computational time of anomaly detection techniques, also helping to effectively extract useful information. Lin et. al [45]
121 proposed symbolic aggregate approximation (SAX) as a method for the reduction of a timeseries and its transformation

122 in a symbolic sequence for an easy detection of relevant symbolic strings (i.e., motifs and discords). This method, which
123 introduces a simple and low-computational cost process to reduce a timeseries while preserving the key information, is
124 an extension of the piecewise aggregate approximation (PAA) technique and was employed in the literature also for the
125 recognition of frequent/infrequent patterns in energy consumption timeseries of buildings [46,47]. Miller et al. [13]
126 through a SAX-based analysis identified the most infrequent symbolic sequences referred to daily load profiles of non-
127 residential buildings and furtherly characterized those patterns carrying out a cluster analysis. SAX and temporal
128 association rule mining (TARM) were employed in [14] to extract discords in the energy consumption timeseries, assess
129 building system performance and suggest the implementation of possible energy conservation measures. An adaptive
130 SAX method (aSAX) was employed in [15] to optimize the dimensionality reduction of an energy consumption timeseries
131 and to enhance the detection of frequent and infrequent patterns through a classification tree model [12].

132 Despite SAX introduced a lot of opportunities in the field of ADD in timeseries, it is an approximation-based pattern
133 recognition technique and the dimensionality reduction provided by PAA method coupled with the symbolic encoding
134 always lead to information loss from the original timeseries [47]. In addition, the information loss is particularly sensitive
135 respect to the setting of input parameters such as the time window length and the number of symbols for the timeseries
136 encoding [13].

137 One of the most promising timeseries analytics techniques, that is not subjected to information loss, is matrix profile
138 (MP). Introduced by [25], MP is a novel exact pattern recognition algorithm that performs all-similarity-join-search (i.e.,
139 full-join) among timeseries, i.e. given a collection of data objects it retrieves the nearest neighbour for every object (where
140 the considered object is intended to be a timeseries subsequence). MP is an unsupervised distance-based anomaly
141 detection algorithm that proposes a fast similarity search under the z-normalized Euclidean distance, does not reduce
142 dimensionality, calculates the full-join among timeseries and eliminates the need of setting a threshold making the method
143 almost parameter-free and exact. MP algorithm allows the method to be incrementally maintainable, deterministic in time
144 and so parallelizable on multicore processor and distributed systems.

145 The MP method has been successfully applied in different fields of anomaly detection. Alshaer et al. [48] proposed a
146 real time MP-based anomaly detection method which was tested on electro-cardiogram (ECG) timeseries. The method
147 employs the concept of shaplet [49] to perform continuous learning of anomalies, which are extracted through MP
148 algorithm, stored in an anomaly library and then used for sliding-window based anomaly detection. An industrial
149 application was reported in [50], where the MP was combined with the hamming distance to automatically detect
150 intrusions in the network of a water processing facility. PanMP is a generalization of MP algorithm able to discover
151 anomalous subsequences of different lengths that was introduced in [51], where it was employed to perform anomaly
152 detection in automated pedestrian counting system developed in Taipei.

153 MP has been also largely employed to identify anomalies in IT field. Herath et al. [52] introduced a real time anomaly
154 detection framework based on MP, called real-time aggregated matrix profile (RAMP), that can identify anomalies in
155 scientific workflows. De Paepe et al. [53] applied a noise elimination technique on real Yahoo! internet traffic timeseries
156 and detected anomalous behaviours through MP; while in [54] was demonstrated how the elimination of noise can help
157 in anomaly detection of noisy dataset by testing the algorithm on Numenta Benchmark [55].

158 In the energy field few implementation studies of MP algorithm are available. Nichiforov et al. [56] used the MP to
159 provide insights about the dominant energy usage patterns in large academic buildings. The authors applied the MP
160 approach with daily, weekly, and monthly time window length as a feature extraction method to identify unusual
161 behaviours in energy consumption timeseries of a large academic building dataset with a length of one year. The process
162 was tested on 422 buildings whose end-uses were classrooms, offices, laboratories and dormitory. Zhu et al. [57]

163 demonstrated how MP can be useful in detecting rare anomalous electricity consumption occasionally produced by a
164 meter swapping events. The algorithm was tested on a synthetic meter swapping event built on top of two timeseries of
165 household electrical demand and was proven to be effective to discover the suspicious similarity between the two
166 timeseries. Park et al. [21] applied MP as a part of an automated load profile discord identification (ALDI) based on
167 statistic comparison between normal and anomalous energy consumption patterns in a large portfolio of buildings. The
168 MP method was used to quantify the similarities of daily electrical load subsequences under z-normalized distance. The
169 computed MP values were then compared with typical-day MP distribution and it was effective in the identification of
170 unique load shape patterns and discords.

171 Despite MP proved to be effective in several application fields, in its full-join formulation, it may compare regions of
172 timeseries characterized by different operating and boundary conditions (i.e., contexts) and, in specific domains of
173 application, may conduct to misleading results in terms of pattern similarity. For this reason, in [26] was introduced a
174 matrix profile-based algorithm called contextual matrix profile (CMP). CMP allows to identify contexts in the timeseries,
175 in which it is possible to compare between each other only subsequences that are characterized by homogeneous boundary
176 conditions, avoiding the identification of meaningless similarity matches. In brief, while MP searches for unique patterns
177 (i.e., discords) in the whole timeseries, the CMP performs the same search within a reference context set by the analyst.

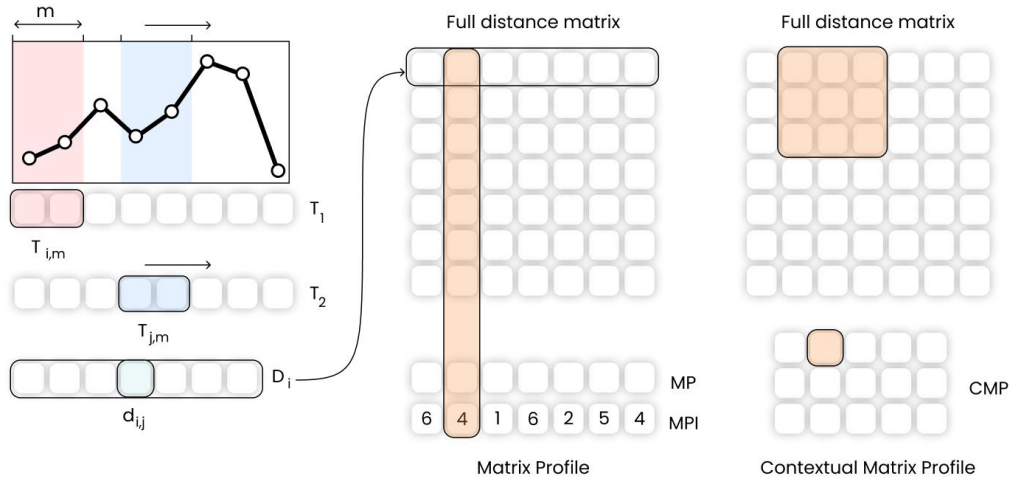
178 An application of the CMP method, for anomaly detection in buildings, is presented in [26], where the method was
179 applied to a dataset including different variables related to indoor air quality (e.g., temperature, humidity, CO₂ etc.)
180 measured in residential built environments. CMP was applied on CO₂ concentration timeseries, enabling the identification
181 of six behavioural anomalies (i.e., subsequences) occurring only during weekend mornings. Moreover, the authors were
182 able to identify periodic behavioural patterns even if there was not a time alignment among them. The case study
183 demonstrated the flexibility of the anomaly detection method and its effectiveness when coupled with domain knowledge.
184 Therefore, it emerges a great potential in the application of CMP in the building energy field, in which the definition of
185 anomaly is strongly related to the expert definition of contexts and boundary conditions, also leading to the robust
186 recognition of patterns not easy to be detected. To the best of author knowledge, other applications of CMP in buildings,
187 including ADD in energy consumption timeseries, are so far missing in the literature. In the next section, the main
188 theoretical aspects related to the MP and CMP algorithms are introduced and discussed to better highlight the
189 contributions brought by this study in the field of ADD in energy and buildings.

190 *1.2. Matrix profile and contextual matrix profile method for anomaly detection*

191 As defined in [25], given a collection of data objects (e.g., timeseries), MP performs the so-called similarity-join-
192 search which is able to find the nearest neighbour for every object (e.g., identification of nearest pairs of subsequences)
193 from the two collections.

194 Given two timeseries and a subsequence length, the MP algorithm produces two new series: the MP and matrix profile
195 index (MPI). MP is a one-dimensional timeseries that stores the z-normalized Euclidean distance values between each
196 subsequence, of the first series, and the closest matching subsequence (i.e., nearest neighbour) of the second timeseries.
197 MPI is a one-dimensional timeseries that contains the position index of the nearest neighbour in the second timeseries.
198 By joining information of MP and MPI useful knowledge could be extracted. By finding the minimum value of the MP
199 it is possible to identify the best matching subsequences in a series (i.e., motif discovery), on the other hand, by finding
200 the maximum value of the MP it is possible to identify the subsequence with the largest distance from its nearest match
201 (i.e., discord discovery) [25]. In this sense, discord discovery may be interpreted as an anomaly detection procedure that
202 identifies the most unique subsequences in a dataset.

203 With reference to Figure 1, some fundamental concepts and definitions need to be introduced before going deeper into
 204 the topic.



205

206 **Figure 1.** Description of MP and CMP calculation steps in case of self-join of a timeseries $T_1 = T_2$. From left to right is explained the calculation of
 207 the element $d_{i,j}$ of the distance vector D_i given the query $T_{i,m}$. By calculating the distance vector for the *all-subsequences-set* of T , and storing those
 208 values in a matrix, the *full distance matrix* is obtained. MP is the row wise minimum while the CMP is the minimum over rectangular regions.

209 First, a *timeseries* $T \in R^n$ is a sequence of real-valued numbers $t_i \in R : T = \{t_1, t_2, \dots, t_n\}$ with $1 \leq i \leq n$ where n
 210 is the length of T . Since the focus is on local properties of timeseries (i.e., portions of timeseries), a *subsequence* $T_{i,m} \in$
 211 R^m is defined as a continuous subset of values from T of length m starting in position i ; formally defined as
 212 $t_i \in R : T_{i,m} = \{t_i, t_{i+1}, \dots, t_{i+m-1}\}$ with $1 \leq i \leq n - m + 1$.

213 An ordered set of all possible subsequences of T , obtained by sliding a window of length m across T , is called *all-*
 214 *subsequences-set* A of a timeseries T and is formally defined as follows: $A = \{T_{1,m}, T_{2,m}, \dots, T_{n-m+1,m}\}$ where m is a
 215 user-defined subsequence length.

216 By computing the distance between a given query (i.e., subsequence $T_{i,m}$) and each subsequence in the *all-*
 217 *subsequences-set* A , it is possible to define a vector of distances called *distance profile* D_i of a timeseries T . Formally,
 218 $D_i = [d_{i,1}, d_{i,2}, \dots, d_{i,n-m+1}]$, where $d_{i,j} = \text{dist}(T_{i,m}, T_{j,m})$ for all $j \in [1, 2, \dots, n - m + 1]$, where $i \neq j$ and dist is the
 219 distance metric applied. It is possible to adopt different kind of distances to compute the *distance profile* as reported in
 220 [57–59] but the original MP method makes use of Euclidean distance between the z-normalized subsequences.

221 If the *distance profile* is calculated between a query in $T_{i,m}$ and the *all-subsequences-set* of T (i.e., self-join), by
 222 definition the i^{th} location of the *distance profile* D_i is zero since the distance is calculated between the query and itself,
 223 $d_{i,i} = \text{dist}(T_{i,m}, T_{i,m}) = 0$. Moreover, the distance is close to zero just before and after this position. Those matches are
 224 called *trivial matches* and are usually avoided during similarity-search by imposing an *exclusion zone* before and after
 225 this location (as function of m , usually set to $m/4$).

226 By calculating all the *distance profiles* for the *all-subsequences-set* of T , the *full distance matrix* is obtained. Formally
 227 defined as $\mathcal{M} = [D_1, D_2, \dots, D_{n-m+1}]$, it includes $n - m + 1$ *distance profiles* arranged in a square matrix.

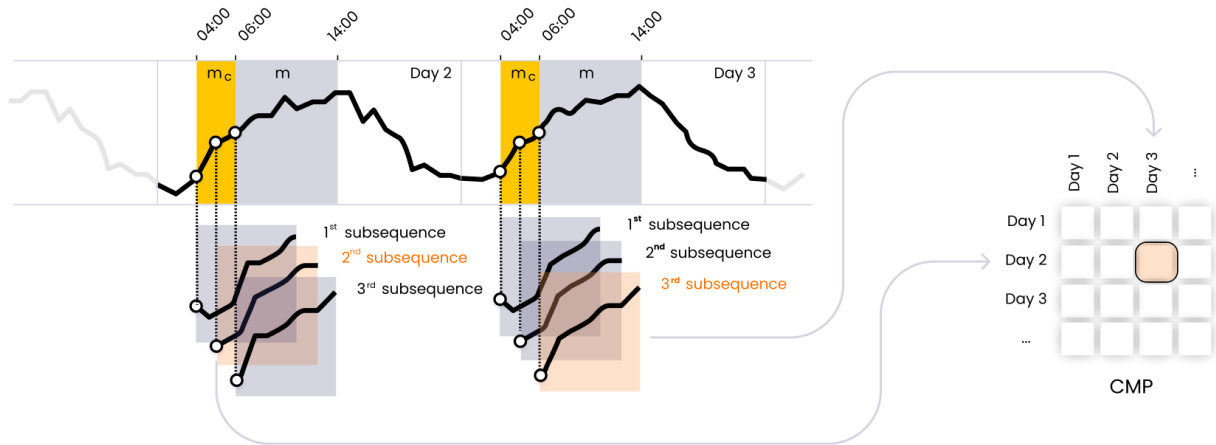
228 Eventually it is possible to define MP as the vector that stores the z-normalized Euclidean distance values between
 229 each subsequence $T_{i,m}$ and its nearest neighbour. Formally, $MP = [\min(D_1), \min(D_2), \dots, \min(D_{n-m+1})]$, where D_i is
 230 the *distance profile* corresponding to query $T_{i,m}$ and timeseries T . In other words, it can be generated by extracting the
 231 smallest value in each row/column of the *full distance matrix* \mathcal{M} . With reference to Figure 1, the MP is the column wise
 232 minimum over the entire *full distance matrix*, meaning that it finds the best matching T_1 subsequence for any subsequence

233 in T_2 . Of course, the construction of the *full distance matrix* is the most straightforward method but even the less
234 computational efficient, this is the reason why many algorithms have been proposed for the MP calculation to reduce time
235 and dimensionality complexity such as STAMP, STAMPI and STOMP based on MASS algorithm [60], approximated
236 AMPSA and AMP [61] and multidimensional mSTAMP [62].

237 MP represents a reference algorithm in the field of timeseries analytics and in particular among similarity-join-search
238 algorithms. However, performing a similarity-join-search considering all the possible subsequence pairs in a timeseries
239 is not always the most robust way to discover motifs and discords. In fact, when the domain of application constraints the
240 meaningfulness of the similarity search in a timeseries, it could be useful to retrieve the most similar pattern of each
241 subsequence also exploiting user's knowledge. The concept of Annotation Vector (AV) can be then used to introduce
242 domain knowledge in the process of motif and discord discovery [63], which allows to find results that respect user
243 defined constraints and produce robust results, closer to expectations of the analyst. AV is a meta timeseries used to
244 correct a-posteriori the values of the original MP, changing its shape and then manipulating the motif/discord search [64].
245 However, this method does not modify the MP calculation itself: *all-pairs-similarity-search* is always performed and then
246 a downstream processing is conducted.

247 As a possible solution, in [26] was introduced the concept of contextual matrix profile (CMP), defined as the minimum
248 over rectangular regions of the *full distance matrix* (see Figure 1), allowing then the search of the best matching
249 subsequences over T_1 and T_2 only considering a-priori determined set of timeseries segments and excluding others.

250 This allows to group data in the timeseries in a custom way comparing only portions of T_1 with portions of T_2 . The
251 CMP calculation is based on the definition of contexts m_c , intended as lapse of time, during which a subsequence of
252 length m may start. Figure 2 shows an example of context and subsequence length definition. Suppose that in an electrical
253 load timeseries T with $1h$ timestep an analyst wants to find for every day, included in the timeseries, the nearest neighbour
254 between the 3 subsequences of length 8h that could start at 04:00, 05:00 or 06:00 a.m., avoiding then similarity search in
255 other time intervals. To this aim the subsequence length must be set to $m = 8$, (grey regions in Figure 2) and the context
256 length $m_c = 3$ (yellow regions in Figure 2) and the resulting self-join CMP will have a row/column for each day. Each
257 point of the CMP would display the distance between the best matching 8h-long subsequences occurring in two different
258 days: lower the distance better the match and vice versa. While context is suitable to define a-priori the subsequence
259 matching in the timeseries, once the CMP is calculated it is even possible to further divide the CMP in fragments, so-
260 called *groups*. For instance, in the case of building energy consumption timeseries, it is possible to fragment the CMP by
261 keeping only weekdays or weekends in order to consider the existence of different occupancy-based operational schedules
262 in the discovery of motif and discords. Considering the main features and parameters (i.e., context and subsequence
263 lengths, groups) of the CMP algorithm, in the next section are introduced the contributions of this study and the research
264 gaps that this research aims to bridge.



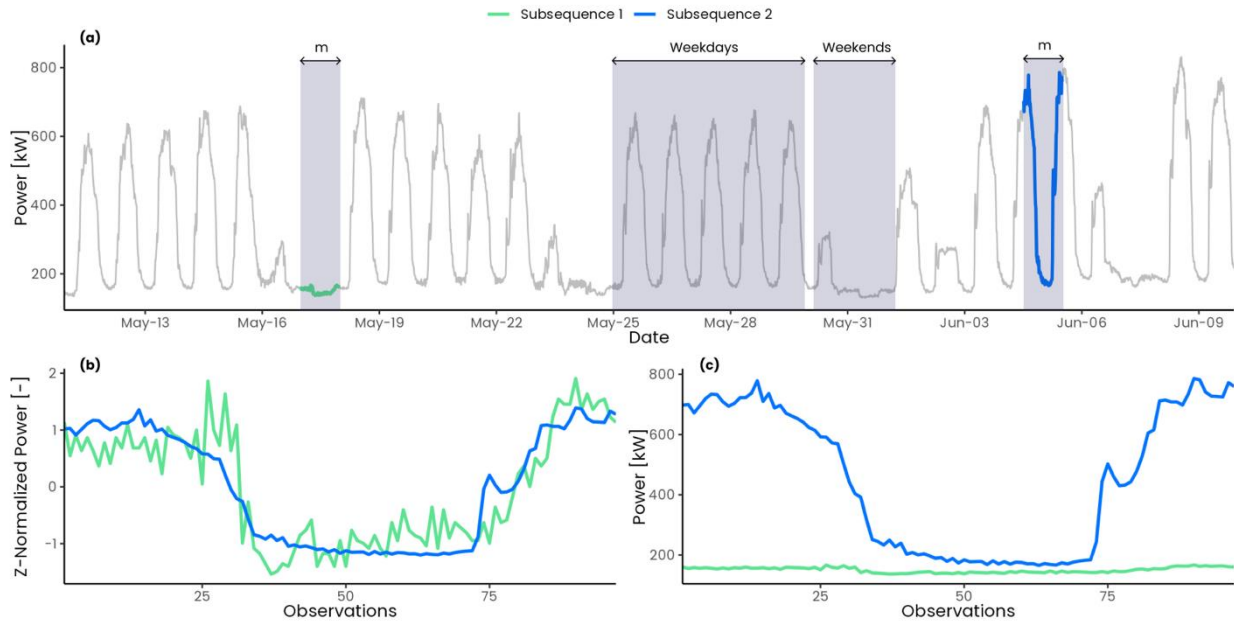
265

266 **Figure 2.** Focus on the definition of subsequence length m and contexts m_c in an electrical load timeseries T and their importance in the CMP
 267 calculation. The timeseries has one hour frequency and for each day a context of $m_c = 3$ (yellow) and a subsequence length of $m = 8$ (grey) are defined.
 268 When calculating the CMP, the nearest neighbour search is performed only on the subsequences that start within the defined context. The nearest
 269 neighbour subsequences between Day 2 and Day 3 are highlighted in orange. The Euclidean distance between the nearest neighbour subsequences is
 270 calculated and stored in the CMP shown on the right side.

271 1.3. Research gap and contribution of the paper

272 The MP algorithm has been successfully employed in different anomaly detection applications and researchers have
 273 proposed different MP implementations according to specific constraints and boundary conditions that are characteristic
 274 of the domain of interest.

275 In buildings, anomalous subsequences in energy consumption timeseries are defined as unexpected behaviours that
 276 result in subsequences with atypical shape and/or magnitude. The original MP algorithm, through z-score normalization
 277 of the subsequences, searches the nearest neighbour based on shape similarity by losing information related to the
 278 magnitude. Figure 3(a) shows a real electrical load timeseries for a non-residential building (i.e., university campus) in
 279 May and June. It is possible to observe how the electrical load changes significantly from weekdays to weekends when
 280 the load profile is almost flat. Applying the MP algorithm with a subsequence length of one day, the two subsequences
 281 highlighted respectively in blue and green are identified as nearest neighbours. As shown in Figure 3(b) under z-score
 282 normalization they are almost overlapping. However, Figure 3(c) shows that the not-normalized subsequences have very
 283 different amplitudes and are referred to distinct energy consumption patterns that are typical of weekdays and weekends
 284 respectively. This is a clear example of how the subsequence normalization could led to misleading results in terms of
 285 subsequence similarity.



286

287
288
289

Figure 3. Effect of z-score normalization on two electrical load timeseries subsequences (blue, green) of length 24h and a timestep of 15-min (96 observations per day): (a) full electrical load timeseries; (b) comparison between z-score normalized subsequences; (c) comparison between not normalized subsequences.

290

291

292

293

294

295

296

297

298

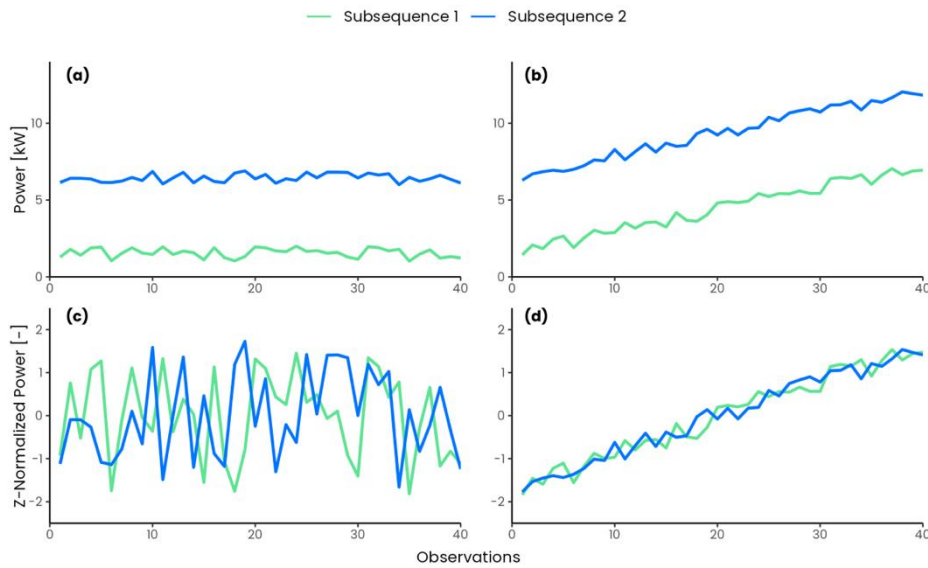
299

300

301

302

Z-score normalization not only minimizes magnitude effects in the research of motifs and discords but also tends to emphasize any fluctuation and noise in the timeseries. By comparing two relatively flat subsequences under z-score normalization the resulting z-normalized Euclidean distance is higher if compared to non-flat subsequences, leading to high values of MP in flat regions of the timeseries. In Figure 4 a comparison between two noisy synthetics timeseries is shown. In Figure 4(a) the two timeseries are relatively flat while in Figure 4(b) the two timeseries present a positive slope. In the first case, the normalization led to an amplification of the noise of the original subsequences (Figure 4(c)), resulting into a high Euclidean distance $d = 9.25$, while in the second case (Figure 4(d)), the Euclidean distance is much lower ($d = 1.5$). This issue has been largely analysed in [53], considering as possible solution to discard flat regions from the analysis, to change the subsequence length, or to smooth the timeseries. As a clear consequence of normalization, in an electrical load timeseries (Figure 3(a)), the MP method would identify the weekends as discords since they are typically associated to flat profiles compared to weekdays subsequences. This is a critical aspect when dealing with energy consumption timeseries of buildings that, by their nature, often present a strong pattern periodicity among working and not working days.



303

304

Figure 4. Effect of z-score normalization on noisy flat subsequences (a-c) and on noisy subsequences with positive slope (b-d). Adapted from [54]

305

306

307

308

309

In this perspective, comparing subsequences belonging to different energy patterns could be not so useful, making the introduction of domain knowledge, extremely important to drive the motif and discord discovery in timeseries. As previously discussed, a possible solution has been proposed by [26] introducing the CMP algorithm. In some applications CMP can be useful, since it allows to exclude some regions, or to split subsequences into different groups, and then perform the similarity search.

310

311

312

313

314

315

The application of CMP, coupled with domain knowledge, can be effectively exploited for the accurate discovery of anomalies in energy consumption timeseries, reducing energy wastes and enhancing energy management in buildings. Domain expertise should be used for an aware setting of the parameters of CMP algorithm (i.e., subsequence length, context length and groups) to avoid undesirable results and to make the ADD process as robust as possible. However, this is not a task that can be easily generalizable and still requires research effort for applications in the energy and building field.

316

317

To this purpose, according to the previous literature review and to the reasoning on MP and CMP methods, this paper intends to contribute as follows:

318

319

320

321

322

323

324

325

326

327

328

329

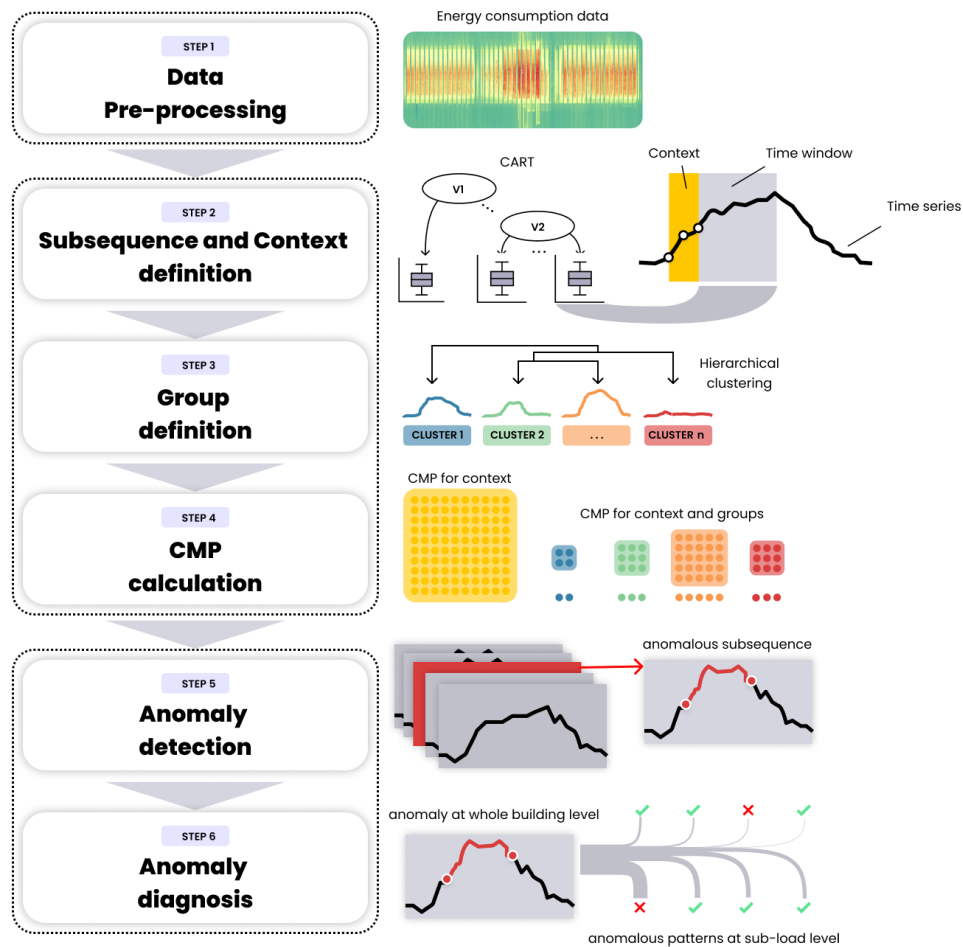
330

- 1) Develop a contextual ADD methodology by introducing a CMP-based process that employs supervised and unsupervised algorithms (i.e., clustering and decision tree) for the identification of parameters such as subsequence length, groups and contexts, systematically introducing domain expertise in the whole analytical framework.
- 2) Develop an ADD methodology based on the concept of *similarity-join* without any loss of information in the discovery of motifs and discords (e.g., avoiding dimensionality reduction and transformation of timeseries as in SAX algorithm) that is flexible enough to identify anomalous/typical patterns even if they are not perfectly aligned in time (i.e., time-tolerant analysis according to the length of the contexts).
- 3) Perform a discord discovery that is sensitive to both shape and magnitude of energy consumption patterns. To this purpose normalization aspects are considered upstream the computation of subsequence distances in CMP allowing a knowledge-based comparison between energy patterns in the analysed time series.
- 4) Introduce a robust severity score, based on different statistical methods and domain knowledge, that allows potential meter-level anomalies to be identified, ranked and diagnosed, within particular boundary conditions.

331 The rest of the paper is organized as follows. Section 2 provides the description of the methodological framework
 332 introduced. Section 3 presents the case study and the obtained ADD results while Section 4 critically discusses the
 333 outcomes and includes the concluding remarks.

334 2. Methodology

335 In this section the methodological framework employed to conduct ADD in energy consumption timeseries of
 336 buildings is presented. The framework is based on the application of the CMP coupled with unsupervised and supervised
 337 data analytics techniques, such as cluster analysis and classification and regression trees (CART), to perform the tuning
 338 of CMP parameters (subsequence length, context length, groups). The procedure, depicted in Figure 5, unfolds over of
 339 six steps described in detail in the following paragraphs.



340

341 **Figure 5.** Graphical representation of the methodological framework.

342 2.1. Pre-processing.

343 The first step consists in data pre-processing and is performed through univariate statistical approaches. Punctual
 344 statistical outliers and inconsistencies (e.g., negative values of electrical load) in the timeseries are identified and removed
 345 by means of boxplot analysis. Then all the missing values are replaced by means of linear interpolation.

346 2.2. Contextual matrix profile.

347 By following the methodological framework in Figure 5, the application of the CMP method goes through the
 348 following steps: subsequence and context definition (step 2), group definition (step 3) and CMP calculation (step 4).

349 Within the energy consumption timeseries of buildings, different characteristic sub-daily load patterns (e.g., base load,
 350 peak load, ramp-up and ramp-down period) can be identified and set as relevant subsequences to perform similarity-join-
 351 search by means of CMP algorithm. Such subsequences, with their relative lengths, can be statistically defined or inferred
 352 from the typical building operational and occupational schedule [65,66] The methodology proposed in this paper
 353 identifies sub-daily subsequences (t_w) through the implementation of recursive partitioning Classification and Regression
 354 Tree (CART) [15]. Starting from the root node (that contains all the available observations) this method proceeds through
 355 a set of binary decisions to split the instances in purer subsets (nodes) in a forward stepwise fashion, minimizing at each
 356 decision step the impurity/variance of each node [15,27,67], yielding local optimum [68] once a stopping condition is
 357 satisfied. To address this task, the decision tree was set as a regression model using the electrical load values as numeric
 358 target attribute and the hour of the day as explanatory attribute. This allows to identify, through a cost complexity process,
 359 a set of leaf nodes that group together hours with similar electrical load values in n non-overlapping sub-daily
 360 subsequences in the time domain $t_{w,i} \in \{t_{w,1}, t_{w,2}, \dots, t_{w,n}\}$. As a reference, the choice to not setting the time variable as
 361 a categorical one, allows CART to not indiscriminately group together electrical load values pertaining to the early
 362 morning with those referring to the night.

363 Thus, for each leaf node it is possible to extract a start and end-time of the subsequence and its length in terms of time
 364 duration ($m_i = t_{w,i}$). Once, all the interesting subsequences are identified, for each of them a context period is set to
 365 drive their *similarity-join-search* in the analysed time series of energy consumption. As reported in Figure 2 the context
 366 is defined as the time lapse during which can be located the start-time of a subsequence. This allows to investigate
 367 similarity of subsequences of the same length but not necessarily aligned in the time domain. Specifically, the context i
 368 ends at the beginning of the subsequence $t_{w,i}$, assuming as reference start time of $t_{w,i}$ the one extracted from the leaf node
 369 of the regression tree. Instead, the context length is set equal to the half of the smallest subsequence length
 370 ($m_c = \min(t_{w,i}) / 2$) identified through the regression tree. For instance, if the smallest subsequence is two hours long
 371 $t_w = 2h$ (e.g., from 6:00 to 8:00 a.m.) the context is defined as one hour long $m_c = 1h$ (i.e., from 5:00 to 6:00 a.m.).

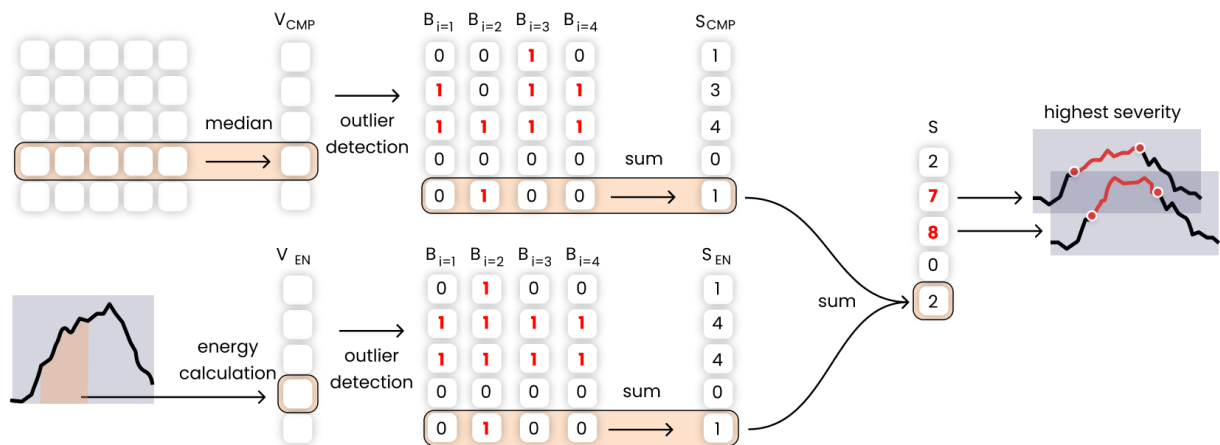
372 The following step is then aimed to define groups in CMP. This step is based on daily load profile clustering. An
 373 expert-based approach is firstly applied to split working days, holidays/not working days and partially working days (e.g.,
 374 Saturdays). Then a nested hierarchical clustering analysis is performed on the daily profiles of working days to better
 375 characterize them. The hierarchical clustering generates non-overlapping groups by splitting instances based on a
 376 geometrical distance measure and each cluster can be further divided into subclusters and so on, creating a tree structure
 377 (i.e., dendrogram). The grouping of daily load profiles in representative clusters of specific energy consumption patterns,
 378 allows the CMP results to be evaluated in groups with high internal similarity and to make the anomaly detection process
 379 more robust.

380 Once contexts and subsequence lengths are defined, the CMP is calculated for each context under the hypothesis of
 381 not-normalized Euclidean distance. According to previous definitions, each day has n not overlapping contexts and the
 382 resulting CMP contains one row/column for each day. Eventually, the obtained overall CMP matrix is split into different
 383 sub-matrices according to the membership of each day to a specific group identified by means of the cluster analysis.
 384 After this, the anomaly detection process can be carried out for each context in each group.

385 2.3. Anomaly detection

386 The anomaly detection process, shown in Figure 6, is carried out in each group after performing the similarity-join-
 387 search of subsequences in their relative context, and consists in the definition of a severity score through the use of four
 388 statistical univariate outlier detection methods. The considered four methods are applied on two vectors i.e., V_{CMP} and
 389 V_{EN} . The first vector V_{CMP} contains the median value of each row/column of the CMP i.e., the median of the Euclidean
 390 distances between a specific subsequence and its nearest neighbours in the associated context considering only the days
 391 in the same group (i.e., cluster of daily load profiles). The second vector V_{EN} , contains for each subsequence the
 392 corresponding energy consumption. In this way it is possible to label as anomalous, subsequences with Euclidean distance
 393 and energy consumption values that are out-of-range respect to the domain that is considered normal for both V_{CMP} and
 394 V_{EN} . With reference to Figure 6, as a first step the CMP matrix is reduced into a vector V_{CMP} by calculating the median
 395 of each row/column, then the four outlier detection methods are applied to each component of V_{CMP} producing 4 new
 396 vectors that define whether an element is anomalous or not in a boolean form $B_i = \{0,1\}$. Then the severity score
 397 associated to a subsequence is calculated counting the number of positive detections $S_{CMP} = \sum_{i=1}^4 B_i$. To make more
 398 robust the anomaly detection process and consider only positive anomalies (e.g., anomalies that resulted into an over-
 399 consumption of energy) the energy consumption for each subsequence is calculated and stored in a vector V_{EN} which
 400 undergoes to the same process described before: outlier methods are applied and then the severity is calculated. With
 401 reference to Figure 6, given $n = 4$ outlier detection methods the severity vector S_{CMP} and S_{EN} range from a minimum of
 402 0 (four methods out of four did not tag the subsequence as out-of-range observation) to a severity of 4 (four methods out
 403 of four found tagged the subsequence as out-of-range observation). By summing S_{CMP} and S_{EN} the resulting overall
 404 severity S ranges from 0 (four methods out of four did not tag the subsequence as out-of-range observation for both energy
 405 and Euclidean distance vector) to 8 (four methods out of four tagged the subsequence as out-of-range observation for
 406 both energy and Euclidean distance vector).

407 To avoid spurious alerts and reduce the number of feedbacks only the subsequences with severity 6-7-8 were
 408 considered and further analysed in the diagnosis phase.



409

410

Figure 6. Graphical representation of the anomaly detection phase.

411 The four univariate outlier detection methods are: (i) inter quartile range analysis, (ii) Z-score, (iii) elbow method and
 412 (iv) Generalized Extreme Studentized Deviate. The basic principles behind these statistical methods are introduced and
 413 described below.

414 *Inter quartile range analysis* defines as outlier any observation that falls below $Q_1 - 1.5 * IQR$ and above $Q_3 + 1.5 *$
 415 IQR where the interquartile range (IQR) is defined as the difference between the third quartile and the first quartile $IQR =$
 416 $Q_3 - Q_1$ of the population of observations. In this study only positive outliers are considered i.e., values above $Q_3 + 1.5 *$
 417 IQR .

418 *Z-score, or standard score*, is the output of data standardization process and is intended as the number of standard
 419 deviations by which an object, in a population of observations, is above or below the mean value. Observations above the
 420 mean have positive Z-scores, while those below the mean value have negative Z-scores. This method can be employed
 421 to identify extremal observations (i.e., outliers) by setting a z-score threshold value. For instance, define a Z-score
 422 threshold equal to 2 means that the all the observations that differ from the mean more than 2σ are tagged as outliers, and
 423 in the case of normally distributed observations they represent the 2.3% of the whole population.

424 *Elbow method* is a heuristic method that allows to find the elbow (or knee) of a curve. The implemented method
 425 employs the kneedle algorithm [69] to identify the point of maximum curvature and thus locate the elbow. By finding the
 426 elbow of a univariate vector ordered in descending values it is possible to identify two different regions: the region below
 427 the elbow and the one above the elbow, where observations are tagged as outliers.

428 *Generalized Extreme Studentized Deviate (GESD)*: is an iterative method that progressively evaluates the presence of
 429 outliers in a univariate timeseries $T = \{ t_1, t_2, \dots, t_n \}$ through a statistical test [8]. The method initialization requires
 430 ($i = 1$) a presumed number of outliers r and confidence interval α , then the following statistical test is performed:

- 431 • H_0 There are no outliers in the timeseries
- 432 • H_a There are up to r outliers in the timeseries

433 The hypotheses test is performed by calculating the R_i statistic and the critical value λ_i as follows:

$$434 \quad R_i = \frac{\max|t_i - \mu_T|}{\sigma_T} \quad \lambda_i = \frac{(n - i) * t_{p, n-i-1}}{\sqrt{(n - i - 1 + t_{p, n-i-1}^2) * (n - i + 1)}}$$

435 Where μ_T and σ_T denote sample mean and sample standard deviation of the timeseries, n is the timeseries length,
 436 $i = \{1, 2, \dots, r\}$ is the iteration number, $t_{p, \nu}$ is the 100p percentage point from the t distribution with ν degrees of
 437 freedom and $p = 1 - \frac{\alpha}{2(n-i+1)}$

438 2.4. Anomaly diagnosis phase

439 Once the anomaly at meter-level is identified (i.e., anomalous subsequence with at least a severity score S equal to 6),
 440 the diagnosis phase is enabled with the aim to identify if also sub-loads are anomalous in the same time interval and
 441 contribute to the diagnosis definition. The diagnosis process involves only those subsequences that resulted to be
 442 anomalous at meter-level and extend the analysis at sub-load level.

443 By keeping the same hyperparameter settings (i.e., contexts, subsequence length and group) as previously described,
 444 each sub-load timeseries undergoes through CMP calculation and anomaly detection process using the 4 outlier methods
 445 previously presented. In this way, for each subsequence of the sub-loads, it is possible to assess if the sub-load is
 446 anomalous by calculating the severity score that allows to rank the sub-loads accordingly (i.e., from low to high severity).
 447 The sub-loads that present severity scores in the range (6 - 8) are assumed as anomalous and contribute to the diagnosis
 448 of the anomaly detected at higher level. Otherwise, if no sub-load, for a specific anomalous subsequence at meter level,
 449 reaches an anomaly score at least of 6, the anomaly is considered as undiagnosed.

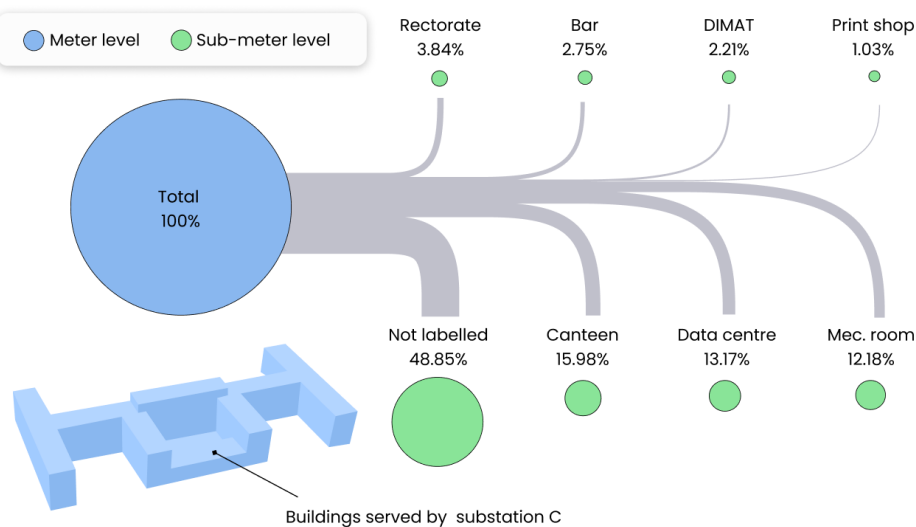
450 3. Results

451 The analysis was carried out using the R statistical software [70] for the pre-processing, regression analysis, cluster
 452 analysis and visualization, and Python [71] for the CMP calculation and anomaly detection process.

453 The presented methodology has been tested on the energy consumption timeseries of a MV/LV transformer cabin (i.e.,
 454 substation C) that serves part of the university campus of Politecnico di Torino located in the norther Italy. The monitoring
 455 infrastructure continuously provides both meter-level and sub-meter level measurements of the mean electrical load with
 456 a timestamp of 15 min. The hierarchical structure of the installed monitoring infrastructure is shown in Figure 7: the first
 457 level (blue) refers to the total electrical load of substation C, while the second level (green) shows the associated sub-
 458 loads. In addition, the breakdown in terms of average annual energy consumption was provided for each sub-load.

459 The substation serves a bar and a staff canteen that are at disposal of students and campus staff and, on average,
 460 represent the 2.75% and 15.98% of the yearly electrical energy consumption of substation C respectively. The university
 461 data centre accounts for the 13.17% of the total energy consumption. The administration offices (Rectorate) correspond
 462 to the 3.84% of energy consumption while the mathematics department (DIMAT) to the 2.21%. A relevant amount of
 463 energy consumption (12.18%) is related to the mechanical room. The equipment located in this room includes hot and
 464 chilled water circuits with their associated auxiliaries such as recirculation pumps. The chilled water is produced by two
 465 chillers with nominal electrical power of 220 kW and a rated cooling capacity of 1120 kW, and a reversible water-water
 466 heat pump, with nominal a power and cooling capacity of 165 kW and 590 kW, respectively.

467 The remaining energy consumption is aggregated under a single synthetic sub-load tagged as “Not labelled” as showed
 468 in Figure 7. It accounted for 48.85% of the yearly energy consumption of the substation and since it is not directly
 469 measured cannot be labelled, as the other sub-loads, respect to the energy service provided or building served. The authors
 470 tested the presented methodology on a one-year dataset that spans from 01/01/2019 to 12/31/2019. The year 2019 was
 471 selected, even if more recent data were available, mainly because the pandemic COVID-19 completely changed the
 472 operational patterns due to the closure of the university campus from February 2020. The raw dataset includes 35040
 473 observations with 15-min timestep. The dataset is characterised by a missing value ratio of less than 0.1%. Inconsistencies
 474 were removed and missing values replaced through linear interpolation following the first step of the methodology.



475

476

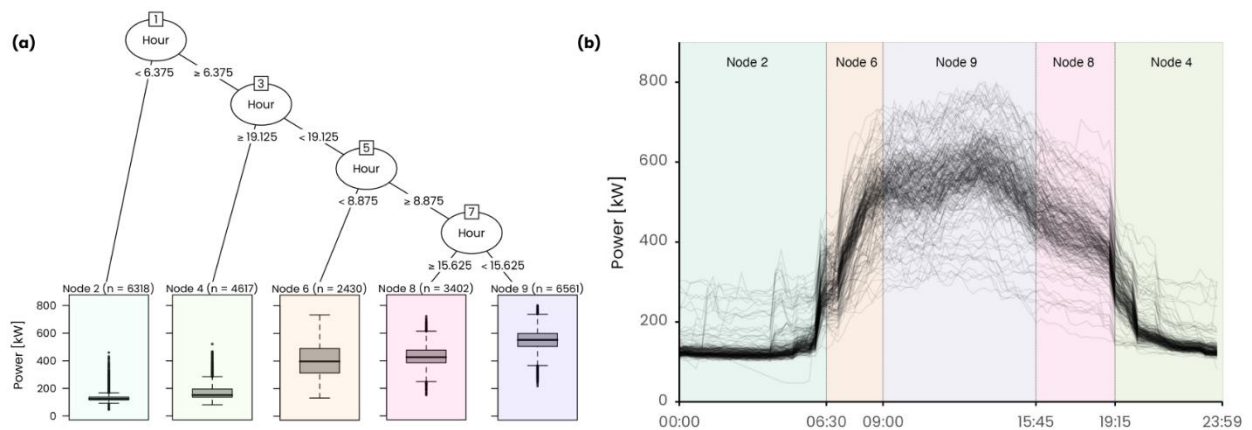
Figure 7. Hierarchical structure of the electrical load database under study.

477 3.1. Contextual matrix profile results

478 To identify the subsequences to be used for the *similarity-join-search* through CMP, not overlapping time periods
 479 were evaluated through CART algorithm using the meter-level 15-min electrical load as target variable, and hour of the
 480 day as numeric predictive variable. To preserve the accuracy of the model in the leaf nodes, pertaining the operation hours
 481 of the energy systems (e.g., from 07:00 to 19:00), only working days were considered and days with a low standard
 482 deviation of the electricity demand (e.g., Sundays, holidays) were excluded.

483 In fact, for days with a low standard deviation of the electricity demand, that are almost flat, the ADD results would
 484 be less sensitive to the subsequence definition. On the other hand, for days with high variability of electricity demand it
 485 is particularly important to extract the best patterns that can properly characterize the load variation over time. In particular
 486 the CART is extremely effective to identify leaf nodes (in this case associated with subsequence length) with low variable
 487 load and others that instead include all the load variability (i.e., morning ramp-up, evening ramp-down). Indeed, even
 488 though the accuracy is not the highest achievable, the results that can be obtained are exactly the ones of interest.

489 The adopted stopping criterion, to avoid overfitting of the regression tree, was based on the minimum number of
 490 objects in a child node in order to identify subsequences with a length of at least 2.5 hours. The tree was subjected to
 491 cross validation and cost-complexity pruning, resulting in the five-leaf tree shown in Figure 8(a). For completeness the
 492 leaf nodes resulting from the tree are represented as sub-daily subsequences also in the Figure 8(b), where are reported
 493 with their duration on a daily scale.



494

495 **Figure 8.** (a) Identification of sub-daily subsequence by means of the CART algorithm (b) graphical representation of the daily load profiles and the
 496 identified subsequences on daily scale.

497 The model was able to effectively separate the night hours (subsequences 1 and 5 extracted from leaf node 2 and 4
 498 respectively of the tree in Figure 8(a)) from the diurnal operation and was able to identify the ramp-up (subsequence 2
 499 extracted from leaf node 6 of the tree in Figure 8(a)), mid-day operation (subsequence 3 extracted from leaf node 9 of the
 500 tree in Figure 8(a)) and ramp-down (subsequence 4 extracted from leaf node 8 of the tree in Figure 8(a)). The smallest
 501 subsequence lasts 2.5 hours. Following the methodological process, once the reference sub-daily subsequences are
 502 evaluated the contexts for the definition of the CMP can be identified. In this study the context length was defined as the
 503 half of the smallest subsequence length identified ($m_c = 2.5h/2 \cong 1h$). The outcome of this preliminary step led to
 504 the definition of the lengths of 5 subsequences and 5 contexts (i.e., the number of observations in each subsequence and
 505 context) with their durations for the CMP calculation (Table 1). For instance, the third sub-daily subsequence, that ranges
 506 in [09:00 – 15:45) (as extracted from the tree in Figure 8(a)), has a length of 27 observations (i.e., $27 * 15 \text{ min} =$

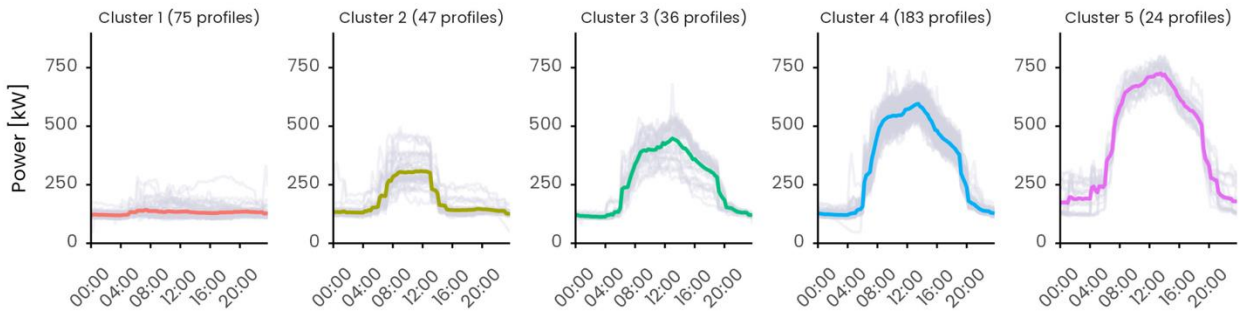
507 6 h 45 min) and in the similarity-join-search with CMP all the subsequences of length 27 that starts in the 1h context
 508 [08:15 – 09:15) are compared to identify pairs of nearest neighbours among days in the same group.

509 **Table 1.** Context and subsequence lengths.

Subsequence				Context			
ID	Time interval	Duration	Length	ID	Time interval	Duration	Length
$t_{w,1} = m_1$	[00:00 - 06:30)	6 h 30 min	26	$m_{c,1}$	[23:15 - 00:15)	1 h	4
$t_{w,2} = m_2$	[06:30 - 09:00)	2 h 30 min	10	$m_{c,2}$	[05:45 - 06:45)	1 h	4
$t_{w,3} = m_3$	[09:00 - 15:45)	6 h 45 min	27	$m_{c,3}$	[08:15 - 09:15)	1 h	4
$t_{w,4} = m_4$	[15:45 - 19:15)	3 h 30 min	14	$m_{c,4}$	[15:00 - 16:00)	1 h	4
$t_{w,5} = m_5$	[19:15 - 24:00)	4 h 45 min	19	$m_{c,5}$	[18:30 - 19:30)	1 h	4

510
 511 The second step consists in the definition of the groups and was performed using a semi-supervised approach applied
 512 on the set of 365 daily load profiles available in the dataset. Firstly the 75 daily load profiles corresponding to public
 513 holidays, university closures and Sundays were extracted and grouped together in a cluster labelled as Cluster 1. Secondly,
 514 the 47 daily load profiles of half-working days and Saturdays were extracted and assigned to Cluster 2. The remaining
 515 243 daily load profiles, corresponding to working days, were organized into a MxN matrix 243x96 where M is the number
 516 of days considered and N the number of observations per day (i.e., 96 observations considering a 15-min timestep). Then
 517 a hierarchical clustering algorithm with Ward method was implemented on the not normalized daily load profiles. The
 518 silhouette index, implemented in the package NbClust [72], was used to search the optimal number of clusters in a range
 519 between 2 and 6. Three clusters were identified as the optimal partition and were labelled as follows: Cluster 3 (36
 520 profiles), Cluster 4 (183 profiles), Cluster 5 (24 profiles).

521 Figure 9 shows the five final clusters identified: the grey lines represent the daily load profiles belonging to the cluster
 522 while the coloured line represents the cluster centroids. The clustering process led to well-defined groups of load profiles
 523 each one representing a typical load condition useful to split the CMP for a given context into homogeneous groups for
 524 conducting the anomaly detection process.

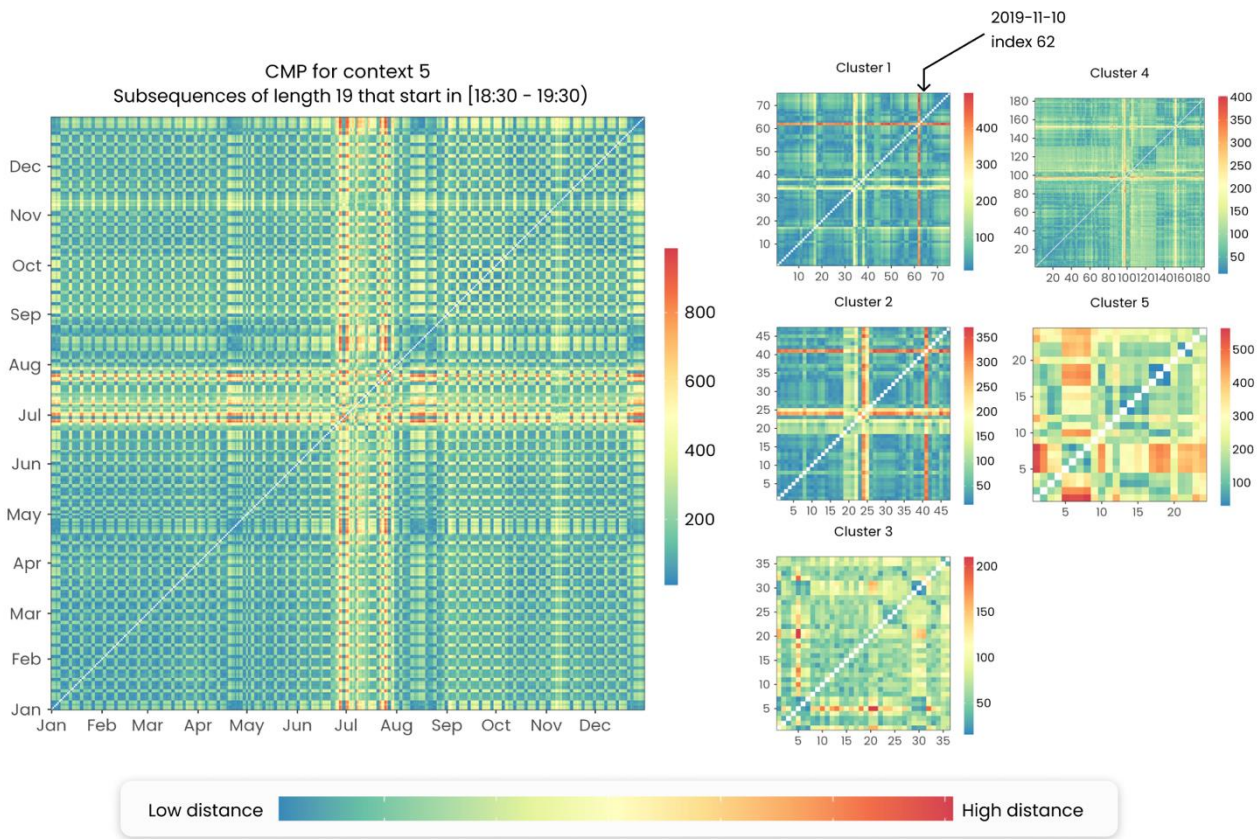


525
 526 **Figure 9.** Clusters of daily electrical load profiles with the corresponding centroids.

527 The CMP was then calculated by self-joining the data for each of the 5 contexts using the Euclidean distance between
 528 not normalized subsequences. The calculation was performed using the open-source Python code [26] implementing the
 529 Series Distance Matrix framework to calculate the CMP. As a representative example in Figure 10 is reported the CMP
 530 for context 5. Since the dataset contains 365 days and a specific context occurs one time per day the resulting CMP is a
 531 365x365 symmetric matrix where each value represents the Euclidean distance between the best matching subsequences

532 among two days considering all the possible subsequences that start during the same context. The higher the distance
 533 value the higher the dissimilarity between a pair of nearest neighbours.

534 The overall CMP (on the left of Figure 10) shows a weekly regularity: there are typically 5 days with the same
 535 behaviour (green) followed by two days with different behaviour (yellow). Moreover, a change of typical patterns during
 536 summer can be observed, especially during July and August, corresponding to holidays and summer closures of the
 537 facilities of university campus. By further splitting the CMP in the previously defined groups (i.e., clusters), days expected
 538 to behave in a similar manner are grouped together to perform a more robust inspection of anomalous patterns. For
 539 instance, the day 10th of November 2019 (highlighted in Figure 10 at index 62 of cluster 1 sub-matrix of CMP)
 540 to be remarkably different from all the other days in the same cluster, which is not so evident by only visualizing the full
 541 CMP. This highlights the impact that the group definition has, after the CMP computation, for the identification of
 542 contextual anomalies, given the relative importance that a pattern could have according to the group in which it is
 543 compared with its nearest neighbours. As a reference, performing clustering after the definition of the subsequences (and
 544 so contexts) makes it possible to extracted more general patterns (i.e., occurrence time of an anomaly) that can be
 545 compared among groups regardless to their composition (even in the case they change over time).



546

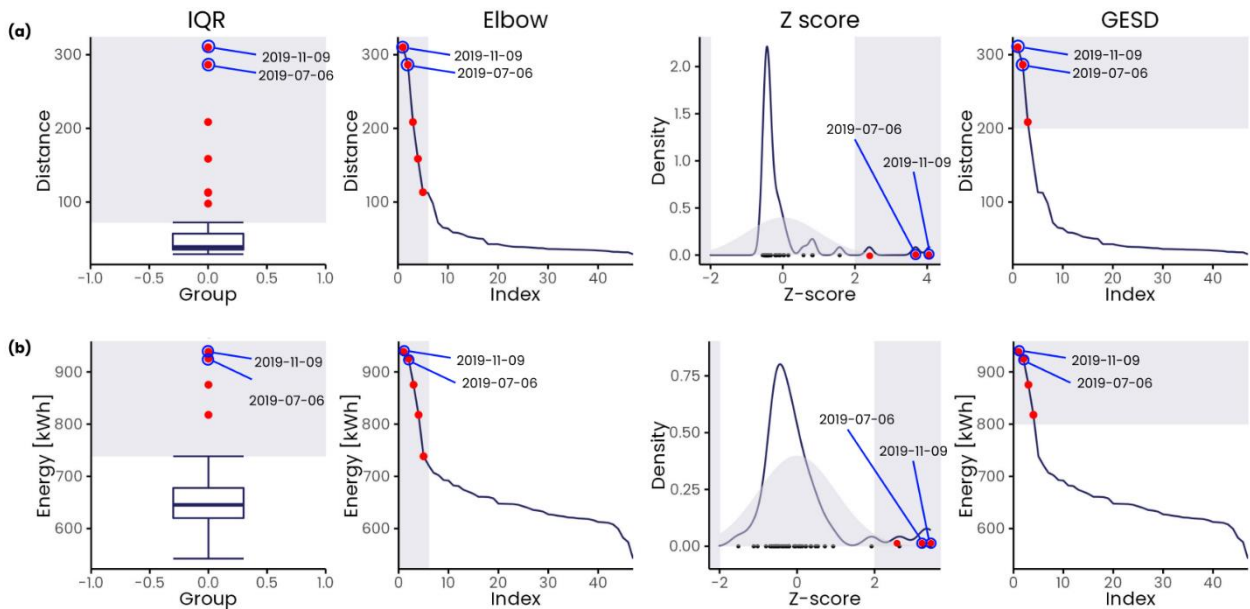
547 **Figure 10.** On the left the CMP result is reported for the similarity-join search of subsequences that start in context 5. Each value of the matrix shows
 548 the Euclidean distance between the best matching subsequences. Lower the distance better the match. The CMP is divided on the right side into
 549 submatrices according to the membership of each day to one of the predefined groups (i.e., clusters).

550 3.2. Anomaly detection results

551 The anomaly detection is performed among the subsequences that start in the same context considering a group at
 552 time. For each context-group CMP matrix, median values of columns/rows of the CMP and the energy consumption
 553 related to the subsequences are stored in the vector V_{CMP} and V_{EN} respectively, and the four univariate outlier detection

554 methods are applied. The methods are tuned as follows: the IQR method is calibrated to consider only the positive outliers
 555 over $1.5 * IQR + Q_3$, the z-score method considers only the positive observations over a value of 2, GESD method is
 556 initialized with presumed number of outliers $r = 10$ and confidence interval $\alpha = 0.05$ and eventually, the elbow method,
 557 since it is a pure graphical method, considers as outliers all the observations above the evaluated knee point. Each method
 558 defines whether a subsequence is anomalous or not in a Boolean form $B_i = \{0,1\}$. Then, the severity is obtained summing
 559 up the number of positive detections. By summing the two resulting severity vectors S_{CMP} and S_{EN} an overall severity S
 560 ranging from 0 to 8 was calculated to robustly rank anomalies from the most to the least severe one.

561 Figure 11 reports the anomaly detection results obtained for subsequences pertaining the context 5 considering days
 562 in cluster 2. In detail, Figure 11(a) shows the case of the outlier methods applied on the vector of the median Euclidean
 563 distances while the Figure 11(b) on the vector of the energy consumption. It is possible to easily verify that two days, 9th
 564 of November 2019 and 6th of July 2019, were detected as anomalous by all the methods in both the Euclidean distance
 565 and energy consumption vector and this resulted into an overall severity score of 8.



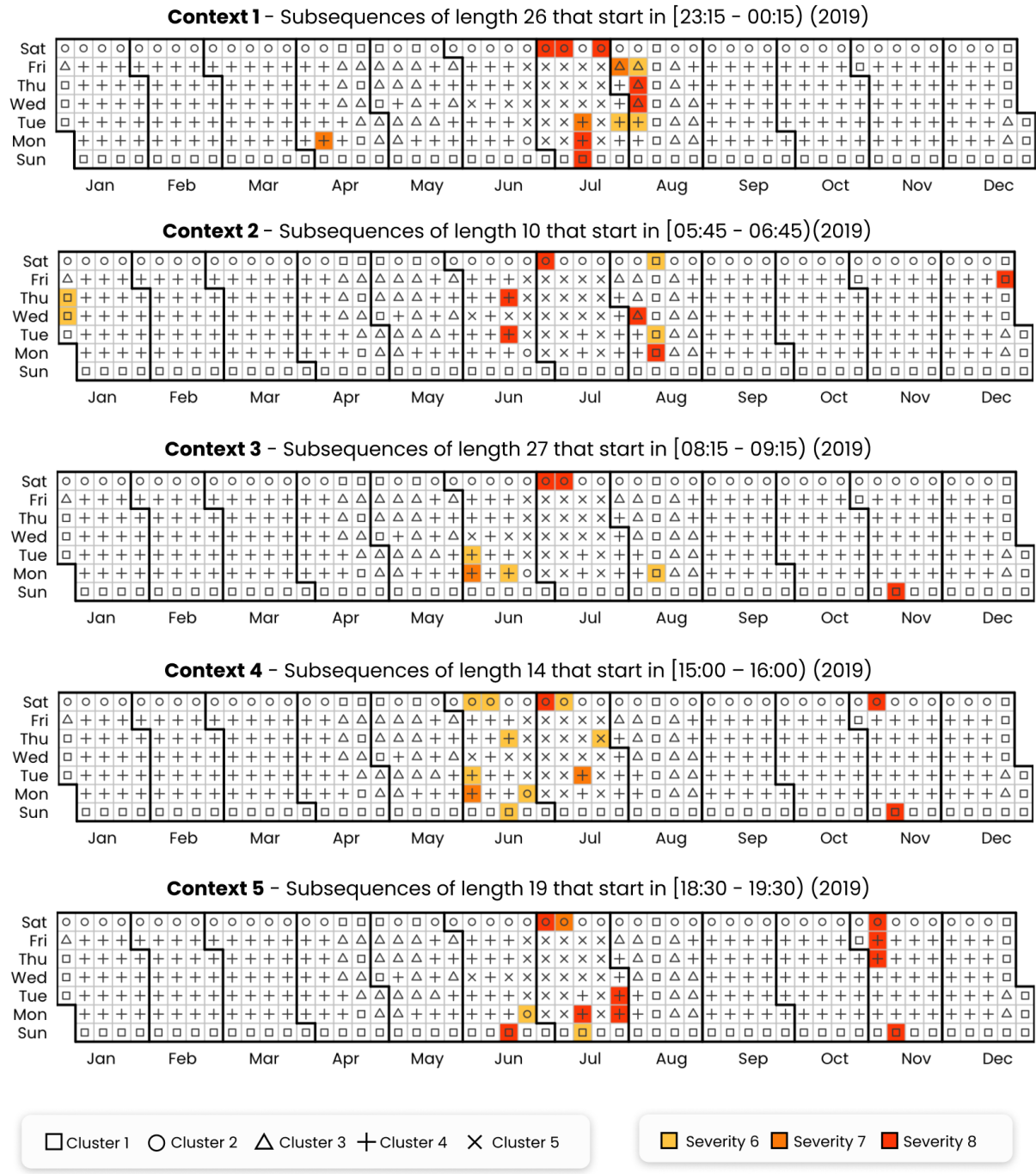
566

567 **Figure 11.** Example of severity calculation on (a) median Euclidean distance vector and (b) energy consumption vector using the four anomaly detection
 568 methods respectively IQR method, Z-score method, elbow method and GESD.

569 To reduce the number of spurious alerts, only the severities 6-7-8 were considered as relevant, resulting in 55
 570 anomalous subsequences detected: 20 with severity 6, 7 with severity 7 and 28 with severity 8. A high severity denotes a
 571 significant difference in terms of shape and magnitude of a subsequence within the relative group and context.

572 Results are summarized in Figure 12 through a calendar plot which shows anomalies over the year for the 5 contexts
 573 considered in the CMP *similarity-join-search*. According to Figure 12, anomalies resulted to be more frequent at the
 574 beginning and end of the day: 13 anomalous subsequences started in the 1st context, 10 in the 2nd context, 7 in the 3rd
 575 context, 13 in the 4th context and 12 in the 5th and last context. Moreover, the calendar plot clearly shows as 44 anomalies
 576 out of 55 detected, occurred during summer, from June to August, compared to the rest of the year where only 11
 577 anomalous subsequences were detected.

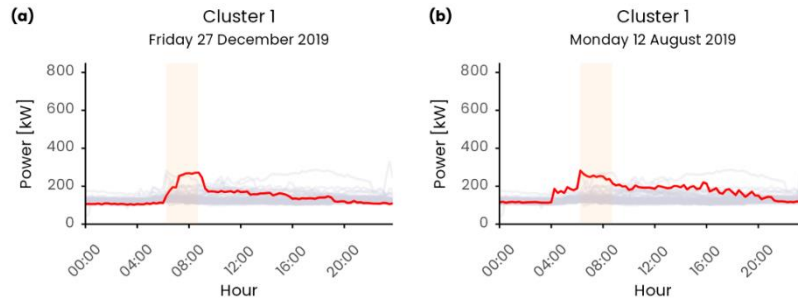
578 The definition of the 5 contexts and related subsequences, with diversified lengths among contexts, allowed the ADD
 579 process to be effective in finding anomalous subsequences that are limited to one context during the day (i.e., spot
 580 anomalies) and others that instead persisted among subsequent contexts and last until the end of the day (i.e., persistent
 581 anomalies).



582

583 **Figure 12.** Calendar plot reporting the occurrence of the detected anomalies through the whole year dataset. Each row corresponds to a different context.
 584 For each day, the colour fill, represents the severity of the anomaly while the symbol represents the cluster membership.

585 An example of a spot anomaly is the subsequence starting at 05:45 in the second context of Friday 27th of December
 586 2019 that was tagged as anomalous with a severity score of 8. Referring to Figure 13(a) this is a holiday belonging to
 587 cluster 1 where a flat profile is expected to be found. However, a rise of the electrical load after 6:00 and an abrupt switch
 588 off at 9:00 was detected, resulting in a deviation from the centroid of about 285 kWh. The same pattern, shown in
 589 Figure 13 (b), was detected during summer season in the subsequence starting at starting at 05:45 of Monday 12th of
 590 August 2019, when despite the university campus was closed for the summer break, an abnormal increase in electrical
 591 load in the second context was detected resulting into a deviation of more than 280 kWh from the cluster centroid. These
 592 two examples are symptoms of a wrong schedule of the energy systems pertaining to the substation C.



593

594
595
596
597

Figure 13. Anomalous subsequences identified as spot anomalies, respectively: (a) Monday 12th August 2019 (context 2, cluster 1); (b) Friday 27th December 2019 (context 2, cluster 1). The daily load profile which includes the anomalous subsequence is represented with the red line while the grey lines correspond to the daily load profiles included in the relative cluster. The vertical orange band denotes the anomalous subsequence while the band width is related to the subsequence length.

598

599

600

601

602

603

604

605

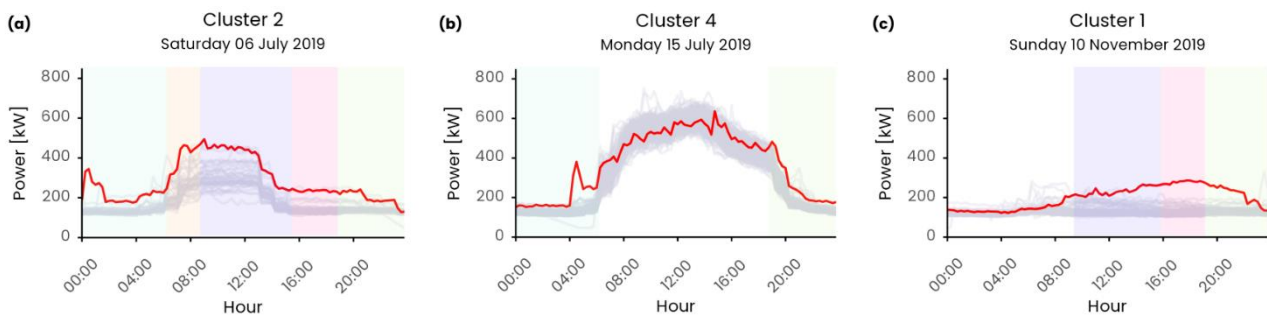
606

607

608

609

An example of persisting anomaly is Saturday 6th of July 2019 that presents anomalous subsequences of severity 8 for the whole day. With reference to Figure 14 (a), starting from an unexpected peak during night hours, the load profile was anomalous also during the following periods of day by keeping an average offset of almost 80 kW compared with the cluster centroid, leading to an overall deviation compared to the cluster centroid of about 2490 kWh at the end of the day. This behaviour is a symptom of energy systems running under unusual conditions, that may be related to a fault, a wrong schedule, or an exceptional outdoor boundary condition (e.g., high external temperature). A similar behaviour was observed on Monday the 15th of July 2019 where a wrong schedule of start-up and switch-off resulted in anomalies starting at 00:00 in contexts 1 and at 19:15 in context 5 leading to a deviation of about 800 kWh compared to the average cluster energy consumption, see Figure 14 (b). Another example of persisting anomaly is shown Figure 14 (c) where during winter season on Sunday the 10th of November 2019 an anomalous energy consumption involving subsequences that respectively start in contexts 3-4-5, is detected, leading to a deviation of about 1620 kWh from the reference behaviour of cluster centroid.



610

611

612

613

614

Figure 14. Anomalous daily load profiles identified as persisting anomalies, respectively: (a) Saturday 6th July 2019 (context 1, 2, 3, 4, 5 - cluster 2); (b) Monday 15th July 2019 (context 1, 5 - cluster 4); (c) Sunday 10th November 2019 (context 3, 4, 5 - cluster 1). The anomalous load profile is represented with the red line while the grey lines correspond to the load profiles contained in the relative cluster. The vertical band denotes the anomalous subsequences, the band colours are referred to the contexts considered, while the band widths are related to the subsequence lengths.

615

3.3. Anomaly diagnosis results

616

617

618

619

620

The diagnosis process aims to assess if during anomalous subsequences at meter-level also anomalies at sub-load level occurred, in order to provide an explanation to the detected infrequent energy consumption patterns. A sub-load with high severity score (i.e., higher than 6) is considered in the diagnosis of meter-level anomalies. The 55 anomalies discovered at meter-level are investigated and results are presented in the following. In particular, for 8 out of 55 anomalous subsequences, any sub-load does not reach a severity score at least of 6. It means that those anomalies are considered as

621 undiagnosed given that it is not possible to clearly recognize abnormal patterns among sub-loads according to the
 622 developed severity score. However, for those subsequences, the “Mechanical room” resulted to be the sub-load with the
 623 highest severity score (ranging between 3 and 5) in any case supporting the end user in the process of interpreting the
 624 detected anomalous energy consumption pattern. For what concerns the 47 remaining anomalies, 30 times are associated
 625 with the occurrence of a single anomalous sub-load (i.e., with severity score at least of 6), 12 times with two anomalous
 626 sub-loads, and 5 times with three anomalous sub-loads. Additional results are reported in Table 2. The table reports the
 627 number of times in which a sub-load or group of sub-loads was identified as a potential root cause to the meter-level
 628 anomaly due to severity scores at least of 6. “Not labelled” and “Mechanical room” are most frequently included in the
 629 lists of sub-loads contributing to the diagnosis of the detected anomalies at meter level, occurring respectively 24 times
 630 and 26 times (also considering the co-occurrence with other sub-loads). Other sub-loads, such as “Canteen”, “Data centre”
 631 were found to be anomalous 4 times while the “Print shop” exhibited an abnormal pattern 11 times. It is worth to note,
 632 that the severity score is formulated to reveal the occurrence of infrequent shapes and magnitudes of the subsequences by
 633 means of analysis that are internal to the timeseries of the single sub-load. It means that, even though some sub-loads
 634 have a relatively low impact in terms of yearly energy consumption on the total one of the whole substation C, the
 635 diagnosis process is able to highlight if they are normally behaving during an anomalous subsequence at meter-level
 636 regardless to their contribution to the total energy consumption.

637

Table 2. Summary of the diagnosis results.

Number of occurrences	Sub-loads contributing to the diagnosis
16	Not labelled
11	Mechanical room
5	Not labelled, Mechanical room
4	Mechanical room, Print shop
3	Not labelled, Mechanical room, Print shop
1	Canteen
1	Data centre
1	Print shop
1	Canteen, Print shop
1	Mechanical room, Data centre
1	Canteen, Data centre
1	Canteen, Mechanical room, Print shop
1	Data centre, Mechanical room, Print shop

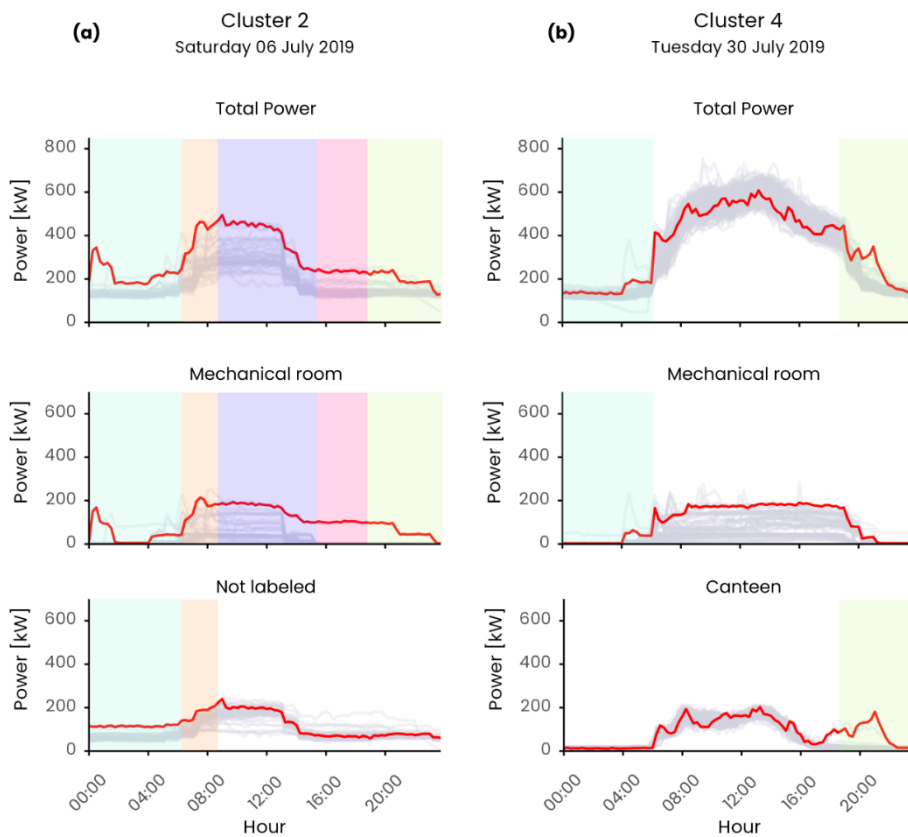
638

639 Figure 15 reports two representative examples in which the diagnosis approach was able to clearly identify and isolate
 640 the anomalous sub-loads responsible for the anomaly detected at meter-level. Figure 15(a) shows the load profiles of
 641 Saturday 6th of July 2019 during which both the “Mechanical room” and “Not labelled” sub-load contributed to the
 642 anomalous trend detected on the total electrical load of substation C, starting in contexts 1-2-3-4-5. For subsequences
 643 starting in contexts 1-2 the “Mechanical room” and “Not labelled” are both responsible for the anomaly detected. In detail,
 644 to the “Mechanical room” sub-load, is associated a severity score of 8 in both contexts while the “Not labelled” sub-load
 645 has a severity score of 8 for the subsequence starting in context 1 and severity 6 for the one starting in context 2. The
 646 “Mechanical room” was characterized by an unexpected switch-on/off during night hours (00:00 to 01:30) followed by a
 647 second switch-on at 04:00 that, superposed to the infrequent high night load of the “Not labelled”, led to an abnormal
 648 behaviour at meter-level. Then, for the remaining subsequences starting in contexts 3-4-5 only the “Mechanical room”
 649 was identified as responsible with a severity of 8 while the “Not labelled” load showed no impact at meter-level (i.e., very

650 low severity, less than 1). Figure 15 shows that the anomalous load of “Mechanical room” after 12:00 led to a deviation
 651 of the “Total Power” load at meter-level until the end of the day.

652 Figure 15(b) shows the load profiles of Tuesday 30th of July 2019 in which a single sub-load contributed to the meter-
 653 level anomaly detected. An early morning start-up at 4:00 of the “Mechanical room” identified by a severity 7 contributed
 654 to the “Total power” anomaly of severity 6 starting in the first context. While an unexpected activity of the “Canteen”
 655 load after 19:00, tagged with severity 8, contributed to the occurrence of an anomalous subsequence of severity 8 at meter-
 656 level that started at 19:15 in context 5.

657 It is very likely that the anomalous behaviour of the “Mechanical room” during night hours, the early start-up and late
 658 switch-off was caused by both incorrect operation and wrong schedule of the chillers. This kind of anomaly can be easily
 659 detected, diagnosed and rectified by facility managers by reviewing the system ON-OFF logics. On the other hand, the
 660 correction of energy management procedures pertaining the “Not labelled” load can be more challenging since there is
 661 no detail about the energy systems served by this sub-load of the substation C.



662

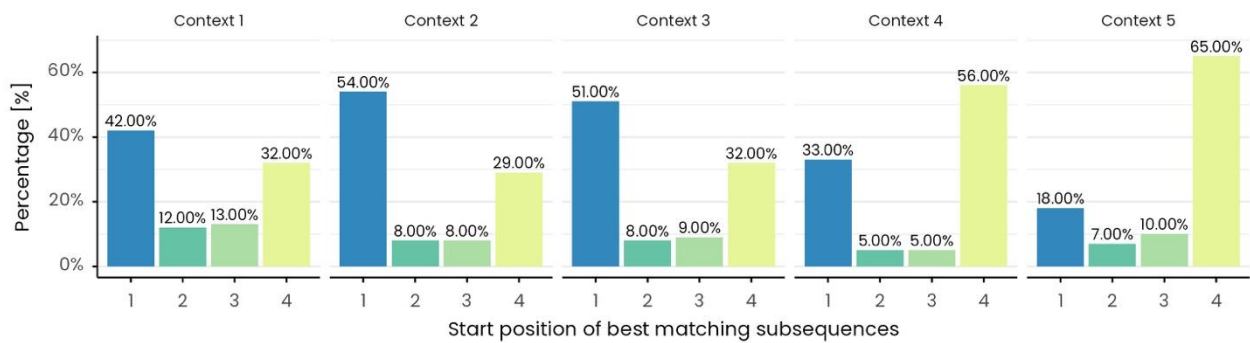
663 **Figure 15.** Comparison between anomalous meter-level daily electrical load profile and responsible sub-load identified through the diagnosis process.
 664 The figure shows the load profiles of (a) Saturday 6th of July 2019 belonging to cluster 2 and the load profiles of (b) Tuesday 30th of July 2019 belonging
 665 to cluster 4. The vertical bands denote the anomalous subsequences at meter and sub-load level, the band colours are referred to the contexts considered,
 666 while the band widths are related to the subsequence lengths.

667 4. Discussion and conclusion

668 In this paper a meter-level ADD process is proposed in order to demonstrate (i) the capabilities of the CMP algorithm
 669 as unsupervised method for detecting anomalies in energy consumption timeseries and (ii) the potential of integrating
 670 such tool in a building energy management system to gain insights on infrequent energy consumption patterns of sub-
 671 loads, promoting the correction of anomalies and reduction of energy wastes. The methodology is based on the coupling
 672 of unsupervised, supervised and timeseries analytics methods with domain knowledge to detect infrequent/anomalous
 673 patterns in energy consumption timeseries of building collected at meter-level.

674 The capability to identify anomalies that start in specific sub-daily periods (i.e., context) represents an opportunity for
 675 the early-stage identification and prompt correction of incorrect operations that can help prevent energy wastes over time.
 676 The recognized 55 anomalous subsequences start at different times in the identified contexts demonstrating the
 677 importance that the relaxation of temporal constraints (i.e., time alignment in the search of frequent/infrequent patterns)
 678 has in the detection of frequent and infrequent patterns in energy consumption timeseries of buildings.

679 As a matter of fact, in Figure 16 is reported the percentage of best matching subsequences that begin in each of the 4
 680 possible start-positions of the considered 5 contexts with a duration of 1h, where for instance position 3 is located three
 681 15-min timesteps apart from the beginning of the context. In particular, it is possible to observe that the context length
 682 was fully exploited in the *similarity-join-search*, allowing to identify pairs of similar subsequences that have
 683 approximately the same overall shape and magnitude even though they are slightly shifted between each other in the time
 684 axis.



685

686 **Figure 16.** Relative frequency associated to the start positions of best matching subsequences identified through CMP for each context. The start
 687 positions corresponding to the four 15-min timesteps in a context are numbered from 1 to 4.

688 Correctly defining the number and the length of contexts and subsequences for tuning the *similarity-join-search* is a
 689 complex task and their wrong setting may negatively affect the capability to isolate important features of the daily load
 690 profile such as the periods related to the start-up or shut down of building energy systems. The regression tree model
 691 (developed using the CART algorithm) has proven to be flexible enough to adapt to different load conditions (also with
 692 reference to different building end-uses as reported in [15]) thanks to the easily generalizable criteria in the tuning of
 693 model parameters. The introduction of the concept of context, i.e., the time interval during which a subsequence of interest
 694 may start, allows to compare subsequences that are not perfectly aligned in the time domain and at the same time to avoid
 695 the computation of similarity-joins that are meaningless considering the operation schedule of the analysed building.

696 The methodology was tested off-line on a static dataset but was conceptualized to work in semi-real time, by enabling
 697 the ADD process at the end of each reference time interval of subsequences that were extracted from the regression tree
 698 (reported in Table 1). However, the practical implementation requires to tackle different technical challenges.

699 Firstly, the ADD process must follow the time constraint of the data stream, meaning that the execution time must be
 700 lower than the interval between two subsequent triggers of the procedure itself. The proposed methodology was not
 701 intended to be a pure real-time process performed upon every new incoming datapoint, rather it was conceived as a semi
 702 real-time batch process that can be performed at the end of each reference time interval; where for batch process is
 703 intended a process in which the whole timeseries must be available before calculating the CMP [25]. In this perspective,
 704 the process execution is triggered only when the last subsequence, that may start in a context, has been fully collected
 705 (i.e., at the end of the reference time intervals). In the analysed case study, the minimum interval of time between
 706 subsequent triggers of the ADD process is 2.5 hours (i.e., the length of the smallest subsequence), meaning that the

707 execution time of the whole ADD procedure including CMP calculation on all the timeseries (meter-level and sub-meter
708 level) must be lower than 2.5 hours. The actual execution time on the offline one-year dataset, tested on 32 GB RAM
709 Intel(R) Core(TM) i7-4790, is of 1 min and 45 seconds for each load which is far less than 2.5 hours proving that a semi-
710 real-time batch approach is suitable for this kind of analysis and further scalable to a higher number of sub-loads.
711 However, a streaming approach is more desirable and would enable an even earlier recognition of anomalies. The
712 computational burden could be reduced performing the CMP calculation by using an incremental approach like the
713 STAMPI algorithm [25] that adjusts the CMP rather than re-compute it once a new observation is collected.

714 The second challenge, that exists in the deployment phase of the ADD process is related the so-called cold-start
715 problem, intended as scarcity of initial data useful to perform an accurate definition of subsequences, contexts, groups for
716 the CMP calculation. One advantage of the MP algorithm is that, differently from CMP, it does not require a minimum
717 length of the timeseries to be calculated. On the other hand, since the CMP objective is to perform a *similarity-join-search*
718 to discover frequent and infrequent patterns within specific segments of timeseries that may be characterized by a certain
719 periodicity, with a small dataset CMP may fail to recognize and highlight relevant patterns. A dataset containing at least
720 4 weeks of observations is desirable in order to have a minimum set of daily load profiles on which it is possible to
721 perform a more robust discord/motif discovery. With respect to the hyperparameters of CMP (subsequences, contexts and
722 groups), in the first deployment period, they can be easily defined based on domain knowledge, setting a reasonable
723 context length and by grouping daily load profiles according to typical building operational schedules. As a first
724 configuration it is possible to define the subsequences by simply split the 24h into N not overlapping time windows of
725 fixed length. As a reference, this approach was followed by [13] in the definition of SAX parameters for the
726 characterization of daily load profiles, where after a sensitivity analysis the suggested setting for N was between 3 and 4.
727 Instead, the initial number of groups (i.e., clusters) can be based on the weekly operational calendar. For example, a
728 possible hyperparameter setting may consist in 3 groups (weekday, Saturday and Sundays/holidays), 4 fixed-length
729 subsequences of 6 hours with the corresponding context of 1h length.

730 To better discuss this scenario, the authors performed the same experiment on the entire one-year dataset, using this
731 hyperparameter setting. The results obtained, although not comparable with the ADD results presented in the paper since
732 no ground truth is available, shows that 40 anomalous sub-sequences were detected at meter-level (almost in the same
733 periods of the presented results). However, it was possible to noticed that, despite the a-priori grouping and a-priori
734 subsequence definition is a good compromise, it was not effective as the proposed process, since the context-group CMPs
735 exhibited very high internal variability of distance values, meaning that the subsequences of daily load profiles within the
736 same group were not so homogeneous. In this way, a wider distribution of the distances made the identification of isolated
737 and extremal observations, by means of statistical methods, more difficult to be performed.

738 Along with the data scarcity, for the CMP algorithm even the abundance of data may represent a critical issue, not
739 only under the computational point of view but mainly from the conceptual one. Frequent and rare subsequences in the
740 original concept of CMP are defined as the ones with smallest/largest 1st nearest neighbour distance [25]. This implies
741 that if a rare subsequence occurs more than once in the timeseries it may be considered as common or even frequent
742 pattern [73] and tackle this issue is of paramount importance in building energy management since an anomaly if not
743 promptly detected may persists in time and must not be considered as a motif. A repeated anomaly would cause false
744 negatives due to the previous anomaly instance being part of all subsequence set: this issue is recognized in the literature
745 as the twin freak problem [73].

746 To address the twin freak problem, in the future, the implementation of the kNN distance instead of 1st Nearest
747 Neighbour distance in the CMP calculation presented in this paper can be included as a process improvement. One of the

748 possible algorithmic implementation was discussed in [73] where the authors proposed a density based approach for the
 749 kNN calculation applied in the MP algorithm. A further step would be to enhance the robustness of the ADD
 750 methodology, through the dynamic adjustment of parameters and weights, based on human users' feedbacks to report
 751 anomaly using a human-in-the-loop training scheme.

752 Given the variety of building load conditions and their variability throughout the day due to multiple forcing factors
 753 such as the building occupancy, user behaviour, system failures and external boundary conditions it is necessary to assist
 754 energy managers during building operation by means of easily interpretable and effective informative tools. To this aim,
 755 the detection and diagnosis output of the proposed ADD process has been represented with an easy interpretable anomaly
 756 severity score, coupled with an effective visualization that allows to immediately compare the detected anomalous
 757 behaviour with the reference one for the same context and group (i.e., cluster centroid profile). Thanks to the proposed
 758 hierarchical methodological framework, the user is not overwhelmed by the stream of information and data but instead is
 759 supported in the ADD process and involved only in the case the sub-loads determine the occurrence of an anomalous
 760 energy trend at meter-level de facto rationalizing the quality and number of feedbacks he/she receive during the day-by-
 761 day energy management operations.

762 References

- 763 [1] Directorate-General for Energy (European Commission), Clean energy for all Europeans - Publications Office of the EU, Publ. Off. EU. 14
 764 (2019) 3. <https://doi.org/10.2833/9937>.
- 765 [2] T. Ramesh, R. Prakash, K.K. Shukla, Life cycle energy analysis of buildings: An overview, *Energy Build.* 42 (2010) 1592–1600.
 766 <https://doi.org/10.1016/j.enbuild.2010.05.007>.
- 767 [3] L. Erhan, M. Ndubuaku, M. Di Mauro, W. Song, M. Chen, G. Fortino, O. Bagdasar, A. Liotta, Smart anomaly detection in sensor systems:
 768 A multi-perspective review, *Inf. Fusion.* 67 (2021) 64–79. <https://doi.org/10.1016/j.inffus.2020.10.001>.
- 769 [4] S. Rinaldi, P. Bellagente, A.L. Camillo Ciribini, L. Chiara Tagliabue, T. Poli, A. Giovanni Mainini, A. Speroni, J.D. Blanco Cadena, S.
 770 Lupica Spagnolo, A cognitive-driven building renovation for improving energy efficiency: The experience of the elisir project, *Electron.* 9
 771 (2020). <https://doi.org/10.3390/electronics9040666>.
- 772 [5] M. Molina-Solana, M. Ros, M.D. Ruiz, J. Gómez-Romero, M.J. Martín-Bautista, Data science for building energy management: A review,
 773 *Renew. Sustain. Energy Rev.* 70 (2017) 598–609. <https://doi.org/10.1016/j.rser.2016.11.132>.
- 774 [6] H. Kramer, G. Lin, J. Granderson, C. Curtin, E. Crowe, Synthesis of Year One Outcomes in the Smart Energy Analytics Campaign
 775 Building Technology and Urban Systems Division, Lawrence Berkeley National Laboratory: Berkeley, CA, USA. (2019).
 776 <https://doi.org/10.2172/1545159>
- 777 [7] C. Fan, Y. Sun, Y. Zhao, M. Song, J. Wang, Deep learning-based feature engineering methods for improved building energy prediction,
 778 *Appl. Energy.* 240 (2019) 35–45. <https://doi.org/10.1016/j.apenergy.2019.02.052>.
- 779 [8] C. Fan, F. Xiao, S. Wang, Development of prediction models for next-day building energy consumption and peak power demand using data
 780 mining techniques, *Appl. Energy.* 127 (2014) 1–10. <https://doi.org/10.1016/j.apenergy.2014.04.016>.
- 781 [9] N. Somu, G.R. M R, K. Ramamritham, A hybrid model for building energy consumption forecasting using long short term memory
 782 networks, *Appl. Energy.* 261 (2020) 114131. <https://doi.org/10.1016/j.apenergy.2019.114131>.
- 783 [10] T. Liu, Z. Tan, C. Xu, H. Chen, Z. Li, Study on deep reinforcement learning techniques for building energy consumption forecasting,
 784 *Energy Build.* 208 (2020). <https://doi.org/10.1016/j.enbuild.2019.109675>.
- 785 [11] M.S. Piscitelli, S. Brandi, A. Capozzoli, F. Xiao, A data analytics-based tool for the detection and diagnosis of anomalous daily energy
 786 patterns in buildings, *Build. Simul.* (2020) 1–17. <https://doi.org/10.1007/s12273-020-0650-1>.
- 787 [12] R. Chiosa, A. Capozzoli, M.S. Piscitelli, A Data Analytics-Based Energy Information System (EIS) Tool to Perform Meter-Level
 788 Anomaly Detection and Diagnosis in Buildings, *Energies.* 14 (2021) 1–28. <https://doi.org/10.3390/en14010237>.
- 789 [13] C. Miller, Z. Nagy, A. Schlueter, Automated daily pattern filtering of measured building performance data, *Autom. Constr.* 49 (2015) 1–17.
 790 <https://doi.org/10.1016/j.autcon.2014.09.004>.
- 791 [14] C. Fan, F. Xiao, H. Madsen, D. Wang, Temporal knowledge discovery in big BAS data for building energy management, *Energy Build.* 109
 792 (2015) 75–89. <https://doi.org/10.1016/j.enbuild.2015.09.060>.
- 793 [15] A. Capozzoli, M.S. Piscitelli, S. Brandi, D. Grassi, G. Chicco, Automated load pattern learning and anomaly detection for enhancing energy
 794 management in smart buildings, *Energy.* 157 (2018) 336–352. <https://doi.org/10.1016/j.energy.2018.05.127>.
- 795 [16] P. Arjunan, K. Poolla, C. Miller, BEEM: Data-driven building energy benchmarking for Singapore, *Energy Build.* 260 (2022) 111869.
 796 <https://doi.org/10.1016/j.enbuild.2022.111869>.
- 797 [17] K. Li, Y. Sun, D. Robinson, J. Ma, Z. Ma, A new strategy to benchmark and evaluate building electricity usage using multiple data mining

- 798 technologies, *Sustain. Energy Technol. Assessments*. 40 (2020) 100770. <https://doi.org/10.1016/j.seta.2020.100770>.
- 799 [18] A. Capozzoli, M.S. Piscitelli, F. Neri, D. Grassi, G. Serale, A novel methodology for energy performance benchmarking of buildings by
800 means of Linear Mixed Effect Model: The case of space and DHW heating of out-patient Healthcare Centres, *Appl. Energy*. 171 (2016)
801 592–607. <https://doi.org/10.1016/j.apenergy.2016.03.083>.
- 802 [19] X. Luo, T. Hong, Y. Chen, M.A. Piette, Electric load shape benchmarking for small- and medium-sized commercial buildings, *Appl.*
803 *Energy*. 204 (2017) 715–725. <https://doi.org/10.1016/j.apenergy.2017.07.108>.
- 804 [20] M.S. Piscitelli, S. Brandi, A. Capozzoli, Recognition and classification of typical load profiles in buildings with non-intrusive learning
805 approach, *Appl. Energy*. 255 (2019) 113727. <https://doi.org/10.1016/j.apenergy.2019.113727>.
- 806 [21] J.Y. Park, E. Wilson, A. Parker, Z. Nagy, The good, the bad, and the ugly: Data-driven load profile discord identification in a large building
807 portfolio, *Energy Build.* 215 (2020) 109892. <https://doi.org/10.1016/j.enbuild.2020.109892>.
- 808 [22] J.Y. Park, X. Yang, C. Miller, P. Arjunan, Z. Nagy, Apples or oranges? Identification of fundamental load shape profiles for benchmarking
809 buildings using a large and diverse dataset, *Appl. Energy*. 236 (2019) 1280–1295. <https://doi.org/10.1016/j.apenergy.2018.12.025>.
- 810 [23] K. Nweye, Z. Nagy, MARTINI: Smart Meter Driven Estimation of HVAC Schedules and Energy Savings Based on WiFi Sensing and
811 Clustering, (2021) 0–13. <http://arxiv.org/abs/2110.08927>.
- 812 [24] A. Capozzoli, M.S. Piscitelli, A. Gorrino, I. Ballarini, V. Corrado, Data analytics for occupancy pattern learning to reduce the energy
813 consumption of HVAC systems in office buildings, *Sustain. Cities Soc.* 35 (2017) 191–208. <https://doi.org/10.1016/j.scs.2017.07.016>.
- 814 [25] C.-C.M. Yeh, Y. Zhu, L. Ulanova, N. Begum, Y. Ding, H.A. Dau, D.F. Silva, A. Mueen, E. Keogh, Matrix Profile I: All Pairs Similarity
815 Joins for Time Series: A Unifying View That Includes Motifs, Discords and Shapelets, (2017) 1317–1322.
816 <https://doi.org/10.1109/icdm.2016.0179>.
- 817 [26] D. De Paepe, S. Vanden Haute, B. Steenwinckel, F. De Turck, F. Ongenaes, O. Janssens, S. Van Hoecke, A generalized matrix profile
818 framework with support for contextual series analysis, *Eng. Appl. Artif. Intell.* 90 (2020). <https://doi.org/10.1016/j.engappai.2020.103487>.
- 819 [27] P. Tan, M. Steinbach, V. Kumar, *Introduction to data mining*, Addison Wesley, 2011. https://doi.org/10.1007/978-3-642-19721-5_1.
- 820 [28] Y. Himeur, K. Ghanem, A. Alsalemi, F. Bensaali, A. Amira, Artificial intelligence based anomaly detection of energy consumption in
821 buildings: A review, current trends and new perspectives, *Appl. Energy*. 287 (2021) 116601.
822 <https://doi.org/10.1016/j.apenergy.2021.116601>.
- 823 [29] N.R. Prasad, S. Almanza-Garcia, T.T. Lu, Anomaly detection: A Survey, *Comput. Mater. Contin.* 14 (2009) 1–22.
824 <https://doi.org/10.1145/1541880.1541882>.
- 825 [30] J.E. Zhang, D. Wu, B. Boulet, Time Series Anomaly Detection for Smart Grids: A Survey, 2021 IEEE Electr. Power Energy Conf. EPEC
826 2021. (2021) 125–130. <https://doi.org/10.1109/EPEC52095.2021.9621752>.
- 827 [31] H. Chen, X. Fei, S. Wang, X. Lu, G. Jin, W. Li, X. Wu, Energy Consumption Data Based Machine Anomaly Detection, *Proc. - 2014 2nd*
828 *Int. Conf. Adv. Cloud Big Data, CBD 2014*. (2015) 136–142. <https://doi.org/10.1109/CBD.2014.24>.
- 829 [32] Y. Zhao, S. Wang, F. Xiao, Pattern recognition-based chillers fault detection method using Support Vector Data Description (SVDD), *Appl.*
830 *Energy*. 112 (2013) 1041–1048. <https://doi.org/10.1016/j.apenergy.2012.12.043>.
- 831 [33] K. Kammerer, B. Hoppenstedt, R. Pryss, S. Stöckler, J. Allgaier, M. Reichert, Anomaly detections for manufacturing systems based on
832 sensor data—insights into two challenging real-world production settings, *Sensors*. 19 (2019). <https://doi.org/10.3390/s19245370>.
- 833 [34] J.S. Chou, A.S. Telaga, Real-time detection of anomalous power consumption, *Renew. Sustain. Energy Rev.* 33 (2014) 400–411.
834 <https://doi.org/10.1016/j.rser.2014.01.088>.
- 835 [35] R. Wu, E. Keogh, Current Time Series Anomaly Detection Benchmarks are Flawed and are Creating the Illusion of Progress, *IEEE Trans.*
836 *Knowl. Data Eng.* (2021) 1–9. <https://doi.org/10.1109/TKDE.2021.3112126>.
- 837 [36] C. Fan, F. Xiao, Y. Zhao, J. Wang, Analytical investigation of autoencoder-based methods for unsupervised anomaly detection in building
838 energy data, *Appl. Energy*. 211 (2018) 1123–1135. <https://doi.org/10.1016/j.apenergy.2017.12.005>.
- 839 [37] D.J. Bartholomew, Principal components analysis, *Int. Encycl. Educ.* (2010) 374–377. [https://doi.org/10.1016/B978-0-08-044894-7.01358-](https://doi.org/10.1016/B978-0-08-044894-7.01358-0)
840 [0](https://doi.org/10.1016/B978-0-08-044894-7.01358-0).
- 841 [38] P. Esling, C. Agon, Time-series data mining, *ACM Comput. Surv.* 45 (2012). <https://doi.org/10.1145/2379776.2379788>.
- 842 [39] A. Capozzoli, M.S. Piscitelli, S. Brandi, Mining typical load profiles in buildings to support energy management in the smart city context,
843 *Energy Procedia*. 134 (2017) 865–874. <https://doi.org/10.1016/j.egypro.2017.09.545>.
- 844 [40] Y. Zhao, C. Zhang, Y. Zhang, Z. Wang, J. Li, A review of data mining technologies in building energy systems: Load prediction, pattern
845 identification, fault detection and diagnosis, *Energy Built Environ.* 1 (2020) 149–164. <https://doi.org/10.1016/j.enbenv.2019.11.003>.
- 846 [41] V. Chandola, D. Cheboli, V. Kumar, Detecting anomalies in a time series database, (2009). <https://hdl.handle.net/11299/215791>
- 847 [42] C. Zhang, F. Wang, Multi-feature fusion based anomaly electro-data detection in smart grid, *Proc. - 2018 15th Int. Symp. Pervasive Syst.*
848 *Algorithms Networks, I-SPAN 2018*. (2019) 54–59. <https://doi.org/10.1109/I-SPAN.2018.00018>.
- 849 [43] K. Li, Z. Ma, D. Robinson, J. Ma, Using Evidence Accumulation-Based Clustering and Symbolic Transformation to Group Multiple
850 Buildings Based on Electricity Usage Patterns, *Sustain. Energy Build.*, (2019) 61–67. <https://doi.org/10.1007/978-981-32-9868-2>.
- 851 [44] B. Rossi, S. Chren, B. Buhnova, T. Pitner, Anomaly Detection in Smart Grid Data: An Experience Report, (2016) 2313–2318.
- 852 [45] J. Lin, E. Keogh, L. Wei, S. Lonardi, Experiencing SAX: a Novel Symbolic Representation of Time Series, *Cs.Gmu.Edu.* 15 (2007) 107–
853 144.
- 854 [46] E. Keogh, K. Chakrabarti, M. Pazzani, S. Mehrotra, Dimensionality Reduction for Fast Similarity Search in Large Time Series Databases,
855 *Knowl. Inf. Syst.* 3 (2001) 263–286. <https://doi.org/10.1007/pl00011669>.
- 856 [47] H. Ren, M. Liu, Z. Li, W. Pedrycz, A Piecewise Aggregate pattern representation approach for anomaly detection in time series,

- 857 Knowledge-Based Syst. 135 (2017) 29–39. <https://doi.org/10.1016/j.knosys.2017.07.021>.
- 858 [48] M. Alshaer, S. Garcia-Rodriguez, C. Gouy-Pailler, Detecting Anomalies from Streaming Time Series using Matrix Profile and Shapelets
859 Learning, Proc. - Int. Conf. Tools with Artif. Intell. ICTAI. 2020-Novem (2020) 376–383.
860 <https://doi.org/10.1109/ICTAI50040.2020.00066>.
- 861 [49] L. Ye, E. Keogh, Time series shapelets: A new primitive for data mining, Proc. ACM SIGKDD Int. Conf. Knowl. Discov. Data Min. (2009)
862 947–955. <https://doi.org/10.1145/1557019.1557122>.
- 863 [50] S.D.D. Anton, H.D. Schotten, Intrusion Detection in Binary Process Data: Introducing the Hamming-distance to Matrix Profiles, Proc. -
864 21st IEEE Int. Symp. a World Wireless, Mob. Multimed. Networks, WoWMoM 2020. (2020) 347–353.
865 <https://doi.org/10.1109/WoWMoM49955.2020.00065>.
- 866 [51] F. Madrid, S. Imani, R. Mercer, Z. Zimmerman, N. Shakibay, E. Keogh, Matrix profile XX: Finding and visualizing time series motifs of
867 all lengths using the matrix profile, Proc. - 10th IEEE Int. Conf. Big Knowledge, ICBK 2019. (2019) 175–182.
868 <https://doi.org/10.1109/ICBK.2019.00031>.
- 869 [52] J. Dinal Herath, C. Bai, G. Yan, P. Yang, S. Lu, RAMP: Real-Time Anomaly Detection in Scientific Workflows, Proc. - 2019 IEEE Int.
870 Conf. Big Data, Big Data 2019. (2019) 1367–1374. <https://doi.org/10.1109/BigData47090.2019.9005653>.
- 871 [53] D. De Paepe, D.N. Avendano, S. Van Hoecke, Implications of Z-Normalization in the Matrix Profile, Lect. Notes Comput. Sci. (2020) 95–
872 118. https://doi.org/10.1007/978-3-030-40014-9_5.
- 873 [54] D. De Paepe, O. Janssens, S. Van Hoecke, Eliminating noise in the matrix profile, ICPRAM 2019 - Proc. 8th Int. Conf. Pattern Recognit.
874 Appl. Methods. (2019) 83–93. <https://doi.org/10.5220/0007314100830093>.
- 875 [55] S. Ahmad, A. Lavin, S. Purdy, Z. Agha, Unsupervised real-time anomaly detection for streaming data, Neurocomputing. 262 (2017) 134–
876 147. <https://doi.org/10.1016/j.neucom.2017.04.070>.
- 877 [56] C. Nichiforov, G. Stamatescu, I. Stamatescu, I. Fagarasan, Learning Dominant Usage from Anomaly Patterns in Building Energy Traces,
878 IEEE Int. Conf. Autom. Sci. Eng. (2020) 548–553. <https://doi.org/10.1109/CASE48305.2020.9216794>.
- 879 [57] Y. Zhu, S. Gharghabi, D.F. Silva, H.A. Dau, C.C.M. Yeh, N. Shakibay Senobari, A. Almaslukh, K. Kamgar, Z. Zimmerman, G. Funning,
880 A. Mueen, E. Keogh, The Swiss army knife of time series data mining: ten useful things you can do with the matrix profile and ten lines of
881 code, Springer US, 2020. <https://doi.org/10.1007/s10618-019-00668-6>.
- 882 [58] S. Alaei, K. Kamgar, E. Keogh, Matrix profile XXII: Exact discovery of time series motifs under DTW, Proc. - IEEE Int. Conf. Data
883 Mining, ICDM. 2020-Novem (2020) 900–905. <https://doi.org/10.1109/ICDM50108.2020.00099>.
- 884 [59] S. Gharghabi, S. Imani, A. Bagnall, A. Darvishzadeh, E. Keogh, An ultra-fast time series distance measure to allow data mining in more
885 complex real-world deployments, Springer US, 2020. <https://doi.org/10.1007/s10618-020-00695-8>.
- 886 [60] A. Mueen, Y. Zhu, M. Yeh, K. Kamgar, K. Viswanathan, C. Gupta, E. Keogh, The Fastest Similarity Search Algorithm for Time Series
887 Subsequences under Euclidean Distance, (2017). <http://www.cs.unm.edu/~mueen/FastestSimilaritySearch.html>
- 888 [61] C. Onwongsa, C. Ratanamahatana, An enhanced time series motif discovery using approximated matrix profile, ACM Int. Conf. Proceeding
889 Ser. (2020) 180–189. <https://doi.org/10.1145/3421558.3421586>.
- 890 [62] C.C.M. Yeh, N. Kavantzias, E. Keogh, Matrix profile VI: Meaningful multidimensional motif discovery, Proc. - IEEE Int. Conf. Data
891 Mining, ICDM. 2017-Novem (2017) 565–574. <https://doi.org/10.1109/ICDM.2017.66>.
- 892 [63] H.A. Dau, E. Keogh, Matrix profile V: A generic technique to incorporate domain knowledge into motif discovery, Proc. ACM SIGKDD
893 Int. Conf. Knowl. Discov. Data Min. Part F1296 (2017) 125–134. <https://doi.org/10.1145/3097983.3097993>.
- 894 [64] G.E.A.P.A. Batista, E.J. Keogh, O.M. Tataw, V.M.A. De Souza, CID: An efficient complexity-invariant distance for time series, Data Min.
895 Knowl. Discov. 28 (2014) 634–669. <https://doi.org/10.1007/s10618-013-0312-3>.
- 896 [65] J.L. Mathieu, P.N. Price, S. Kiliccote, M.A. Piette, Quantifying changes in building electricity use, with application to demand response,
897 IEEE Trans. Smart Grid. 2 (2011) 507–518. <https://doi.org/10.1109/TSG.2011.2145010>.
- 898 [66] J. Zhu, Y. Shen, Z. Song, D. Zhou, Z. Zhang, A. Kusiak, Data-driven building load profiling and energy management, Sustain. Cities Soc.
899 49 (2019) 101587. <https://doi.org/10.1016/j.scs.2019.101587>.
- 900 [67] R. Yan, Z. Ma, Y. Zhao, G. Kokogiannakis, A decision tree based data-driven diagnostic strategy for air handling units, Energy Build. 133
901 (2016) 37–45. <https://doi.org/10.1016/j.enbuild.2016.09.039>.
- 902 [68] T. Grubinger, A. Zeileis, K.P. Pfeiffer, Evtree: Evolutionary learning of globally optimal classification and regression trees in R, J. Stat.
903 Softw. 61 (2014) 1–29. <https://doi.org/10.18637/jss.v061.i01>.
- 904 [69] V. Satopää, J. Albrecht, D. Irwin, B. Raghavan, Finding a “kneedle” in a haystack: Detecting knee points in system behavior, Proc. - Int.
905 Conf. Distrib. Comput. Syst. (2011) 166–171. <https://doi.org/10.1109/ICDCSW.2011.20>.
- 906 [70] R Core Team, R: A Language and Environment for Statistical Computing, (2017). <https://www.r-project.org/>.
- 907 [71] G. Van Rossum, F.L. Drake Jr, Python reference manual, Centrum voor Wiskunde en Informatica Amsterdam, (1995).
- 908 [72] M. Charrad, N. Ghazzali, V. Boiteau, A. Niknafs, NbClust : An R Package for Determining the, J. Stat. Softw. 61 (2014) 1–36.
909 <https://doi.org/10.18637/jss.v061.i06>.
- 910 [73] Y. He, X. Chu, Y. Wang, Neighbor profile: Bagging nearest neighbors for unsupervised time series mining, Proc. - Int. Conf. Data Eng.
911 (2020) 373–384. <https://doi.org/10.1109/ICDE48307.2020.00039>.
- 912