

# Emergent Deep Learning for Anomaly Detection in Internet of Everything

Youcef Djenouri, Djamel Djenouri, Asma Belhadi, Gautam Srivastava\*\*, and Jerry Chun-Wei Lin\*

**Abstract**—This research presents a new generic deep learning framework for anomaly detection in the Internet of Everything (IoE). It combines decomposition methods, deep neural networks, and evolutionary computation to better detect outliers in IoE environments. The dataset is first decomposed into clusters, while similar observations in the same cluster are grouped. Five clustering algorithms were used for this purpose. The generated clusters are then trained using Deep Learning architectures. In this context, we propose a new recurrent neural network for training time series data. Two evolutionary computational algorithms are also proposed: the genetic and the bee swarm to fine-tune the training step. These algorithms consider the hyper-parameters of the trained models and try to find the optimal values. The proposed solutions have been experimentally evaluated for two use cases: 1) road traffic outlier detection and 2) network intrusion detection. The results show the advantages of the proposed solutions and a clear superiority compared to state-of-the-art approaches.

**Index Terms**—Internet of Everything, Intrusion Detection, Smart Transportation, Deep Learning.

## I. INTRODUCTION

In this research work, we focus on the new offshoot of the Internet of Things (IoT), the Internet of Everything (IoE). The IoE extends the IoT by placing a greater emphasis on machine-to-machine (M2M) communication to describe more complex systems that can include people and processes, while considering intelligent connectivity and data processing. This concept enables the accumulation of an enormous amount of data. Effective processing and analysis of such Big Data, while challenging, will drive innovative applications in various fields such as cloud services [1], smart healthcare [2], smart buildings [3], robotics [4], and others. Anomaly detection refers to the process of filtering out anomalies from collected data. The term anomaly is general and can be used to refer to many problems, depending on the application erroneous data that may occur due to faulty sensors or during the data fusion process [5], road traffic outliers, or computer network intrusions [6], [7]. This research work is in this direction

Y. Djenouri is with the Dept. of Mathematics and Cybernetics, SINTEF Digital, Oslo, Norway, youcef.djenouri@sintef.no

D. Djenouri is with the CSRC, Dept. of Computer Science and Creative Technologies, University of the West of England, Bristol, UK, djamel.djenouri@uwe.ac.uk

A. Belhadi is with the Dept. of Technology, Kristiania University College, Oslo, Norway, asma.belhadi@kristiania.no

G. Srivastava is with the Dept. of Mathematics & Computer Science, Brandon University, Canada, and Research Centre for Interneural Computing, China Medical University, Taichung, Taiwan, srivastavag@brandonu.ca (\*\*Co-corresponding author)

J. C. W. Lin is with the Dept. of Computer Science, Electrical Engineering and Mathematical Sciences, Western Norway University of Applied Sciences, Bergen, Norway, jerrylin@ieec.org (\*Corresponding author)

and proposes a new intelligent framework to efficiently and accurately identify anomalies in IoE environments.

Most current anomaly detection solutions in IoE [6]–[8] are time consuming and have low accuracy. Deep learning based solutions [6], [7] provide relatively better accuracy compared to traditional solutions [8], but the improvement is still limited. The main reason is that they need to build a complex model with a high number of parameters to be specified. For example, the recurrent neural network (RNN) [9] requires a large number of states, and each state has parameters that need to be set. Evolutionary computation [10] is also widely studied for anomaly detection, but these solutions are limited only by exploration of the observation space and evaluate each observation separately. Motivated by the success of decomposition, deep learning (DL) and evolutionary computation in solving many real-world applications [11], [12], this research proposes a hybrid framework for inferring anomalies from IoE.

In this paper, we propose deep learning-based decomposition and evolutionary computation framework for anomaly detection networks (D2E-ADN) that aims to build targeted learning models for inferring anomalies in IoE. The data collected from the IoE environment is first divided into several small but as independent clusters as possible, minimizing the number of shared data between the clusters. The generated clusters are used to train the DL models, with each cluster used to train its own model. A hyperparameter optimizer is also investigated to accurately find the relevant parameters of the DL models. In this sense, the main contributions of this work are as follows:

- 1) We propose five decomposition algorithms for clustering data while extracting the relevant features from the IoE. The data clusters are then identified using clustering algorithms whose goal is to minimize the number of the shared data between clusters.
- 2) We propose a new DL model that uses the knowledge gained in the decomposition step. It is based on the recurrent neural network developed for processing time series data.
- 3) We propose two evolutionary computational algorithms to tune the parameters of the different steps of the D2E-ADN system, including the number of clusters in the decomposition step, the number of epochs, the learning error rate, and the activation functions for the DL models. The first evolutionary computational algorithm explores genetic optimization, while the second considers the behavior of the bees in exploring the possible configuration of the hyperparameters of the

D2E-ADN system.

- 4) We evaluate D2E-ADN by comparing its computation time and accuracy with basic anomaly detection algorithms in two areas: intelligent transport (detecting outliers in traffic flow) and network security (detecting intrusions). This evaluation shows that D2E-ADN outperforms the baseline algorithms in both runtime and accuracy.

We give an outline of the remainder of this paper here. Section II gives an in-depth literature survey of existing solutions in anomaly detection. Next, in Section III, we present our proposed approach detailing all of its main components. Section IV gives our experimental analysis, discussion, and results. Lastly, Section V terminates our paper with some closing ideas.

## II. RELATED WORK

Zhong *et al.* [6] proposed a hybrid DL model for intrusion detection in a large network. The set of relevant features is first extracted using the damped incremental statistics algorithm. Then, the autoencoder algorithm is implemented to generate the training data, which is finally used to train the recurrent neural network model. Pawar *et al.* [13] proposed a DL framework for intrusion detection in the context of video-based activity recognition. An intensive comparative study of existing traditional machine learning techniques and advanced DL intrusion detection algorithms was conducted. Roberto *et al.* [7] developed a model for a convolutional neural network to identify abnormal traffic flows. The authors also provided a strategy for generating the labeled data used in the learning process. Khan *et al.* [14] proposed a novel two-stage DL algorithm for network intrusion detection. Network traffic is first classified into two classes (normal vs. abnormal) based on a probability score, which is then used as an additional feature to identify normal behavior or attack classes. Jallad *et al.* [15] used long-term memory (LSTM) to identify different types of intrusion detection such as point anomalies, collective anomalies, and contextual anomalies. The solution was tested on a large network for several million packets using the Spark platform. The results confirm the usefulness of the methods over traditional methods such as kNN.

Abdurrahman *et al.* [16] proposed a hybrid model that derives botnet in network. It combines convolutional networks and recurrent neural networks in the overall process. The relevant features are extracted based on a graph structure strategy. The extracted features are then converted into feature vectors and considered as training data for the hybrid recurrent neural convolutional network model. Garg *et al.* [17] developed a model (hybrid) using the Boltzmann machine, which has been constrained as well as the SVM (Support Vector Machine) in identifying abnormal activities in social media (multimedia) networks. The approach uses an incremental strategy and includes a self-learning mechanism where the anomalies already detected are fed into the DL model. Pektas *et al.* [16] combined the convolutional neural network and the LSTM using spatiotemporal features of network flows. Specifically, the convolutional neural network learns the spatial features of the network, while the long-term memory

learns the temporal features. Ujjan *et al.* [18] presented an adaptive pooling-based sampling method to accurately infer distributed denial-of-service attacks in IoT. It integrates the snort intrusion detection system with the stacked autoencoder DL model to optimize detection accuracy in the control plane. Papamartzivanos *et al.* [19] developed a semi-supervised self-adaptive algorithm by integrating sparse autoencoder and feed-forward autoencoder to train the unlabeled data. Ferrag *et al.* [20] provided an overview of DL -based algorithms for detecting intrusions on 35 datasets. The DL models used in this study are based on neural networks (convolutional, recurrent), self-learning, and deep-belief networks. The detailed results show that the convolutional NN performs better than the models in both runtime and accuracy. Boukela *et al.* [21] developed the modified local outlier factor to mitigate the malfunction of security systems in IoT devices. This approach takes into account the handling of high-dimensional data, determining the reachability distance for all features of the selected neighbors. Edje *et al.* [22] developed a clustering-based algorithm for identifying fault and event outliers in IoT sensors. The event outliers are considered when there are problems in sensor readings. Noshouhi *et al.* [23] presented a new machine learning-based solution for predicting fires using spatiotemporal measurements. Relevant data such as temperature and humidity are trained, and the model attempts to separate abnormal cases from normal behaviors. A refinement process is also performed to ensure that the predicted anomalies are not due to outliers. Zhang *et al.* [24] seeks to ensure the confidentiality of Industrial Internet of Things customers by combining blockchain and federated learning. The fault detection system is developed to provide complete verification of customer data. Lin *et al.* [25] developed a multi-objective algorithm based on ant colony optimization metaheuristics for privacy preservation in IoT environment. The ant colony solution space is encoded and represented by hiding sensitive information. An external archive is used to preserve the extracted Pareto solutions. Chou *et al.* [26] proposed a taxonomy of intrusion detection datasets used for evaluation in the last two decades. In addition, future directions are proposed by extending intrusion detection to a cloud environment and creating ground truth based data in real network environments.

From this extensive literature review, it is clear that traffic anomaly detection solutions are often weak in terms of detection rate because the entire database must be considered during the learning process. Moreover, it is not clear how to tune the hyperparameters for DL models. In this work, we investigate a hybrid approach that combines PSO, decomposition, and CNN to efficiently find outliers and anomalies in traffic databases. We use both cluster-based algorithms and swarm-based approaches to tune CNN.

## III. DEEP LEARNING-BASED DECOMPOSITION AND EVOLUTIONARY COMPUTATION FOR ANOMALY DETECTION NETWORK

### A. Principle

Here, we present the proposed D2E-ADN framework that integrates decomposition, DL, and evolutionary computational

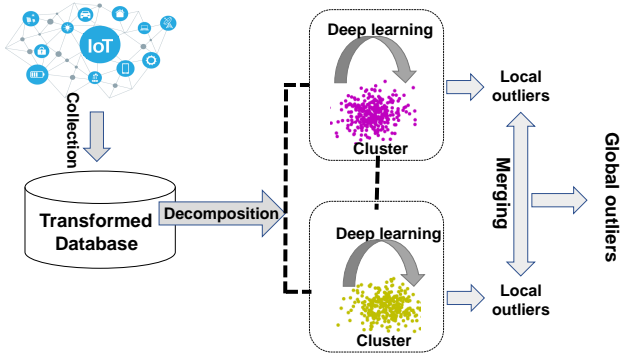


Fig. 1. Illustration of the D2E-ADN framework

201 optimization to identify anomalies in the data environment.  
 202 As shown in Fig. 1, the AD2E-ADN consists of three steps:  
 203 i) decomposition, which divides the data into clusters such  
 204 that each cluster contains similar data. ii) DL model, whose  
 205 goal is to apply the DL process to each cluster to identify  
 206 local anomalies. A merging strategy is used to combine  
 207 the local anomalies into global anomalies. iii) Evolutionary  
 208 computation, which is used to learn the hyperparameters of  
 209 the models of the clusters. In the following, each step is explained  
 210 in detail.

### 211 B. Decomposition

212 The main required aim to this step is for dividing the whole  
 213 data into  $k$  clusters,  $C = \{C_1, C_2, \dots, C_k\}$ , where each cluster  
 214  $C_s = \{D_1^{(s)}, D_2^{(s)}, \dots, D_{|C_s|}^{(s)}\}$  is the subset of the data  $D$ . The  
 215 overlapping data is minimized within clusters, and overlapping  
 216 data in each and every cluster is maximized. In other words,  
 217 using Eq. 1:

$$\begin{cases} \arg \min_C \left| \bigcup_{i=1, j=1}^k ((C_i) \cap (C_j)), i \neq j \right| \\ \wedge \\ \arg \max_C \left| \bigcup_{C_s} (D_i^{(s)} \cap D_j^{(s)}) \mid \forall (i, j) \in [1..|C_s|], i \neq j \right| \end{cases} \quad (1)$$

218 It is necessary to use different clustering algorithms than  
 219 in previous work [27]–[30] to minimize the number of shared  
 220 data between clusters and maximize the number of shared data  
 221 in each cluster. The following concepts should be introduced  
 222 here:

223 1) **Similarity computation.** The distance measure between  
 224 two data  $D_i$  and  $D_j$  is calculated by subtracting the  
 225 number of shared items from the number of all items  
 226 between  $D_i$  and  $D_j$ , as given in Eq. 2.

$$Dist(D_i, D_j) = \max(|D_i|, |D_j|) - (|D_i \cap D_j|) \quad (2)$$

227 2) **Centroids updating.** Here, we should consider datasets  
 228 of each cluster  $C_i = \{D_1^{(i)}, D_2^{(i)}, \dots, D_{|C_i|}^{(i)}\}$ , the aim is  
 229 to find a gravity center of this set which is also a datum.  
 230 The centroid,  $\mu_i$ , is computed based on the centroid  
 231 formula developed in [31]. Each item's frequency can

232 be calculated for all the data in a cluster  $C_i$ . Data center  
 233 length given as  $l_i$  is connected to the avg. number of  
 234 datum within  $C_i$ , as shown in Eq. 3.

$$l_i = \frac{\sum_{j=1}^{|C_i|} |D_j^{(i)}|}{|C_i|} \quad (3)$$

235 Afterwards, the data within  $C_i$  can be sorted by fre-  
 236 quency, and the frequency datum  $l_i$  is then assigned to  
 237  $\mu_i$ , as  $\mu_i = \{j \mid j \in l_i\}$ .

238 3) **data neighborhoods.** Data neighbourhoods of  $D_i$ , de-  
 239 noted as  $\mathcal{N}_{D_i}$ , are defined by the set of all observations  
 240 that are similar to  $D_i$  with a given threshold  $\epsilon$ . It is  
 241 computed as shown in Eq. 4.

$$\mathcal{N}_{D_i} = \{D_j \mid Dist(D_i, D_j) \leq \epsilon \vee j \neq i\} \quad (4)$$

242 4) **Core data.** Datum  $D_i$  is known as core data if and only  
 243 if here is some minimum number of data  $\sigma_D$ , such that  
 244  $|\mathcal{N}_{D_i}| \geq \sigma_D$ .

245 5) **Shared data determination.** Upon the construction of  
 246 the clusters of data, the shared set has to be determined  
 247 of data between clusters. In Eq. 5, the shared sets for  
 248 data denoted as  $S$ , are defined.

$$S = \bigcup_{i=1, j>i}^k C_i \cap C_j, \quad (5)$$

249 where  $S^{i,j}$  is the shared set between clusters  $C_i$  and  $C_j$ .

250 1) **Naive grouping for data decomposition:** For naive  
 251 groupings, the main aim is to be able to group data into  $k$   
 252 clusters that are disjoint without the need for any processing.  
 253 With  $m$  datum,  $\{D_1, D_2, \dots, D_m\}$ , the first  $\frac{m}{k}$  datum are  
 254 assigned to  $C_1$ , the second  $\frac{m}{k}$  to  $C_2$ , and so until assigning  
 255 all that datum to the  $k$  clusters.

256 2) **Hierarchical agglomerative clustering for data de-**  
 257 **composition:** HAC (Hierarchical Agglomerative Clustering)  
 258 [27] for data decomposition which has the main aim  
 259 in the creation of tree-like nested structure partitions,  
 260  $\mathcal{H} = \{\mathcal{H}_1, \mathcal{H}_2, \dots, \mathcal{H}_h\}$ , of the data such that,  $\forall (i, j) \in$   
 261  $[1, \dots, k]^2, \forall (m, l) \in [1, \dots, h]^2, C_i \in \mathcal{H}_m, C_j \in \mathcal{H}_l, m \geq$   
 262  $l \Rightarrow C_i \in C_j \wedge C_i \cap C_j = \emptyset$ . First, there is a starting  
 263 point with all data points in separate clusters. Next, we keep  
 264 connecting two clusters that can be agreed to be very similar  
 265 until we reach the point of a single cluster. We can define the  
 266 similarity between any two clusters  $C_i$  and  $C_j$  by determining  
 267 the number of common elements between them, or  $|C_i \cap C_j|$ .

268 3) **K-means for data decomposition:** We know that K-  
 269 means [28] is trying to optimize the function:  $J =$   
 270  $\sum_{j=1}^k \sum_{D' \in C_j} |D' - \mu_j|^2$ , where  $\mu_j$  is the centroid of the  
 271 data in  $C_j$ . A centroid is computed for each cluster, and  
 272 then the data are randomly distributed among  $k$  clusters.  
 273 Then, each datum is assigned to a cluster based on which  
 274 centroid is closest to it. These steps are repeated until no more  
 275 assignments to clusters are made, at which point the procedure  
 276 terminates itself.

277 4) *Bisecting k-means for data decomposition*: In the bisecting  
 278 *k*-means algorithm, when [29] decomposes the data, it does  
 279 so using both hybrid partitioning and a divisive hierarchical  
 280 methodology. We start with a single cluster and then split a  
 281 cluster into 2 in each individual step, using the standard *k*-  
 282 means approach. Looking more closely at the approach, the  
 283 process of bisecting clusters can be repeated many times, with  
 284 higher similarity achieved in the division.

285 5) *DBSCAN for data decomposition*: In the DBSCAN  
 286 algorithm [30], the main goal in data decomposition is to  
 287 be able to search for clusters in each  $\epsilon$  neighborhood per  
 288 datum. Once the core data is found, DBSCAN is responsible  
 289 for iteratively collecting all density-reachable data directly  
 290 from the core data. This process may result in some density  
 291 reachable clusters being merged individually. We can stop the  
 292 process if no new data is added to a cluster.

### 293 C. DL model

294 Here is presented a new DL model for detecting anomalies  
 295 in data. It is based on a recurrent neural network and considers  
 296 time series as input. The input of the recurrent neural network  
 297 is the set of clusters generated in the previous step. As a result,  
 298 different models are generated, each of which is associated  
 299 with a data cluster. Our model network is a (many-to-many)  
 300 architecture. The problem of the model is binary classification,  
 301 i.e., outputting a class label indicating whether the data is  
 302 anomalous or not. This is done for each datum in the cluster. A  
 303 multilayer feedforward network is applied to each data cluster,  
 304 consisting of multiple neurons arranged in layers. Each neuron  
 305 of layer  $l$  is connected to each neuron of layer  $(l - 1)$  with  
 306 a certain weight. Each input datum  $D_{i-1}$  is connected to a  
 307 group of neurons in the input layer. The neurons in the output  
 308 layer are associated with the output of the model (the class  
 309 label 1 for anomalous or 0 for normal). The goal is to reduce  
 310 the error between the output data of the model and the ground  
 311 truth of the data, such as:

$$E(D) = \sum_{i=1}^{|D|} E(D_i), \quad (6)$$

312 where,

$$E(D_i) = \sqrt{\sum_{j=1}^{|D_{ij}|} (D_{ij} - \widehat{D}_{ij})^2} \quad (7)$$

313 The output of the  $m^{th}$  neuron in the layer  $l$ , noted  $o_l^m$  is  
 314 given by Eq. 8. Note that the sum of the outputs of all neurons  
 315 in the given layer should be between 0 and 1. Here, we have  
 316 the following equations as:

$$o_l^m = \sigma\left(\sum_{j=1}^{|l-1|} o_{l-1}^j \omega_{l-1}^{mj} + b_l^m\right), \quad (8)$$

317 with

$$\sum_{m=1}^{|l|} o_l^m = 1, \quad (9)$$

318 where  $\sigma(\cdot)$  is the activation function,  $|l|$  is the number of  
 319 neurons in the layer  $l$ ,  $o_{l-1}^j$  is the output of the  $j^{th}$  neuron

in the  $l - 1$  layer,  $\omega_{l-1}^{mj}$  is the weight value that connects the  
 neurons  $o_l^m$  and  $o_{l-1}^j$ , and  $b_l^m$  is the bias value associated to  
 the neuron  $o_l^m$ .

At each iteration  $i$ , the updating weight rule is given as by:

$$\omega_{l-1}^{mj}(i) = \omega_{l-1}^{mj}(i-1) - \mu \times D_i \times 2 \times E_i, \quad (10)$$

where  $\mu$  is the learning parameter rate, and,

$$E_i = \sum_{j=1}^{|D_i|} (D_{ij} - \widehat{D}_{ij})^2 \quad (11)$$

At the end of the learning step, different models will be  
 designed, and one for each cluster,  $C_i$ . We define a local  
 ranking vector  $Rank_i$  by applying a learning model  $M_i$  on  
 the cluster  $C_i$ , denoted  $Rank_i = M_i(C_i)$ . The process of the  
 global ranking of the data  $D$  is performed as follows:

- 1) Compute the score of each  $D_j$ , say  $Score(D_j)$ .
- 2) Sort the scores of the data,  $D$ , in an ascending order.
- 3) Retrieve the top anomalous according to the scores of  
 $D$ .

### 334 D. Evolutionary Computation

335 In this section, we can show the process by which we  
 336 can determine the optimal set for the D2E-ADN approach to  
 337 finding the set of hyperparameters. Here we can define a set  
 338 of hyperparameters given by  $\mathcal{HP} = \{\mathcal{HP}_1, \mathcal{HP}_2, \dots, \mathcal{HP}_r\}$ .  
 339 Here  $r$  is defined as the total number of hyperparameters. Each  
 340  $\mathcal{HP}_i$  can be represented in a set of possible values for a given  
 341 hyperparameter. Moreover, we define our configuration space  
 342  $\mathcal{CS}$  such that we can say that the set of possible configurations  
 343 where each configuration can be represented as a vector in  
 344 the possible values for all hyperparameters  $\mathcal{HP}$ . Thus, the  
 345 hyperparameter problem for optimization has the main goal  
 346 of finding an optimal configuration that provides the highest  
 347 accuracy for both the regression and classification rates. We  
 348 can also say that the size of the configuration space can depend  
 349 on the number of possible values of the hyperparameters. We  
 350 can use Eq. 12 such that:

$$|\mathcal{CS}| = \prod_{i=1}^r |\mathcal{HP}_i|. \quad (12)$$

351 Here we can clearly see that the configuration space can  
 352 be very large. For example, if only 1,000 possible values  
 353 per epoch parameter and 100 per error rate and 1,000 for  
 354 the number of bounding boxes (i.e. CNN) are considered,  
 355 then the configuration space could be as large as 100 million.  
 356 Therefore, we need to be able to avoid exhaustive search  
 357 approaches as they are inappropriate for this type of problem.  
 358 To solve this problem, evolutionary computational algorithms  
 359 need to be explored. In the following, we discuss the main  
 360 components of such approaches.

361 1) *Population Initialization*: Considering the initial pop-  
 362 ulation represented as *pop\_size*, the individuals must be  
 363 distributed over the configuration space  $\mathcal{CS}$ . This allows ex-  
 364 ploration of different configurations and coverage of most  
 365 regions in  $\mathcal{CS}$ . When generating the initial population, we

366 can start the process by generating a random individual that  
 367 can represent a configuration  $\mathcal{CS}$ . This individual can then  
 368 generate  $pop\_size - 1$  individuals, keeping in mind that each  
 369 new individual can be dissimilar to the already generated  
 370 individuals. The dissimilarity of two configurations can be  
 371 easily determined by the distance between the configurations  
 372 of the individuals in question. We can also say that the  
 373 initial population, given as  $\mathcal{P}$ , should be able to maximize  
 374 the diversification function using the Eq. 13.

$$Diversify(\mathcal{P}) = \sum_{i=1}^{|\mathcal{P}|} \sum_{j=1}^{|\mathcal{P}|} Distance(\mathcal{CS}_i, \mathcal{CS}_j), \quad (13)$$

375 where we note here that  $Distance(\mathcal{CS}_i, \mathcal{CS}_j)$  is defined as  
 376 the distance between  $i^{th}$ , and  $j^{th}$  individuals configurations,  
 377 respectively.

378 2) *Crossover*: For the generation of any new offspring, we  
 379 must ensure that the steps as follows are applied:

- 380 • A crossover point is generated at random which ranges  
 381 from 1 to  $r$ , creating a *left side* and *right side* split.
- 382 • *left side* of first individual can be transferred to *left side*  
 383 of first offspring. However, *right side* of first individual  
 384 can be copied to *right side* of second offspring.
- 385 • *left side* for second individual can be copied to *left side*  
 386 for second offspring. Moreover, *right side* of second  
 387 individual can be copied to *right side* of first offspring.

388 3) *Mutation*: The diversification of the search is increased  
 389 by a mutation operation. By itself, the technique consists only  
 390 in randomly changing the parameter values for each config-  
 391 uration. Once a random mutation point has been generated,  
 392 which can range from 1 to  $r$ , future mutation point values can  
 393 be generated using the crossover operator.

394 4) *Local Search*: The local search tool starts with the  
 395 individuals of the population and returns the neighbors. The  
 396 neighbors are defined by updating the number of a parameter  
 397 to the current setting. This process is repeated for all individ-  
 398 uals of the population, with a high number of repetitions.

399 5) *Fitness Function*: As mentioned earlier, the D2E-ADN  
 400 approach aims to jointly maximize the regression and clas-  
 401 sification ratios. With this in mind, a multicriteria function  
 402 is proposed to be used when evaluating individuals from the  
 403 populations as in Eq. 14.

$$Fitness(\mathcal{CS}_i) = \frac{\alpha \times CR(\mathcal{CS}_i) + \beta \times RR(\mathcal{CS}_i)}{2}. \quad (14)$$

404 We note here that,

- 405 •  $\mathcal{CS}_i$  can be defined as the configuration of  $i^{th}$  individual  
 406 in population.
- 407 •  $CR(\mathcal{CS}_i)$  can be defined as the classification ratio of  
 408 D2E-ADN algorithm using  $\mathcal{CS}_i$ .
- 409 •  $RR(\mathcal{CS}_i)$  can be defined as the regression ratio of D2E-  
 410 ADN algorithm using  $\mathcal{C}_i$ . We note here that  $RR(\mathcal{CS}_i)$   
 411 can be set to 0 for RNN use.
- 412 •  $\alpha$  and  $\beta$  can be defined as 2 user parameters that are set  
 413 between 0.0 and 1.0.

414 Using the above operations, 2 algorithms are proposed for  
 415 the hyperparameter optimization methods. In the first case, a

416 genetic approach is used, and in the second case, a swarm  
 417 optimization method is used. It is shown that both approaches  
 418 are efficient when used with large populations.

TABLE I  
PERCENTAGE (%) OF THE SHARED DATA OF THE CLUSTERING STEP FOR  
THE D2E-ADN FRAMEWORK

Dataset	naive grouping	HAC	kmeans	bisecting kmeans	DBSCAN
Odense	42	40	5	7	30
Beijing	40	39	9	11	31
ICSX2012	39	37	7	18	24
CICIDS2017	45	31	8	10	21

419 6) *Genetic Algorithm*: The initial population of individuals  
 420 of size  $pop\_size$  is first randomly generated. Each individual is  
 421 constructed with respect to the initialization of the population.  
 422 Then, the crossover, mutation, and local search operators are  
 423 applied to generate configurations from  $\mathcal{CS}$ . To maintain the  
 424 same size of the population, all individuals are evaluated using  
 425 the fitness function and only the first  $pop\_size$  individuals (in  
 426 terms of quality) are left while the others are removed. The  
 427 identical procedure is continued until the predefined maximum  
 428 number of iterations is reached.

TABLE II  
DETECTION RATIO OF THE DL STEP FOR THE D2E-ADN FRAMEWORK

Dataset	Epochs 100	Epochs 1,000	Epochs 10,000
Odense	0.65	0.70	0.70
Beijing	0.70	0.72	0.72
ICSX2012	0.70	0.73	0.73
CICIDS2017	0.71	0.72	0.72

TABLE III  
FITNESS COMPUTING OF THE EVOLUTIONARY COMPUTATION STEP FOR  
THE D2E-ADN FRAMEWORK

Dataset	Genetic Algorithm	Bees Swarm Optimization
Odense	0.78	0.79
Beijing	0.77	0.80
ICSX2012	0.80	0.79
CICIDS2017	0.81	0.79

429 7) *Bees Swarm Optimization Algorithm*: First, a bee  
 430 searches for a good feature configuration. After this initial  
 431 configuration is found, a set of configurations *SearchArea*  
 432 in the search space using Eq. 13. Each individual particle  
 433 viewed from the *SearchArea* is the starting point for the search.  
 434 After a local search process is complete, each individual bee  
 435 passes its "best visited" configuration to all neighboring bees  
 436 using a table known as *Dance*. In the *Dance* table, a stored  
 437 configuration then becomes the next reference for the next  
 438 iteration. To ensure that no cycles occur, each new reference  
 439 configuration is added to a tab list that must never be used  
 440 as a starting reference again. If, after several iterations, it is  
 441 determined that the swarm does not improve its configuration,  
 442 the diversification criterion is introduced to avoid trapping the  
 443 local optimum. Usually, the diversification criterion consists  
 444 of a distant configuration that is not stored in the tabu list.  
 445 The algorithm usually ends when the optimal version is found  
 446 or a maximum number of iterations is reached.

#### IV. EXPERIMENTAL EVALUATION

447

448 Several experiments were conducted to validate the useful-  
 449 ness of the proposed framework using two real case studies.  
 450 The first is urban traffic anomalies used in intelligent trans-  
 451 portation and the second is intrusion detection for securing  
 452 World Wide Web technologies. Evaluation measures include  
 453 detection accuracy using the F-measure [32] and runtime.  
 454 All experiments were implemented on a 128 – *bit* Core i9  
 455 processor with UBUNTU 20 and 32GB from RAM used in  
 456 conjunction with a GPU device, an NVIDIA Tesla C2086 with  
 457 534 CUDA cores (16 multiprocessors with 64 cores each) and  
 458 a clock speed of 2.15GHz. There is 3.2GB of global memory,  
 459 59.15KB of shared memory, and a warp size of 64. Both the  
 460 CPU and GPU use single precision.

##### A. Datasets

461 **Urban Traffic Anomaly Detection:** Two real urban traffic  
 462 datasets were used: i) The first was obtained from Odense  
 463 Municipality (Denmark)<sup>1</sup>. This is a set of lines containing  
 464 information about the detection of cars and their locations. The  
 465 flows were observed between 1<sup>st</sup> January 2017 and 30<sup>th</sup> April  
 466 2018 and consist of more than 12 million cars and bicycles.  
 467 ii) The second one is from the Beijing traffic flow and was  
 468 retrieved from Beijing City Lab<sup>2</sup>. It consists of more than 900  
 469 million traffic flow values during two months in one place.  
 470 The anomalies in these two datasets are the set of traffic flows,  
 471 which may be a single traffic value or a sequence of traffic  
 472 values in a given time window.  
 473

474 **Intrusion Detection:** Many intrusion detection datasets,  
 475 such as KDD and DARPA, have been widely used over the past  
 476 two decades. However, these datasets are outdated and do not  
 477 reflect current security attacks in modern computer networks,  
 478 which are characterized by the emergence of IoT-generated  
 479 traffic. The ISCX2012<sup>3</sup> data were recently generated to reflect  
 480 current attack scenarios on networks. They consist of seven  
 481 days of real malicious and normal network activity. The normal  
 482 network traffic is generated by normal operations, while  
 483 the attack scenarios are performed with human assistance  
 484 to minimize misunderstandings with normal network traffic.  
 485 There are four different attack options such as penetrating the  
 486 network from inside, Hypertext Transfer Protocol Denial of  
 487 Service, Distributed Denial of Service using botnets and Brute  
 488 Force Secure Shell. The second data used is CICIDS2017  
 489 [33], which contains labeled network flows in CSV format.  
 490 They were collected over a five-day period and include some  
 491 cutting-edge attack scenarios such as brute force file transfer  
 492 protocol, brute force secure shell, denial of service attack, web  
 493 attack, infiltration, and botnet.

##### B. D2E-ADN Parameter Setting

494 **1) Decomposition:** The first experiment aims to evaluate,  
 495 on different datasets, the quality of the following decomposi-  
 496 tion algorithms: intuitive grouping, HAC, *k*-means, bisecting  
 497

<sup>1</sup><https://www.odense.dk/>

<sup>2</sup><https://www.beijingcitylab.com/>

<sup>3</sup><http://www.unb.ca/cic/datasets/index.html>.

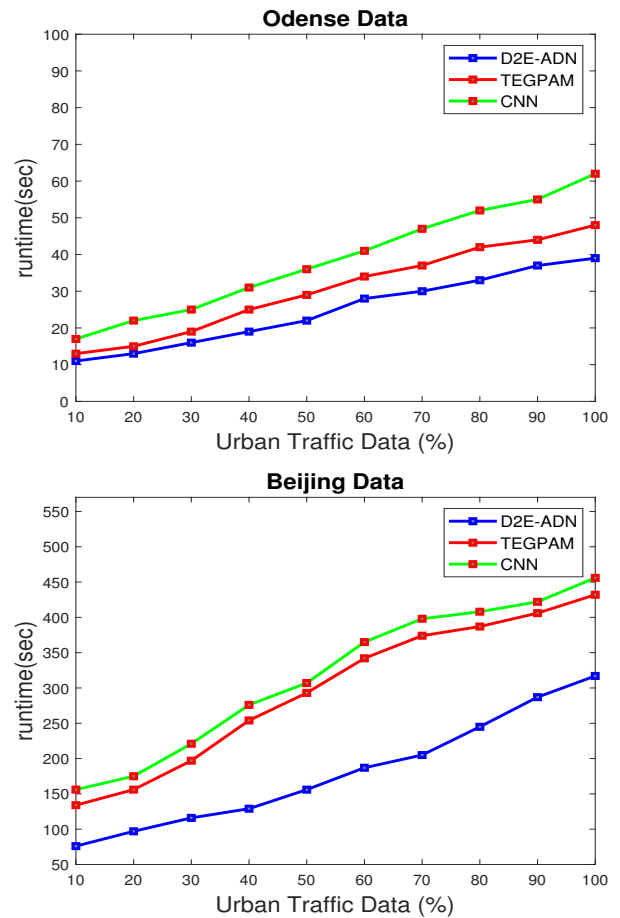


Fig. 2. Runtime in seconds of D2E-ADN versus state-of-the-art urban traffic anomaly detection algorithms

498 *k*-means and DBSCAN. This is determined by the percentage  
 499 of separation data between clusters/groups, while high quality  
 500 is reflected by low values of this percentage. The number of  
 501 clusters was varied from 1 to 50 for the Naive Grouping and  
 502 *k*-means algorithms, and the  $\epsilon$  value was varied from 1 to 10  
 503 for the DBSCAN algorithm. In this experiment, the optimal  
 504 parameter values for each clustering method are used and  
 505 are shown in Table I. Note that the number of clusters 5  
 506 for intuitive clustering, 7 for *k*-means and bisecting *k*-means,  
 507 12 for HAC, and  $\epsilon$  for DBSCAN was set to 4. The number  
 508 of separation data with the best parameter values for each  
 509 database is presented. The results show that *k*-means and  
 510 bisecting *k*-means provide better decomposition into records  
 511 compared to the other three algorithms. These results can be  
 512 explained by the fact that *k*-means and bisecting *k*-means  
 513 are pure partitioning, i.e., both algorithms are oriented to the  
 514 centroids representing the data of the same cluster. DBSCAN,  
 515 on the other hand, is inspired by computing neighborhoods  
 516 to represent dense regions. Consequently, it is conceivable  
 517 that two datasets are comparable and belong to the two  
 518 closest clusters. In the following tests, we use the *k*-means  
 519 decomposition technique of our framework.

520 **2) Performance of DL Model:** Here, we are concerned with  
 521 computing the quality of the DL step of 1) the convolutional  
 522 neural network for urban traffic anomaly detection and 2) the

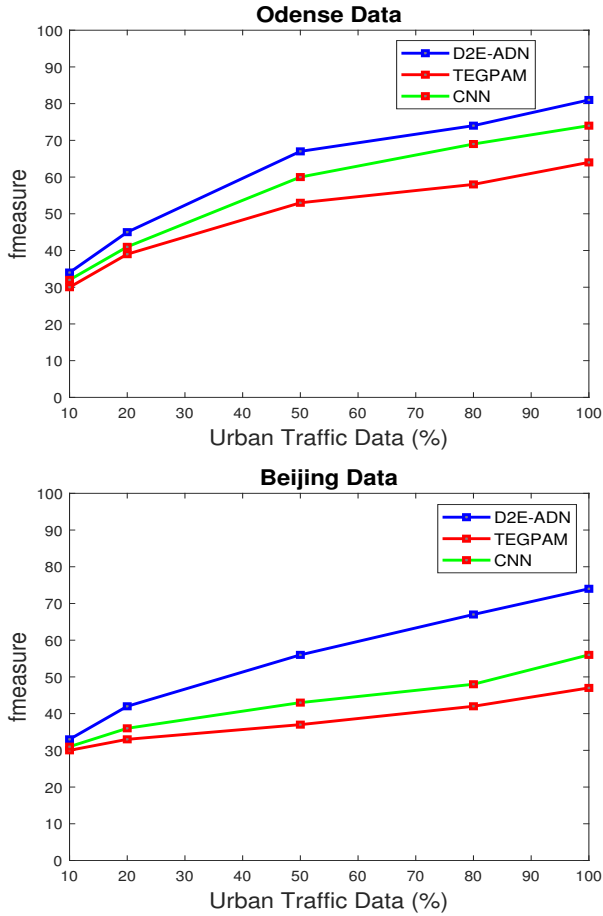


Fig. 3. Accuracy of D2E-ADN versus the state-of-the-art urban traffic anomaly detection algorithms

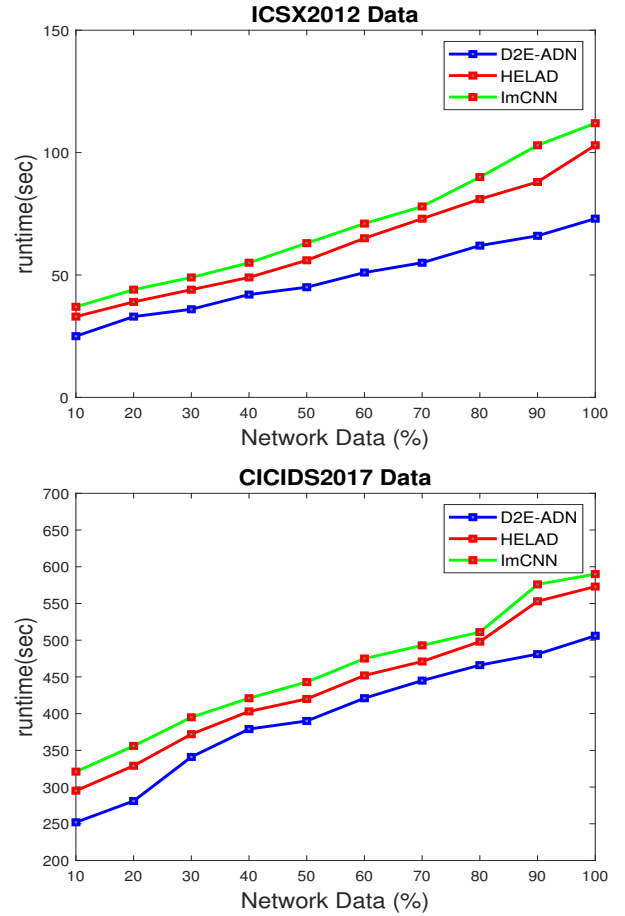


Fig. 4. Runtime in seconds of D2E-ADN versus the state-of-the-art intrusion detection algorithms

523 recurrent neural network for intrusion detection. The quality  
 524 is determined by the detection rate, which is the ratio between  
 525 the number of detected outliers and the number of all outliers.  
 526 If you vary the number of epochs of the network from 100  
 527 to 10,000, Table II shows that the detection rate of both  
 528 algorithms increases up to 1,000 and then converges at this  
 529 value. The reason for these results is that the weights of both  
 530 models became stable after 1,000 iterations. Therefore, the  
 531 best epochs for both algorithms are 1,000, which is used in  
 532 the rest of the experiments.

533 3) *Evolutionary Computation*: In this part, the quality of  
 534 the evolutionary computational step in genetic algorithms  
 535 and bee swarm optimization is evaluated. This quality is  
 536 determined by the best value of the fitness calculation of the  
 537 final population.

538 By varying the number of individuals/bees from 1 to 100  
 539 and the maximum number of iterations from 1 to 100, the best  
 540 parameter values for each evolutionary computation algorithm  
 541 are used in this experiment and listed in Table III. Note that the  
 542 number of individuals and the maximum number of iterations  
 543 are 35 and 47, respectively, for the genetic algorithm, while the  
 544 number of bees and the maximum number of iterations are 43  
 545 and 59, respectively, for the swarm optimization algorithm.  
 546 The results show that the genetic algorithm is better for  
 547 intrusion detection and the bee swarm optimization is better

548 for urban anomaly detection. In the remaining experiments,  
 549 we used the genetic algorithm for intrusion detection and the  
 550 bee swarm optimization for urban traffic anomaly detection.

### C. Results for Urban Traffic Anomaly Detection 551

552 In this experiment, we compare the performance of the D2E-  
 553 ADN algorithm with TEGPAM [8], and CNN [34], as baseline  
 554 urban traffic anomaly detection algorithms.

555 1) *Runtime*: In Fig. 2, the running time in seconds of D2E-  
 556 ADN is shown in comparison to the baseline algorithms. It  
 557 shows that the running time of the three algorithms increases  
 558 with the percentage of data. For 10% of data, all algorithms  
 559 require less than 200 seconds to identify outliers and more  
 560 than 350 seconds to process the entire data. The results also  
 561 show the superiority of our approach compared to the other  
 562 two algorithms, with a difference of more than 100 seconds for  
 563 processing the entire data. These results were obtained thanks  
 564 to the efficient combination of the convolutional neural network  
 565 with the decomposition algorithms in deriving anomalies  
 566 from the urban traffic data.

567 2) *Accuracy*: In Fig. 3, the F-measure of the D2E-ADN is  
 568 shown in comparison with the baseline algorithms. It shows  
 569 that the F-measure increases with the percentage of data in  
 570 the three algorithms. Most importantly, it shows the clear superior-  
 571 ity of D2E-ADN with an advantage of more than 15 points in

572 processing the whole data. These results are obtained thanks to  
 573 the efficient combination of the convolutional neural network  
 574 with the evolutionary computation in the optimization of the  
 575 hyperparameters. Thus, finding the appropriate parameters for  
 576 learning the network can significantly improve the detection  
 577 rate of outliers.

#### 578 D. Results for Intrusion Detection

579 This part compares D2E-ADN with HELAD [6] and Im-  
 580 CNN [35], as two baseline algorithms for network intrusion  
 581 detection.

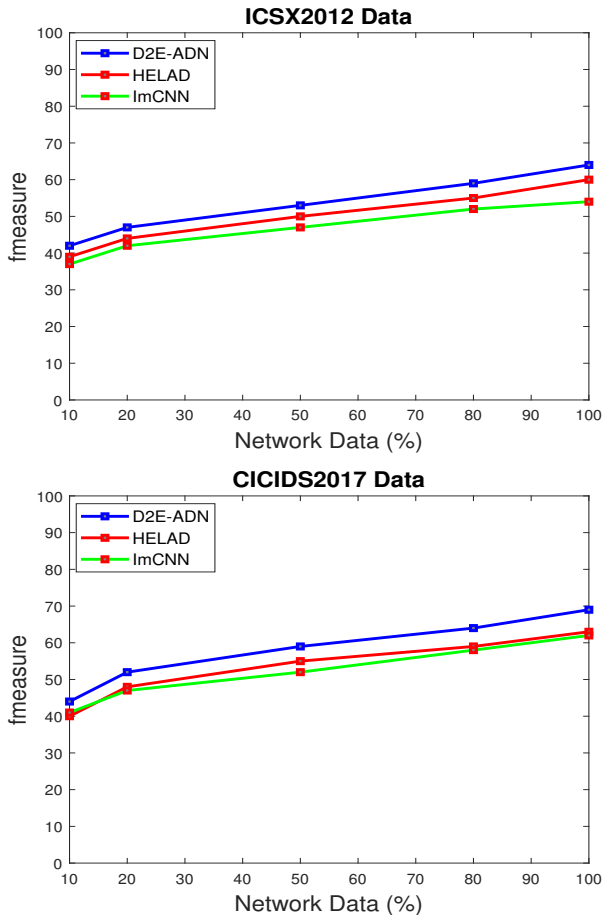


Fig. 5. Accuracy of D2E-ADN versus state-of-the-art intrusion detection algorithms

582 1) *Runtime*: Fig. 4 shows the runtime in seconds of D2E-  
 583 ADN, HELAD, and ImCNN on ICSX2012 and CICIDS2017  
 584 datasets. The results show that the runtime of the three  
 585 algorithms increases with the percentage of data. This has  
 586 significant implications, e.g., all algorithms require less than  
 587 250 seconds to identify an anomaly from 10% of the data, but  
 588 more than 550 seconds to process the entire data. The results  
 589 also show the superiority of the proposed approach (D2E-  
 590 ADN) compared to the other two algorithms, with a difference  
 591 of more than 150 second for processing the whole data. These  
 592 results are obtained thanks to the efficient combination of the  
 593 recurrent neural network with the decomposition algorithms  
 594 in deriving anomalies from the urban traffic data. Any RNN

that learns from homogeneous data can significantly increase  
 the performance in detecting outliers.

2) *Accuracy*: The F-measure of D2E-ADN compared with  
 the baseline algorithms (HELAD and ImCNN) is shown in Fig.  
 5. The results show that the F-measure of the three algorithms  
 increases with the percentage of data. They also reveal the  
 superiority of D2E-ADN, which offers the advantage of more  
 than 12 points for processing the whole data. These results are  
 obtained thanks to the efficient combination of the recurrent  
 neural network with the evolutionary computation in the  
 optimization of the hyperparameters of our algorithm.

#### V. CONCLUSION

In this work, we studied the problem of anomaly detection  
 in IoE and proposed a combination of decomposition, deep  
 neural networks and evolutionary computation to find anomalies  
 from the dataset. In our approach, the dataset is first  
 decomposed into similar clusters using different types of  
 clustering algorithms. The clusters are then trained using an  
 extended recurrent neural network. To perform the training  
 step efficiently, two evolutionary computation algorithms are  
 proposed to take the hyper-parameters of the trained models  
 and try to find the optimal ones. Several experiments in the  
 form of two case studies for two different IoE applications  
 show the advantages of the proposed solution compared to the  
 basic approaches. In perspective, we plan to explore other  
 data representations such as trajectories. We also plan to  
 propose a parallel version that explores high-performance  
 computing to increase the performance of the proposed  
 solution and train the data clusters simultaneously. In  
 addition, the current work can be extended to other subsets  
 of the digital IoT world. Although IoE is a recent  
 development, other areas within IoT can be explored using  
 the concepts presented in this paper. For example, both the  
 Internet of Vehicles (IoV) and the Internet of Smart  
 Infrastructures (III) could be a future home for the  
 research presented here. In this context, in addition to the  
 datasets used here, other novel datasets can be used to  
 further test and refine the work already done.

#### ACKNOWLEDGMENT

This research was partially funded by the Natural Sciences  
 Research Council of Canada (NSERC) Discovery Grant  
 program (RGPIN-2020-05363) held by Dr. Gautam Srivastava.

#### REFERENCES

- 1) Y. Miao, X. Liu, K. R. Choo, R. H. Deng, H. Wu, and H. Li, "Fair and dynamic data sharing framework in cloud-assisted internet of everything," *IEEE Internet Things Journal*, vol. 6, no. 4, pp. 7201–7212, 2019.
- 2) M. N. Bhuiyan, M. M. Rahman, M. M. Billah, and D. Saha, "Internet of things (iot): A review of its enabling technologies in healthcare applications, standards protocols, security and market opportunities," *IEEE Internet of Things Journal*, 2021.
- 3) D. Djenouri, R. Laidi, Y. Djenouri, and I. Balasingham, "Machine learning for smart building applications: Review and taxonomy," *ACM Computing Surveys*, vol. 52, no. 2, pp. 24:1–24:36, 2019.
- 4) B. Ouyang, P. S. Wills, Y. Tang, J. O. Hallstrom, T.-C. Su, K. Namuduri, S. Mukherjee, J. I. Rodriguez-Labra, Y. Li, and C. J. Den Ouden, "Initial development of the hybrid aerial underwater robotic system (haucs): Internet of things (iot) for aquaculture farms," *IEEE Internet of Things Journal*, 2021.



- [5] S. Boukhaboul and D. Djenouri, "DFIOT: data fusion for internet of things," *J. Neww. Syst. Manag.*, vol. 28, no. 4, pp. 1136–1160, 2020.
- [6] Y. Zhong, W. Chen, Z. Wang, Y. Chen, K. Wang, Y. Li, X. Yin, X. Shi, J. Yang, and K. Li, "Helad: A novel network anomaly detection model based on heterogeneous ensemble learning," *Computer Networks*, vol. 169, p. 107049, 2020.
- [7] R. Doriguzzi-Corin, S. Millar, S. Scott-Hayward, J. Martinez-del Rincon, and D. Siracusa, "Lucid: A practical, lightweight deep learning solution for ddos attack detection," *IEEE Transactions on Network and Service Management*, 2020.
- [8] L. Lin, J. Li, F. Chen, J. Ye, and J.-P. Huai, "Road traffic speed prediction: A probabilistic model fusing multi-source data," *IEEE Transactions on Knowledge and Data Engineering*, 2017.
- [9] T. Kieu, B. Yang, C. Guo, and C. S. Jensen, "Outlier detection for time series with recurrent autoencoder ensembles," in *The International Joint Conference on Artificial Intelligence*, 2019, pp. 2725–2732.
- [10] A. Karale, M. Lazarova, P. Koleva, and V. Poulkov, "A hybrid pso-milof approach for outlier detection in streaming data," in *The International Conference on Telecommunications and Signal Processing*. IEEE, 2020, pp. 474–479.
- [11] Y. Djenouri, A. Belhadi, P. Fournier-Viger, and J. C. W. Lin, "Fast and effective cluster-based information retrieval using frequent closed itemsets," *Information Sciences*, vol. 453, pp. 154–167, 2018.
- [12] A. Belhadi, Y. Djenouri, J. C.-W. Lin, C. Zhang, and A. Cano, "Exploring pattern mining algorithms for hashtag retrieval problem," *IEEE Access*, vol. 8, pp. 10 569–10 583, 2020.
- [13] K. Pawar and V. Attar, "Deep learning approaches for video-based anomalous activity detection," *World Wide Web*, vol. 22, no. 2, pp. 571–601, 2019.
- [14] F. A. Khan, A. Gumaiei, A. Derhab, and A. Hussain, "A novel two-stage deep learning model for efficient network intrusion detection," *IEEE Access*, vol. 7, pp. 30 373–30 385, 2019.
- [15] K. Al Jallad, M. Aljndi, and M. S. Desouki, "Big data analysis and distributed deep learning for next-generation intrusion detection system optimization," *Journal of Big Data*, vol. 6, no. 1, p. 88, 2019.
- [16] A. Pektaş and T. Acarman, "A deep learning method to detect network intrusion through flow-based features," *International Journal of Network Management*, vol. 29, no. 3, p. e2050, 2019.
- [17] S. Garg, K. Kaur, N. Kumar, and J. J. Rodrigues, "Hybrid deep-learning-based anomaly detection scheme for suspicious flow detection in sdn: A social multimedia perspective," *IEEE Transactions on Multimedia*, vol. 21, no. 3, pp. 566–578, 2019.
- [18] R. M. A. Ujjan, Z. Pervez, K. Dahal, A. K. Bashir, R. Mumtaz, and J. González, "Towards sflow and adaptive polling sampling for deep learning based ddos detection in sdn," *Future Generation Computer Systems*, 2019.
- [19] D. Papamartzivanos, F. G. Mármol, and G. Kambourakis, "Introducing deep learning self-adaptive misuse network intrusion detection systems," *IEEE Access*, vol. 7, pp. 13 546–13 560, 2019.
- [20] M. A. Ferrag, L. Maglaras, S. Moschoyiannis, and H. Janicke, "Deep learning for cyber security intrusion detection: Approaches, datasets, and comparative study," *Journal of Information Security and Applications*, vol. 50, p. 102419, 2020.
- [21] L. Boukela, G. Zhang, M. Yacoub, S. Bouzeffrane, S. B. B. Ahmadi, and H. Jelodar, "A modified lof-based approach for outlier characterization in iot," *Annals of Telecommunications*, vol. 76, no. 3, pp. 145–153, 2021.
- [22] A. E. Edje, S. M. Abd Latiff, and H. W. Chan, "Enhanced non-parametric sequence-based learning algorithm for outlier detection in the internet of things," *Neural Processing Letters*, vol. 53, no. 3, pp. 1889–1919, 2021.
- [23] M. R. Nosouhi, K. Sood, N. Kumar, T. Wevill, and C. Thapa, "Bushfire risk detection using internet of things: An application scenario," *IEEE Internet of Things Journal*, 2021.
- [24] W. Zhang, Q. Lu, Q. Yu, Z. Li, Y. Liu, S. K. Lo, S. Chen, X. Xu, and L. Zhu, "Blockchain-based federated learning for device failure detection in industrial iot," *IEEE Internet of Things Journal*, vol. 8, no. 7, pp. 5926–5937, 2020.
- [25] J. C.-W. Lin, G. Srivastava, Y. Zhang, Y. Djenouri, and M. Aloqaily, "Privacy-preserving multiobjective sanitization model in 6g iot environments," *IEEE Internet of Things Journal*, vol. 8, no. 7, pp. 5340–5349, 2020.
- [26] D. Chou and M. Jiang, "A survey on data-driven network intrusion detection," *ACM Computing Surveys (CSUR)*, vol. 54, no. 9, pp. 1–36, 2021.
- [27] W. H. Day and H. Edelsbrunner, "Efficient algorithms for agglomerative hierarchical clustering methods," *Journal of classification*, vol. 1, no. 1, pp. 7–24, 1984.
- [28] J. MacQueen *et al.*, "Some methods for classification and analysis of multivariate observations," in *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, vol. 1, no. 14, 1967, pp. 281–297.
- [29] M. Steinbach, G. Karypis, and V. Kumar, "A comparison of document clustering techniques," in *KDD Workshop on Text Mining*, vol. 400, no. 1. Boston, 2000, pp. 525–526.
- [30] M. Ester, H.-P. Kriegel, J. Sander, X. Xu *et al.*, "A density-based algorithm for discovering clusters in large spatial databases with noise," in *Proceedings of KDD*, 1996, pp. 226–231.
- [31] Y. Djenouri, D. Djamel, and Z. Djenouri, "Data-mining-based decomposition for solving MAXSAT problem: Towards a new approach," *IEEE Intelligent Systems*, 2017.
- [32] S. Ramaswamy, R. Rastogi, and K. Shim, "Efficient algorithms for mining outliers from large data sets," in *ACM SIGMOD Record*, vol. 29, no. 2, 2000, pp. 427–438.
- [33] A. Pektaş and T. Acarman, "Classification of malware families based on runtime behaviors," *Journal of information security and applications*, vol. 37, pp. 91–100, 2017.
- [34] L. Zhu, R. Krishnan, A. Sivakumar, F. Guo, and J. W. Polak, "Traffic monitoring and anomaly detection based on simulation of luxembourg road network," in *IEEE Intelligent Transportation Systems Conference*, 2019, pp. 382–387.
- [35] S. Garg, K. Kaur, N. Kumar, G. Kaddoum, A. Y. Zomaya, and R. Ranjan, "A hybrid deep learning-based model for anomaly detection in cloud datacenter networks," *IEEE Transactions on Network and Service Management*, vol. 16, no. 3, pp. 924–935, 2019.