

Online Research @ Cardiff

This is an Open Access document downloaded from ORCA, Cardiff University's institutional repository: <https://orca.cardiff.ac.uk/id/eprint/147062/>

This is the author's version of a work that was submitted to / accepted for publication.

Citation for final published version:

Kayan, Hakan, Majib, Yasar, Alsafery, Wael, Barhamgi, Mahmoud and Perera, Charith ORCID: <https://orcid.org/0000-0002-0190-3346> 2021. AnomL-IoT: an end to end re-configurable multi-protocol anomaly detection pipeline for Internet of Things. Internet of Things 16 , 100437. 10.1016/j.iot.2021.100437
file

Publishers page: <http://dx.doi.org/10.1016/j.iot.2021.100437>
<<http://dx.doi.org/10.1016/j.iot.2021.100437>>

Please note:

Changes made as a result of publishing processes such as copy-editing, formatting and page numbers may not be reflected in this version. For the definitive version of this publication, please refer to the published source. You are advised to consult the publisher's version if you wish to cite this paper.

This version is being made available in accordance with publisher policies.
See

<http://orca.cf.ac.uk/policies.html> for usage policies. Copyright and moral rights for publications made available in ORCA are retained by the copyright holders.



AnoML-IoT: An End to End Re-configurable Multi-protocol Anomaly Detection Pipeline for Internet of Things

Hakan Kayan^a, Yasar Majib^a, Wael Alsafery^a, Mahmoud Barhamgi^b, Charith Perera^a

^aCardiff University, UK

^bClaude Bernard Lyon 1 University, France

Abstract

The rapid development in ubiquitous computing has enabled the use of microcontrollers as edge devices. These devices are used to develop truly distributed IoT-based mechanisms where machine learning (ML) models are utilized. However, integrating ML models to edge devices requires an understanding of various software tools such as programming languages and domain-specific knowledge. Anomaly detection is one of the domains where a high level of expertise is required to achieve promising results. In this work, we present AnoML which is an end-to-end data science pipeline that allows the integration of multiple wireless communication protocols, anomaly detection algorithms, deployment to the edge, fog, and cloud platforms with minimal user interaction. We facilitate the development of IoT anomaly detection mechanisms by reducing the barriers that are formed due to the heterogeneity of an IoT environment. The proposed pipeline supports four main phases: (i) data ingestion, (ii) model training, (iii) model deployment, (iv) inference and maintaining. We evaluate the pipeline with two anomaly detection datasets while comparing the efficiency of several machine learning algorithms within different nodes. We also provide the [source code](#) of the developed tools which are the main components of the pipeline.

Keywords: Internet of Things, Data Science, Pipeline, Data Analytics, Multi-Protocol

1. Introduction

Edge AI which is critical for resource-constrained environments that operates in the Internet of Things (IoT) domain where intelligent tasks are performed has started to become a hot topic with the arrival of Industry 4.0 [1]. It manages the interaction with the physical world that is provided by sensors and actuators. Management of such an environment requires series of tasks (e.g., data collection, anomaly detection) that are operated by microcontrollers running ML models. Data-related professions (e.g., data scientists, ML engineers) define rules/ranges and search for the best practices to increase the operability of edge mechanisms in their relevant scientific disciplines. Finding hidden information from big data can enhance the quality of living but it is not a straightforward task [2].

While for data scientists, being an expert in edge-related infrastructures (e.g., programming languages, microcontrollers, sensors) is not expected, they should be able to utilize data science pipelines which are executable workflows of data-related tasks that automate the desired process. Thus, we developed a reconfigurable data science pipeline based on an IoT sensing infrastructure that utilizes open-source software to facilitate developing an interconnected anomaly detection

16 system that runs on edge, fog, and cloud platforms. We define the edge as the platform where the
17 first interaction between the cyber and physical world happens. Hence, microcontrollers (e.g.,
18 Raspberry Pi Pico) that gather physical data are edge devices. We define fog as the platform
19 where several edge devices can be supervised. Hence, single-board computers (e.g., Raspberry
20 Pi 4B) are fog devices that might act as edge devices as well. Cloud is the platform where real-
21 world data gathered by the edge and fog devices are progressed. We implemented our system
22 based on an example use case scenario to describe how to proposed system works while providing
23 some results.

24 The contributions of this paper are as follows:

- 25 • We provide reconfigurable IoT sensing infrastructure that consists of two main open-
26 source components: (i) The Edge to Cloud Code Generator (EECG) that generates ready-
27 to-deploy codes to enable data circulation from edge to fog. (ii) The Node-RED package
28 is hosted on Node-RED servers that enables accessing and processing to the edge data
29 from anywhere that has access to Node-RED servers while offering visualization via the
30 graphical user interface (GUI). We also provide one Python library and executable shell
31 script that facilitate data training and inference phases.
- 32 • We propose a data science pipeline that interconnects edge, fog and cloud devices/services
33 to provide end-to-end anomaly detection system development. The pipeline contains four
34 main stages: (i) the data collection which is provided by the components mentioned at the
35 first contribution point, (ii) the anomaly detection model training, (iii) model deployment
36 to the edge, fog, and cloud, (iv) inference, and maintaining the model based on the new
37 data. We demonstrate how the proposed tools are utilized during these stages.
- 38 • We provide a dataset that is generated via the utilization of proposed tool. We analyze
39 the performance of Convolutional Neural Network (CNN) [3], Recurrent Neural Network
40 (RNN) [4], Isolation Forest [5] and One-class Support Vector Machines (OC-SVM) [6]
41 on the proposed dataset [7] and the WADI dataset [8]. We also evaluate them according
42 to the platform (edge, fog or cloud) where the anomaly detection model is deployed via
43 the utilization of proposed pipeline. In the edge, we only evaluated CNN due to lack of
44 application programming interface (API).

45 **Structure of the Paper:** This section provides a high-level understanding of what we pro-
46 posed. We outlined the previous commercial and academic works in section 2. Section 3 contains
47 the architecture of an IoT anomaly detection pipeline infrastructure. Section 4 presents details
48 about how the data circulated and progressed within the AnomML-IoT pipeline. We demonstrate
49 our evaluations and results in section 5. Then, we discuss about the results in the section 6 and
50 finally provide our conclusions in section 7.

51 2. Related Work

52 In this section, we introduce the data science pipelines that are offered either by academia
53 or commercial entities, and anomaly detection techniques in time-series sensor data. We also
54 analyze the capabilities of open-source platforms that facilitates data circulation.

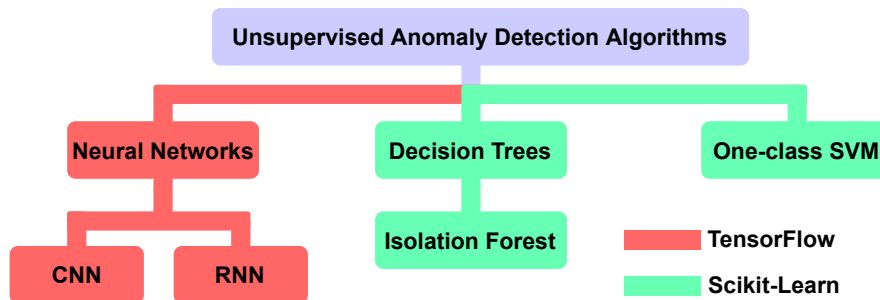


Figure 1: Illustrates the utilized algorithms. TensorFlow also has an API [27] for decision trees, but due to having better documentation we prefer using scikit-learn for implementing Isolation Forest.

55 *2.1. Unsupervised Anomaly Detection in Time Series Sensor Data*

56 Anomaly detection is one of the fundamental fields that utilize the machine learning (ML)
 57 model as the main component. There is extensive research being done in this field [9, 10, 11].
 58 There are three types of anomalies: (i) point anomalies, (ii) contextual anomalies, (iii) collective
 59 anomalies. If the anomalies are contextual where the context is time, the time series anomaly
 60 detection models are applied. For example, in an environment where the weather temperature
 61 decreases at night if the temperature value generated by the sensor acts otherwise, there is a con-
 62 textual anomaly. While point anomalies are easier to detect, contextual and collective anomaly
 63 detection requires additional tasks to identify the normal behavior of the system.

64 The nature of the input data is the core element that determines the efficiency of the ML
 65 model. The features of the data may depend on several complementary terms such as labels, con-
 66 text, and domain. For example, if the input data do not contain any labels that define normality,
 67 unsupervised algorithms [12] are applied, if the data is related to a certain context, context-aware
 68 [13] methods are selected, if the environment is industrial, because of the importance of detection
 69 time, faster models with reduced complexity [14] are preferred.

70 In an interconnected domain such as IoT, cyber-physical systems [15] are utilized to supervise
 71 the environment. These systems observe the behavioral changes (e.g., change in the temperature
 72 or movement) in surroundings through modules that manage sensors [16]. They can also act as
 73 controllers if they contain actuators. In such environments, the anomaly might occur either by
 74 independent or dependent events. If the events are independent, univariate analysis [17] is ap-
 75 plied. For example, the behavioral changes in temperature, loudness, light density, and humidity
 76 can be detected via the related data only, hence require univariate analysis. However, changes in
 77 the angular momentum or acceleration are measured by sensors (e.g., accelerometer, gyroscope)
 78 that generate data per dimension. Hence, the relation between the data points should also be
 79 analyzed to detect anomalies. Then, the multivariate analysis [18] is applied.

80 While academia keeps offering new anomaly detection algorithms [19, 20], most of the time
 81 these are based on the fundamental ones [21]. Hence, for this work, we selected the following
 82 algorithms as they are the most common ones that are utilized for the unsupervised anomaly
 83 detection in time series data and accessible via common ML programming libraries/frameworks
 84 (e.g., scikit-learn [22], TensorFlow [23]) : (i) convolutional neural networks (CNN) [24, 25], (ii)
 85 recurrent neural networks (RNN) [26], (iii) isolation forest (IF) [5], and (iv) one-class support
 86 vector machines (OC-SVM) [6]. Figure 1 demonstrates the used algorithms in this study.

88 **Machine learning platforms.** After showing promising results in a variety of tasks includ-
89 ing speech recognition, image processing, anomaly detection, and medical diagnosis, ML has
90 taken an interest in both academic and commercial entities, hence resulting in the creation of
91 many open-source and proprietary ML platforms and pipelines. ML models can be generated via
92 hand-coding, code generators, or interpreters. *Hand coding.* There are many machine learning
93 libraries available [22, 23, 28] that allows user to create ML models or to deploy and evaluate ML
94 algorithms. A person with the profession might prefer hand-coding as it offers high customiza-
95 tion, allows the development and employment of novel algorithms, and is easy to maintain.
96 However, hand-coding might be very resource-consuming, thus most of the time it is done by a
97 group of programmers. *Code generators.* ML consists of many steps (e.g., data acquisition, data
98 pre-processing, and fitting). Rather than hand-coding all these steps, code generators [29, 30]
99 might be utilized to facilitate the process. Due to the variety of complicated tasks, most code
100 generators provide a specific code for a specific task. *Interpreters.* One of the main challenges
101 of ML is the portability of the generated model. Interpreters provide portability by generating a
102 model file that can be run on other platforms with minimal coding. TensorFlow [23] is the most
103 common one that offers model generation for resource-constrained platforms.

104 **Data science pipelines.** Raw data are needed to be interpreted to be utilized within data
105 science-related tasks. If the data science pipeline contains all the steps that are required to in-
106 terpret the data from data gathering to deployment of a machine learning model, it is called
107 end-to-end. These end-to-end pipelines can be either manual where the user provides many in-
108 puts and sets parameters each time before a new model is generated or automated where little to
109 no input is taken. Due to a variety of data types, automated pipelines put a certain set of rules
110 (e.g., time format) for their system to accept the input data [31]. These pipelines can also be
111 named according to the performed tasks (e.g., anomaly detection pipeline). Now we introduce
112 pipelines that are presented by either industry or academia.

113 *Azure Machine Learning Pipeline* [32]. Microsoft provides an ML pipeline based on run-
114 ning Python scripts on the cloud while automatically handling resource usage. Each step of the
115 pipeline can be independently customized hence offering scalability to the end-user. One of the
116 practical features that Azure Machine Learning Pipeline offers is the automated dependency hand-
117 ling that allows the usage of a variety of hardware and software environments. Microsoft also
118 provides Azure Cognitive Services [33] where you can utilize their ML pipeline and Anomaly
119 Detector [34] service. They apply Graph Attention Network (GAN) [35] for multivariate analy-
120 sis, apply SR-CNN [31] for the univariate analysis.

121 *Amazon Web Services (AWS) Machine Learning Pipeline* [36]. Amazon provides an end-
122 to-end ML pipeline as a service for detecting anomalies in real-time. Inside the pipeline, there
123 are many different services (e.g., database, data formatting) that can be utilized for pipeline
124 tasks. Amazon SageMaker [37] is the main service that provides anomaly detection for both
125 univariate and multivariate data. It allows users to either use a built-in unsupervised anomaly
126 detection algorithm based on Random Cut Forest (RCF) [38] or use a custom algorithm that can
127 be deployed via a Docker image. Now we introduce the pipelines proposed by the academia.

128 Prado et al. [39] propose an end-to-end modular AI pipeline that allows users with less ex-
129 pertise to implement their AI applications such as keyword spotting, image classification, and
130 object detection to systems that contain embedded devices. Their framework relies on Low
131 Power Deep Neural Network (LPDNN) that contains an Inference Engine (LNE) that is compat-
132 ible with Caffe [40]. LNE is a code generator that facilitates the deployment to the embedded

133 devices. The authors use FIWARE [41] for IoT hub integration and Kurento Media server [42]
134 for media streaming which are required to run live inference. The authors define Raspberry Pi
135 devices as edge and evaluate the efficiency of LPDNN compared TF Lite [43] on these devices
136 by running benchmarks that are included in the TF Lite repository.

137 Drori et al. [44] propose an automatic ML (AutoML) system that optimizes the ML pipeline
138 according to the given dataset. Their pipeline utilizes LSTM-RNN as a base ML algorithm.
139 Monte Carlo Tree Search (MCTS) [45] is applied to the predictions generated by the LSTM-RNN
140 to evaluate the performance of the pipeline and decide on the better pipeline. They evaluate the
141 proposed pipeline compared to baseline stochastic gradient descent (SGD) [46] estimators from
142 scikit-learn [22]. They claim their pipeline provides faster run time according to its peers.

143 Sutton et al. [47] propose an open-source ML pipeline that receives physiological data that
144 is used to identify anomalous behaviors as an input in real-time. The authors try to detect Parox-
145 ysmal atrial fibrillation (PAF) by applying Probabilistic Symbolic Pattern Recognition (PSPR)
146 to the Electrocardiogram (ECG) signals. PSPR is used for online feature extraction while they
147 apply random forest (RF) to classify ECG data. The proposed pipeline is based on Spark's ML
148 library (MLlib) [48], hence allows other anomaly detection techniques included within MLlib.

149 Nitsche and Halbritter [49] propose a data science pipeline that is optimized for text classi-
150 fication. The authors benchmark different GPUs to evaluate the performance of their hardware
151 setup which consists of 10 NVIDIA Quadro P6000 and the effect of the number of GPUs on the
152 image processing time. They apply the Naive Bayes classifier that is included in scikit-learn API
153 and achieve above 90% accuracy on Deutsche Presse-Agentur (dpa) dataset.

154 Shaikh et al. [50] focus the challenges of ensuring policy fairness within end-to-end ML
155 pipelines. They claim the ML-based tasks are done by engineers that have a variety of professions
156 including data creators and future engineers. Hence, each step of the ML pipeline might be
157 subjected to a policy violation. The authors provide an end-to-end ML pipeline that is based
158 on log management to prevent these violations as manually ensuring policy fairness is highly
159 resource-consuming.

160 Boovaraghavan et al. [51] propose an adaptive end-to-end ML system for IoT applications.
161 Their pipeline is optimized for activity recognition-based tasks including object recognition. Au-
162 thors claim that the main challenge regarding end-to-end pipeline is due to the heterogeneity of
163 IoT applications. Authors evaluate their pipeline with various hardware platforms and datasets
164 while comparing prediction time and accuracy per each machine learning technique they applied.

165 Molinara et al. [52] propose an end-to-end ML-based indoor air monitoring system for con-
166 taminant classification. Authors compare the performances of Multi Layer Perceptron (MLP) to
167 CNN and LSTM based deep learning techniques while testing the performance of MLP and CNN
168 on ESP32 MCU. They investigate the power consumption of the MCU regarding the utilized ML
169 technique. They claim the proposed system is only lacked to classify alcohol and acetone due to
170 their chemical similarities.

171 Vinzamuri et al. [53] propose an end-to-end context-aware anomaly detection system that
172 requires time-series data. The proposed system utilizes a semi-supervised algorithm with Sparse
173 Gaussian Graphical Models. They benchmark the pipeline on several public datasets. The au-
174 thors claim semantics can improve the Gaussian Graphical Models further beyond other anomaly
175 detection techniques. Their ML comparison is based on F-Score as the authors mention that the
176 proposed pipeline is promising for industrial IoT environments.

177 Li et al. [54] develop an end-to-end automated anomaly detection system. They utilize
178 Apache Spark backend server to run the query-based operations. After the user provides a
179 dataset, the proposed system automatically selects the most appropriate algorithm then applies

Table 1: Comparison of AnoML-IoT with the Previous Works

Related Work	Topic	Environment	Commercial	Open-source	End-to-End	Time-series Data	Adaptability
Azure Machine Learning Pipeline [32]	General	General	✓		✓	✓	
AWS Machine Learning Pipeline [36]	General	General	✓		✓	✓	✓
Prado et al. [39]	Classification	IoT	✓		✓	✓	
Sutton et al. [47]	Anomaly Detection	Medical		✓		✓	
Nitsche and Halbritter [49]	Classification	Linguistics					
Shaikh et al. [50]	Policy	General			✓	✓	
Boovaraghavan et al. [51]	Classification	IoT			✓	✓	
Molinara et al. [52]	Classification	IoT			✓		
Vinzamuri et al. [53]	Anomaly Detection	IoT			✓	✓	
Li et al. [54]	Anomaly Detection	General			✓	✓	✓
AnoML-IoT	Anomaly Detection	IoT		✓	✓	✓	✓

180 anomaly detection. Finally, the results are shown in figures within the pipeline. They benchmark
 181 the proposed system based on several datasets while applying quantification analysis.

182 Our pipeline utilizes scikit-learn and TensorFlow for the anomaly detection while relying on
 183 TensorFlow Lite for the deployment on the fog, TensorFlow Lite Micro for the deployment on the
 184 edge devices. Hence, it allows the lightweight implementation of anomaly detection techniques
 185 while offering a variety of communication protocols (e.g., Bluetooth low energy (BLE)) and
 186 sensor types for the inference. The Table 1 compares our pipeline with the previous works based
 187 on significant features that determine the efficiency of the pipeline.

188 3. The Architecture of the AnoML-IoT Pipeline

189 In this part, we present the overall architecture of our pipeline by defining main components,
 190 describing the workflows that differ according to application scenario, and explaining how to
 191 automate these workflows to maintain the pipeline.

192 3.1. AnoML-IoT Layers and Application Scenarios

193 IoT infrastructures should provide semantic data exchange to be considered as completely
 194 ubiquitous. Current technologies that are utilized in IoT architectures are rapidly evolving to
 195 achieve semantic interoperability, hence causing the debate of what kind of infrastructure is
 196 needed for a certain task. Due to each technology has a variety of pros and cons per IoT en-
 197 vironment, extensive testing is required to decide on IoT elements (e.g., wireless technologies,
 198 edge/fog node types, sensors, and actuators). Building an IoT application from scratch to perform
 199 these tests requires intensive labor. Thus, reconfigurable IoT sensing architecture that allows the
 200 implementation of various ML-based application scenarios including anomaly detection with
 201 minimum user (e.g., data scientist) interaction is required. AnoML-IoT allows the implemen-
 202 tation of a variety of application scenarios where the scenarios are evolved around the platform
 203 that the anomaly detection is implemented. In this context, AnoML-IoT consists of three layers:

- 204 • *Edge Layer*: The edge layer contains edge nodes that consist of microcontrollers, sensors,
 205 and actuators. These nodes are physically observing the IoT environment by gathering data
 206 via sensors while conducting physical operations via actuators (e.g., the fan stops working
 207 when a certain degree is reached).
- 208 • *Fog Layer*: The fog layer contains fog nodes that get the data from edge nodes and process
 209 it according to the end user’s preferences. Usually, fog nodes offer more computing power

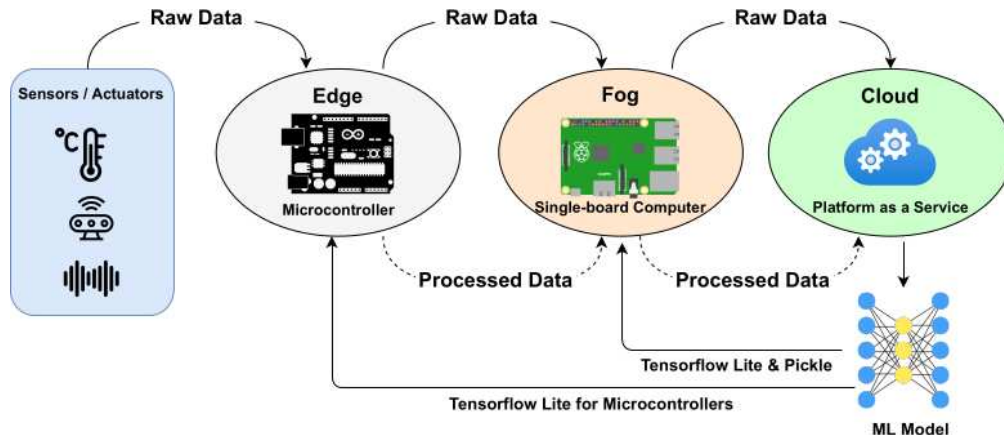


Figure 2: Overview of the data circulation and the pipeline application scenarios based on the location of anomaly detection: (i) If the anomaly detection is done on the cloud, the edge and fog devices only forward data, (ii) if the anomaly detection is done on the fog, the edge sends raw data to fog, then the fog might send processed data to the cloud, (iii) if the anomaly detection is done on the edge, the processed data might be sent to fog and then to cloud. Each scenario has its pros and cons that we introduce in section 5. Here processed data contains information to be used to decide if the data is anomalous or not. Thus, it might be either binary (e.g., 0 for normal data, 1 for anomalous data) or a decimal(float) as an anomaly score that varies in the range of -1 to 1. The model in the cloud will be updated with the new normal data according to the predetermined intervals.

210 than edge nodes while allowing flexible deployment options. Even though fog computing
 211 may be considered as an alternative to cloud computing, fog nodes can also act as an IoT
 212 gateway between the edge nodes and the cloud. In this work, we utilize Raspberry Pi 4 as
 213 a fog device and present our results. We believe, similar Linux-based devices can be used
 214 as fog devices instead of the Raspberry Pi for the proposed pipeline.

- 215 • *Cloud Layer*: The cloud layer provides services ranging from data management, storing,
 216 applying anomaly detection to developing multi-purpose frameworks. Even though, cloud
 217 computing offers many benefits (e.g., automatic service integration, and high accessibil-
 218 ity), edge/fog computing is preferred where time-critical or confidential applications (e.g.,
 219 industrial) are present to reduce reliance on cloud services. While the edge and fog layer
 220 might contain standalone nodes, the cloud layer requires interaction with other layers.

221 Figure 2 summarizes the working principle of the machine learning pipeline. The anomaly
 222 detection can be done within the pipeline on edge, fog, or cloud. The edge and fog devices also
 223 can be utilized just to forward raw data to the platform at one upper level. The initial training
 224 requires prior data. Hence, if there is no dataset to be used for initial training, the edge and fog
 225 will send only raw data until the cloud can generate an efficient ML model. Then the model will
 226 be deployed to edge, fog, or cloud. According to the given intervals, the new model will replace
 227 the old model to keep the system up-to-date. Having an up-to-date ML model is significant to
 228 prevent a decrease in efficiency that depends on the dynamic context.

229 3.2. The AnoML-IoT Open-Source Tools

230 AnoML-IoT consists of three main tools: (i) Edge to Cloud Code Generator (ECCG), (ii)
 231 Node-RED package, (iii) python library. Now we introduce what these tools provide, what are

232 their functionalities, use cases, and roles within the AnoML-IoT pipeline. The source code of
233 the tools is published in GitLab¹.

234 3.2.1. The Edge to Cloud Code Generator

235 The ECCG is a web-based code generator that generates edge code to be either used for
236 inference or to transmit data to the fog node. Currently, it is only compatible with Arduino IDE
237 [55]. It has a user-friendly interface, where even the data scientist with minimal IoT knowledge
238 can design a basic IoT application that contains several sensors where the sensor data can be
239 transferred between layers. The characteristics of ECCG including the justification for design
240 choices are given below:

- 241 • Currently, five sensor types are available: *temperature, humidity, air quality, light, loud-*
242 *ness*. The end-user can simultaneously select all sensors. Selected sensor data will be
243 included within the pipeline. *Justification*. These five sensor types are among the most
244 common sensors that are utilized in IoT applications. The generated code clearly describes
245 how the sensor data is received, and processed, thus allowing the easy adaptation of the
246 code for a similar type of sensor.
- 247 • Four communication protocols are available: *Wi-Fi* [56], *Bluetooth Classic* [57], *BLE* [58],
248 and *Zigbee* [59]. The end-user can only select one. The application also includes additional
249 settings regarding communication protocols for advanced users. *Justification*. The selected
250 four communication protocols occupy the vast majority of the IoT market and offer a
251 variety of topologies (e.g., mesh, star, tree). Understanding the basics of how to establish
252 these wireless technologies enables the implementation of high-range IoT applications.
253 The generated code by ECCG defines how to handle the required network elements (e.g.,
254 MAC address, personal area network (PAN) ID). An advanced user may conveniently
255 adapt the generated code to be used with other IoT communication protocols that are out-
256 of-scope of this project such as WirelessHART [60].
- 257 • The data transfer rate determines the time between two data blocks. It is in milliseconds
258 that ranges from 30000 to 300000. After copying to code the user can set the data transfer
259 rate as desired. However, we do not recommend setting it below 30 seconds as it is the
260 time that is required for the module to initialize. *Justification*. Data generating time differs
261 per sensor module. Controlling the data transfer rate is necessary to ensure the integrity of
262 transmitted data.
- 263 • The sensor locations are determined by the end-user. Identification (ID) number starting
264 from 00 is given per location. The application supports up to 99 locations. *Justification*.
265 Data scientists work with comma-separated value (CSV) files, to handle the further pro-
266 gressing of the data. The data in CSV files mostly in pairs as *text-value* where *text* is the
267 identifier, and *value* is digital presentation of a physical quantity. To generate such files,
268 data is transmitted in data-serialization formats (e.g., JavaScript Object Notation (JSON),
269 Extensible Markup Language (XML)) supported by IoT application standards. Thus, the
270 ECCG allows an end-user to define location. Then, it generates a unique location ID num-
271 ber to be used in data transmission.

¹<https://gitlab.com/IOTGarage/anoml-iot-analytics>

The Edge to Cloud Code Generator

Main Inputs

The screenshot shows the 'Main Inputs' section of the ECCG interface. It is divided into several functional areas:

- Select Sensor Types:** A row of checkboxes for Temperature (checked), Humidity (checked), Air Quality (checked), Light (unchecked), and Loudness (unchecked).
- Select Communication Protocol:** Radio buttons for WiFi (selected), Classic Bluetooth, BLE, and Zigbee.
- Time Interval:** A dropdown menu showing '30000'.
- MAC Address:** A text input field with the placeholder 'Ex. 00:0E:EA:CF:4A:23'.
- Show Advanced Settings:** A button to expand the interface.
- Edge Node Location:** A text input field containing 'Kitchen'.
- Location ID:** A list of three options: '00. Garden', '01. Bathroom', and '02. Kitchen'. There are 'Add' and 'Remove' buttons next to this list.
- Select Location:** A dropdown menu with 'Kitchen' selected.

Figure 3: Illustrates the main input section of the ECCG. Each input is strictly controlled for the following reasons: (i) to prevent cross-site scripting attacks, (ii) to prevent bugs that may occur in the code due to mistyping. Each text input has a tooltip that clarifies what kind of information should be given. In addition, placeholders demonstrate an example input.

- 272
- 273
- 274
- 275
- 276
- 277
- 278
- 279
- Three microcontroller types are available: Arduino Nano 33 BLE Sense, Arduino Nano RP2040 Connect, and Raspberry Pi Pico. The end-user can select one of the microcontrollers to obtain the edge node ID number to be used to identify the microcontroller types when needed. *Justification.* The ECCG supports the top three microcontrollers that are officially supported by TensorFlow Lite for Microcontrollers (TFLM) [61]. Supporting a variety of microcontrollers allows users to design a heterogeneous IoT application, where the environment benefits from different features of these devices (e.g., RAM, flash memory).
 - The ECCG allows sending lowest, mean, or highest sensor data which are generated during the time interval determined by the user. *Justification.* In some cases, the normal range might just be determined by lower or upper limits. Hence, we allow the user to decide on the data to be sent. If the user selects mean, the mean of the number of data points generated by the sensor during the predetermined time interval will be sent.

285

286

287

288

The user interface of the ECCG facilitates usability by navigating users via input controls. Figure 3 demonstrates the main input section of the code generator while the example inputs are given. Minimalist design is preferred to assist a data scientist with no prior IoT knowledge. Thus inputs are controlled, tooltips are included, and example inputs are given as placeholders.

289

290

291

292

293

294

295

296

297

298

After the inputs are given, the final step is clicking on the "Generate Code" button. The ECCG will output the followings: (i) the edge code that is ready to be deployed to the microcontroller via Arduino IDE, (ii) the Python3 script that should be run on the fog device, (iii) the Linux commands to be run via terminal, (iv) example Node-RED setup. The outputs differ according to the given inputs. For example, the Linux commands are only required if the Bluetooth Classic will be used within the pipeline. The Figure 4 illustrates how the generated code blocks are presented. Two main buttons are included for each code block: the first button generates the code shown in the pre-scrollable division to let a data scientist examine the codes before further progression, the second button copies the code without breaking the format to prevent possible errors. The user manually uploads the code into the microcontroller.

299

300

301

302

303

304

We assume the following scenario: the edge node (micro-controller) gathers physically observed data via sensors and sends it to a fog node(a single-board computer (SBC)) where the data is either processed or forwarded to the cloud via Node-RED. Thus, the ECCG consists of three main sections: (i) the input field where the end-user determines the basic characteristics of the desired IoT application, (ii) the transmitter field where the edge node code is generated for a microcontroller, and (iii) the receiver field where the fog node code is generated for a single-board

Generate Code

Transmitter Code

```

/*Please replace with your own correct variables if you see "IR" near the given variable.*/
#include "DHT.h"
#include "Air_Quality_Sensor.h"

#define DHTPIN 4 //!R
#define DHTTYPE DHT11 //!R
DHT dht(DHTPIN, DHTTYPE);
AirQualitySensor sensor(A2); //!R
//send data per interval seconds(end user decides)
unsigned long interval = 30000L;
unsigned long pass_time = millis();
int exit_while = 0;

void setup()
{


```

Copy Transmitter Code

Receiver Code

Node-RED Example Setup

Example Design



TCP Input Node Settings

Type: Listen on port 4448

Output: stream of String payload(s)

delimited by \n

Topic: Topic

Name: Name

Figure 4: Illustrates the generated code block and the example Node-RED setup. In the example scenario, Raspberry Pi Pico is selected as an edge node/transmitter. The generated transmitter code is written in C++ and should be deployed via Arduino IDE. The Raspberry Pi 4B is acting as a fog node/receiver.

305 computer. Table 2 illustrates the specifications of the ECGG in detail.

306 3.2.2. The Node-RED Package

307 Node-RED [62] is an open-source browser-based workflow development tool that runs on
 308 Node.js [63] based runtime. We can develop workflows via utilizing drag-and-drop nodes on
 309 resource-constraint environments such as Raspberry Pi thanks to the non-blocking nature of
 310 Node.js. Hence, we use Node-RED on the fog device to develop and control the workflows
 311 within the AnoML-IoT pipeline. Even though Node-RED contains many open-source packages
 312 developed for the IoT networks, we could not find any up-to-date package that provides the devel-
 313 opment of Bluetooth Low Energy (BLE) modules while offering customization choices. Hence,
 314 we developed and published our package under the name of "node-red-contrib-IoT-procotols".
 315 The package contains the following nodes:

- 316 • **BLE Scanner.** BLE devices can operate in four different roles: broadcaster, observer, cen-
 317 tral, and peripheral. The task of the BLE Scanner node is to scan for the advertising BLE
 318 peripherals. It can specifically search for a local name and output one of the followings:
 319 whole peripheral as an object, MAC address, and advertisement data. So, the user can use

Table 2: The Specifications of The Edge to Cloud Code Generator

Main Inputs	Sensor Types	Temperature
		Humidity
	Communication Protocols	Loudness
		Light
		Air Quality
Edge Node Types	Wi-Fi	
	Bluetooth Classic	
	BLE	
	Zigbee	
Edge Node Identifiers	Arduino Nano 33 BLE Sense	
	Arduino Nano RP2040 Connect Raspberry Pi Pico	
Advanced Inputs	Wi-Fi	Edge Node Location
		Edge Node ID Number
		Service Set Identifier (SSID)
	Bluetooth & BLE	Password
		Host IP Address
		Host Port
		*MAC Address
Zigbee	Module Name	
	Module PIN	
	PAN ID	
Generated Outputs	Node-RED Example Setup	Destination Address High
		Destination Address Low
	Transmitter Code	Node Settings
		Workspace Design
		Example Functions
Receiver Code	Arduino Scripts	
	Python Scripts	
		JavaScript Function Nodes
		Python Scripts

* Among advanced inputs, only the MAC address is obligatory. The default options that are included in the code are explained via comments in detail.

320 this node either to discover the MAC address of the target BLE peripheral or to get the
321 advertisement data.

322 • **BLE Connect.** This node establishes a BLE connection with the target peripheral device.
323 While one peripheral device can only be connected to one central, central devices can
324 manage multiple peripherals. Hence, our pipeline allows the deployment of multiple edge
325 devices via BLE while the maximum number of peripherals that can be connected depends
326 on the system-on-a-chip (SoC) of the microcontroller. For example, the Arduino Nano 33
327 BLE Sense [64] SoC nRF52840 [65] supports up to 20 parallel connections.

328 Figure 5 introduces nodes that offer the establishment of multiple network communication
329 protocols in the Node-RED environment within the AnoML-IoT pipeline. Each communication
330 protocol might have different pros and cons according to the application domain. Hence, it is
331 important to offer multiple options to the end-user. We consider providing more choices in the
332 future.

333 3.2.3. The *AnoML.py* and *SetupAnoML.sh*

334 **AnoML.py.** Library of functions in python which can support most of the tasks in developing
335 ML models with various types of normalization and data points. Generating different models for

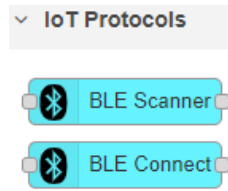


Figure 5: Illustrates the developed nodes to be used within to Node-RED environment hosted by the fog device.

336 the cloud, fog and edge platforms using various normalization methods and data points is a labor-
 337 intensive task. Evaluating different models for unsupervised ML when there is a small amount
 338 of anomalous data is also another challenge. Our Python library supports data preprocessing to
 339 generating ready to deploy ML models for the cloud, fog, and edge for unsupervised ML. It also
 340 provides performance visualization for each data point while including normalization techniques
 341 which helps to select the most efficient model. The proposed library provides data cleaning,
 342 normalization, reduction, scaler, and visualization functions before feeding into machine learning
 343 models. The library also provides various functions that allow training and testing. The inference
 344 part of the library provides functions to evaluate different machine learning models generated
 345 using our library at both fog and cloud platforms so the performance of the platforms can be
 346 compared.

347 **SetupAnoML.sh.** It acts as an installer to prepare both fog and cloud platforms for action.
 348 It installs correct versions of necessary packages and libraries required to run inference. Once
 349 the prerequisites are handled, it can then download, install and configure packages such as Node-
 350 RED and TensorFlow runtime. Also, it sets our custom-developed services to listen on web ports
 351 to order to allow microcontrollers to interact with fog devices and fog devices to interact with
 352 the cloud platform. The script can be configured to download files and configurations either
 353 from Google Drive or a local/remote FTP Server. Users also can manually download, install and
 354 configure all prerequisite requirements and model files if it is desired.

355 4. The Data Circulation

356 The nature of the input data shapes the characteristics of the data science pipeline. Hence,
 357 gathering raw data is the first step in this kind of pipeline. In the IoT environment, the sensors
 358 generate a variety of data in various formats. While most of the temperature sensors generate data
 359 in floating-point numbers (floats) to provide more accuracy, the light sensors that measure light
 360 density usually provide data in integer. This is significant for two main reasons: (i) each data type
 361 occupies memory in different sizes, (ii) there might be different regulations [66, 67] according
 362 to the application domain for certain data types such as floats. Hence, the memory usage should
 363 be optimized for microcontrollers that have very limited memory. The memory usage might also
 364 differ according to the microcontroller architecture. For instance, while Arduino Uno [68] stores
 365 int values in two bytes, Arduino Due [69] stores int values in four bytes. The data circulation
 366 within the AnoML-IoT pipeline consist of four main steps: (i) data ingestion, (ii) data training,
 367 (iii) model deployment, (iv) inference and maintaining.

368 4.1. Data Ingestion

369 Data ingestion is the first step of all kinds of data science pipelines unless the user already
 370 has a ready-to-deploy ML model. In our pipeline, the ECCG is the main tool that utilizes the

371 data ingestion process. Figure 3 illustrates the choices offered by the ECCG. The user can select among the offered choices to obtain a generated code. The generated code can be copied
 372 without disrupting the format. Multiple communication protocols and Node-RED example that
 373 demonstrates how to receive data via fog device are also provided.
 374

375 What the ECCG generates is an example edge code. Besides, in our setup, we use Grove
 376 sensors and shields [70] while utilizing Raspberry Pi 4B as a fog device. The IoT environment is
 377 very heterogeneous hence, the user might have or want to use different sensors, microcontrollers,
 378 or fog devices. The generated code by the ECCG should require minimal editing even when this
 379 is the case as we also provide instructions via comments for the exact lines within the code that
 380 might require editing due to individual preferences. If needed, the user can edit the generated
 381 code after copying it from the ECCG to Arduino IDE. Then, the code can be uploaded to a
 382 microcontroller. To complete these tasks the user needs a computer that can access the internet
 383 and run Arduino IDE. So, the tasks of copying the code and uploading it to a microcontroller are
 384 handled manually. Figure 6 demonstrates the data ingestion process.

385 4.2. Model Training

386 Model training can be done before ingestion if the user already has a training dataset. Other-
 387 wise, the user should follow the steps mentioned in the data ingestion phase. The AnomL.py
 388 library contains all the required functions that are needed to generate a ML model. Uploading
 389 the dataset to a cloud platform is handled manually. After uploading the dataset, the user should
 390 preprocess the data to convert to an appropriate format for the ML algorithms. Then, the ML
 391 models are generated based on the user preferences. User can generate multiple ML models at
 392 the same time. Here, we facilitate the model training by providing user a python library that is
 393 capable of preprocessing the time series data and generating multiple ML models based on given
 394 parameters with only a few line of codes. The required time depends on the capabilities of a
 395 cloud platform that is used during the training. The user should decide on the followings: (i) the
 396 size of the training dataset, (ii) the algorithm specific parameters(e.g., contamination factor for
 397 Isolation Forest). Figure 7 illustrates the model training phase.

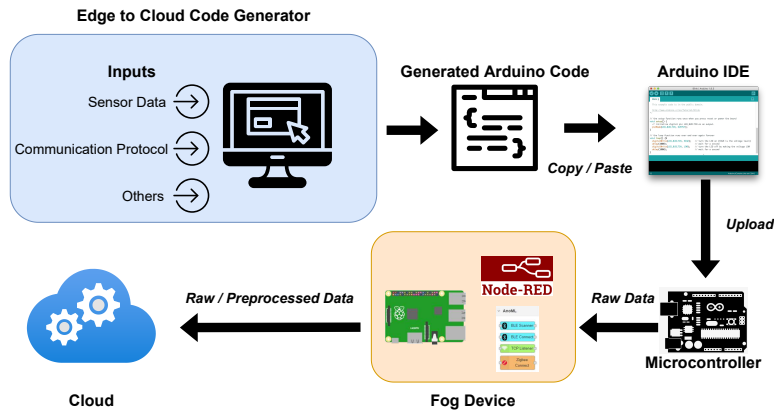


Figure 6: Illustrates data ingestion phase. The data preprocessing can be done on all platforms. However, due to the limited computing power of edge and fog devices, we prefer to utilize the cloud. In addition, if more than one edge device is connected to the fog device, data preprocessing can be done on the fog device to identify the edge devices.

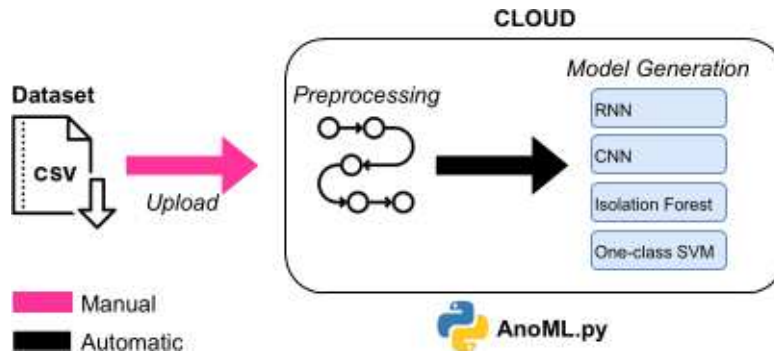


Figure 7: Demonstrates the model training phase. Currently, the built-in algorithms within the AnoML.py are RNN, CNN, isolation forest, and One-class SVM. Multiple datasets can be utilized at the same. The storage place of these models depends on the preferences of the user. During our evaluation we used Google Drive [71] to store our models as it can be easily integrated into Google Colab [72]. Another option would be running an FTP server. The user also can utilize both at the same time.

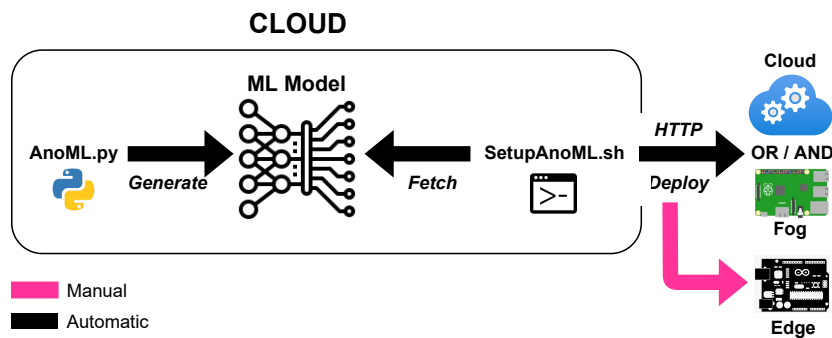


Figure 8: Illustrates the model deployment phase. SetupAnoML is an executable Linux shell script that does the followings in order: (i) fetches generated ML models from Google Drive or FTP server, (ii) installs required libraries and packages, (iii) deploys models over HTTP. The deployment to edge should be handled manually. We are considering to provide automated edge ML deployment function via over-the-air (OTA) transmission in the future versions of AnoML-IoT.

398 4.3. Model Deployment

399 The next step after training the model is model deployment. Here we provide an executable
 400 shell script SetupAnoML.sh that facilitates the model deployment phase by automating several
 401 key tasks. The user should provide the following inputs to the script: (i) the platform (edge
 402 or cloud) where the anomaly detection technique will be deployed, (ii) the details of the place
 403 where the model is stored (e.g., server, username, password, and port names for an FTP server,
 404 Google Drive token or Google Drive). After these inputs are given, the script will automatically
 405 download, install, and configure all the required software packages (e.g., libraries, models, con-
 406 figurations). Then the script automatically will deploy the models to cloud, fog, or both cloud
 407 and fog platforms. The user can deploy multiple models at the same time to different platforms.
 408 Figure 8 demonstrates the model deployment process.

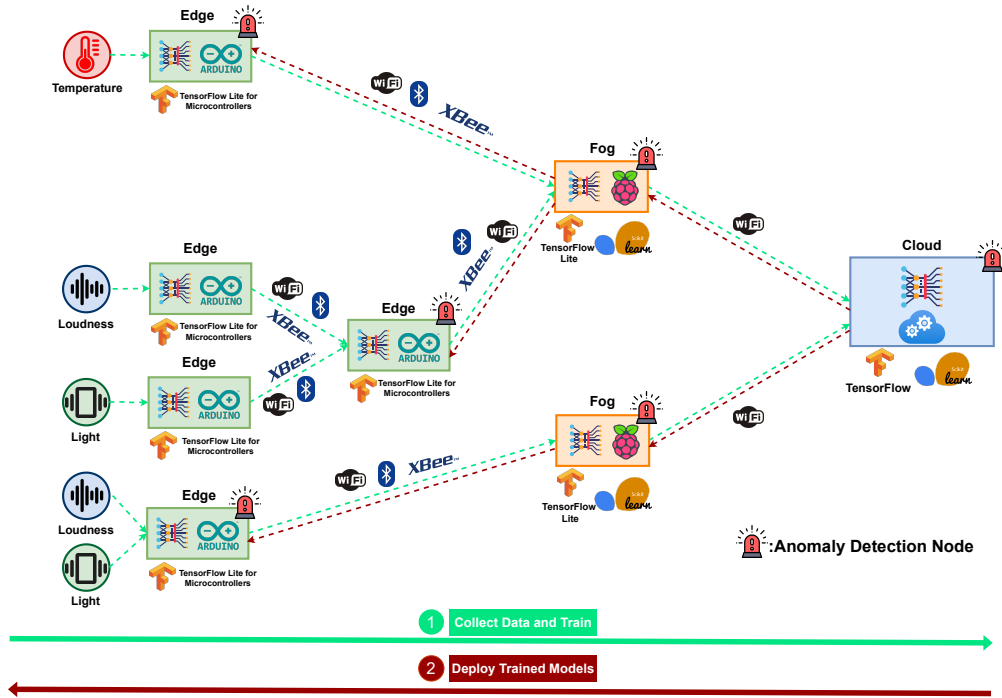


Figure 9: Illustrates the possible scenarios that can be generated within the AnoML-IoT pipeline. We allow the use of various topologies, wireless communication protocols, and execution of use case scenarios shaped around the decision of anomaly detection platform.

409 **4.4. Inference and Maintaining**

410 Successful deployment results in inferring all models on each type of platform. We developed
 411 an end-to-end pipeline which allows a user to generate ML models for anomaly detection
 412 from sensor data at all platforms. We developed a performance monitor which can present visu-
 413 alization of performance and accuracy of all type platforms, ML models, data-points and normal-
 414 ization/reduction techniques. Maintenance of data is a key aspect for evolution, we recommend
 415 that user should make it a practice to visualize performance comparison as a process of decision
 416 making in order to enhance performance and accuracy on all platforms.

417 **5. Evaluation**

418 In this section, we present the evaluations that are done to measure the efficiency of the
 419 pipeline with various configurations enabled. A data science pipeline evaluation is not a trivial
 420 task if the environment is heterogeneous such as IoT. Figure 9 demonstrates possible scenarios
 421 available within our pipeline where different sensors and wireless communication protocols are
 422 used.

423 **5.1. Wireless Protocol Comparison**

424 The decision on the wireless protocol selection depends on the application domain as each
 425 protocol has its benefits and disadvantages. There are several important features to be considered

Table 3: Wireless Protocol Comparison

	Latency (ms)				Power Consumption	Topology*
	Edge to Edge	Edge to Fog	Fog to Fog	Fog to Cloud		
Wi-Fi	18.24	14.566	17.25	21.23	High	Star
Bluetooth Classic	195.13	171.15	187.15	NA	Medium	Point-to-Point
BLE	11.23	13.45	13.21	NA	Low	Mesh
Zigbee	18.56	16.66	14.56	NA	Low	Mesh

NA: Not Applicable. *: The most common topology is given.

426 when designing an anomaly detection pipeline for the IoT environments:

- 427 • **Latency.** The time that takes for one network package to be transmitted from the transmit-
428 ting endpoint to the receiving endpoint. Latency optimization is very significant to achieve
429 near real-time processing and inference. It is one of the significant features that determine
430 the anomaly detection time. Low latency is aimed for the applications (e.g., industrial)
431 where the anomaly detection time is critical [73].
- 432 • **Power consumption.** IoT is a resource-constraint environment, hence power consumption
433 is one of the significant features to be considered.
- 434 • **Network Topology.** The data circulation and the system design depend on the network
435 topology. While there are many topologies (e.g., star, mesh, ring) offered by the current
436 communication protocols, the mesh, and star networks are the most common ones that are
437 established in IoT environments.

438 Table 3 compares the wireless communication protocols utilized in the proposed pipeline
439 based on the aforementioned features. We measured the latency by taking the mean of the passed
440 time for the 1000 transmitted packages. Edge and fog devices were placed next to each other
441 during the tests hence there was no physical barrier between the devices. No packet drop is
442 observed. While we observed similar latency for the BLE, Wi-Fi, and Zigbee, the Bluetooth
443 Classic had the highest latency. However, it is challenging to have a conclusion as many factors
444 that might affect the latency. Hence, to get the most realistic results, the tests should run in a real
445 environment where anomaly detection takes place.

446 5.2. Dataset and Testbed

447 Realistic datasets and testbeds are required to achieve the most promising results. To evaluate
448 the efficiency of the pipeline, we built a testbed that contains components that are low-cost and
449 accessible to most of the IoT community. When designing the IoT testbed, the sensor selection is
450 the initial task that shapes the main futures of the testbed. In cyber-physical environments, there
451 are several behaviors (e.g., temperature, noise, humidity) that can be considered as common.
452 Hence, we observe these common behaviors via our testbed and generate a dataset. Table 4
453 demonstrates the components utilized during the evaluation.

454 We recorded the temperature, humidity, light, loudness, air quality data for around two days.
455 We observed the data at the beginning to form an initial opinion about physical behavioral
456 changes in the test environment. We realized the only temperature and humidity sensors are
457 working as expected and let us creating anomalous behaviors. Hence, we only utilize these two
458 data during the evaluation. Figure 10 provides visualization of the generated dataset. Arduino

Table 4: Testbed Components

Sensors	Grove - Temperature & Humidity Sensor (High-Accuracy & Mini) v1.0 Grove - Light Sensor Grove - Loudness Sensor Grove - Air Quality Sensor v1.3 Grove - UART Wifi V2 Digi XBee 3 Zigbee 3 RF Module
Microcontrollers	Arduino Nano 33 BLE Sense Raspberry Pi Pico Arduino Nano RP2040 Connect
Shields	Grove - Bee Socket Grove Shield for Pi Pico V1.0 Arduino Tiny Machine Learning Shield
Single-board Computer	Raspberry Pi 4 Model B
Software Tools	Putty Raspberry Pi Imager Node-RED Arduino IDE 1.8.15 Visual Studio Code Google Colab

Table 5: Comparison of Datasets

Features	WADI [8]	AnoML [7]
Amount of Data	122,543,744	45,906
Number of Rows	957,374	6,559
Number of Columns	129	7
Number of Sensors & Actuators	123	5
Time-Series	✓	✓
Labelled	✓	✓
File Size	588,906 KB	265 KB

459 Nano RP2040 Connect, Grove sensors, and Putty is used to generate the dataset which is pub-
 460 lished on Kaggle [7] where more details are also provided regarding the testing environment.
 461 In addition, we utilize the WADI [8] dataset during the evaluation to analyze how the key ML
 462 parameters such as accuracy, F1-score, and prediction time differ according to the ML platform.
 463 The dataset contains data from 123 sensors and actuators from a water distribution testbed that
 464 had non-stop run for 16 days while being attacked during the last two days. Table 5 demon-
 465 strates the differences between two utilized datasets. By using these datasets, we also evaluate
 466 the relationship between the power consumption and data volume.

467 5.3. Anomaly Detection Methods

468 We tested the anomaly detection algorithms mentioned in Section 2 on WADI and AnoML
 469 datasets, as well as edge, fog, and cloud platforms by applying several scaling/reduction tech-
 470 niques. During the evaluations, we only used the baseline version of anomaly detection algo-
 471 rithms, hence we did not tune the algorithms to get better results to demonstrate the performance

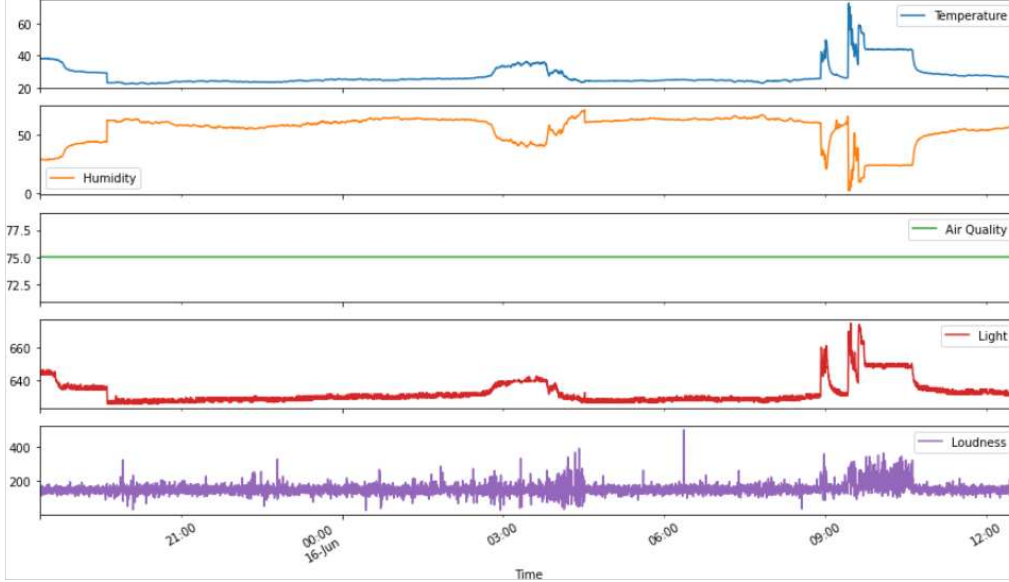


Figure 10: Demonstrates the behaviour per data type. The controlled anomalies for temperature and humidity data are identifiable. The anomalies in temperature and humidity are created via hair dryer. The air quality, light, and loudness can be considered as faulty due to not responding to our anomaly creation attempts. There might be two reasons for the occurrence of faulty data: (i) the sensors are cheap quality, (ii) these sensors are designed for the AVR[74] architecture. However, during the evaluation we used ARM-based microcontrollers.

[75].

472 of baseline versions. We evaluate the following parameters as they are the core elements that deter-
 473 mine the efficiency of an anomaly detection algorithm. *Accuracy* [76] determines the overall
 474 correct prediction ratio. *Precision* [77] defines how close the predictions are. *Recall* [77] demon-
 475 strates how the anomaly detection algorithm is successful at detecting normal data. *F1-Score*
 476 [76] is an evaluation metric that takes class distribution into considering. It might be preferred
 477 as the main metric when false detections matter (e.g., industrial environments). *True positive*
 478 (*TP*). The normal data is detected as normal. *True negative* (*TN*). The anomalous data is detected
 479 as anomalous. *False positive* (*FP*). The normal data is detected as anomalous. *False negative*
 480 (*FN*). The anomalous data is detected as normal. The following equations demonstrate how these
 481 parameters are calculated:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad Precision = \frac{TP}{TP + FP}$$

$$F1 - Score = \frac{2 * TP}{2 * TP + FP + FN} \quad Recall = \frac{TP}{TP + FN}$$

482 The TensorFlow has three different APIs. The main API (TensorFlow) contains all the avail-
 483 able methods and utilizes high-level API Keras [78] to build neural networks. TensorFlow Lite
 484 [43] which is the lightweight version of TensorFlow that is specifically designed for mobile/IoT
 485 devices is generated via the main TensorFlow API. TensorFlow Lite for Microcontrollers [79]
 486 only allows a subset of TensorFlow operations to be run on 32-bit microcontrollers. Hence, cur-

487 rently only CNN is supported. TensorFlow Lite models are converted into micro models via
488 TensorFlow Lite converter Python API [80]. Then, inference is possible at the edge.

489 Table 6 and Table 7 demonstrates the results of anomaly detection tests. We see that baseline
490 CNN at the edge achieves 77% accuracy and 86% F1-score. We did not observe any significant
491 differences regarding accuracy and F1-Score. We also see those reduction methods help to reduce
492 the inference time. If we compare scikit-learn methods and TensorFlow methods, we see that
493 scikit-learn performs better when there is more computing power, but TensorFlow might generate
494 better results on fog rather than the cloud. If we use CNN-AE, we see that the accuracy is higher
495 when the standard deviation is applied, but for one-class SVM, we see that skew performs better
496 than the standard deviation. Also, we see that the TF Lite models are significantly less in size
497 than TF models. During the evaluation of the Isolation Forest, the whole dataset must be fit at
498 once, due to nature algorithm. The Raspberry Pi 4B with 4 GB RAM fails due to high RAM
499 usage that occurs because of fitting the whole dataset at once. Hence, we utilized the 64-bit 8
500 GB version of Raspberry Pi to evaluate the performance of Isolation Forest.

501 6. Lessons Learned

502 IoT infrastructures are expanding thanks to the benefits of ubiquitous computing. The rapid
503 developments in lightweight edge mechanisms made them desirable for labor-intensive tasks
504 such as anomaly detection. The increasing variety of 32-bit microcontrollers allows us to choose
505 a specific device with desired features (e.g., relatively high RAM) that can run ML algorithms.
506 In the edge environment where these microcontrollers are utilized, while ML-based tasks such
507 as keyword spotting, natural language processing, image processing have proven to be promis-
508 ing, anomaly detection is still in a very early phase [81]. Machine learning pipelines are the
509 key elements that let us evaluate these edge ML algorithms. Thus, we develop an end-to-end
510 ML pipeline that facilitates developing anomaly detection systems where the detection can be
511 done on different layers (edge, fog, and cloud). We make the following observations based
512 on our findings: (i) lack of complete multi-protocol edge anomaly detection frameworks, (ii)
513 the lack of control over converted ML models, (iii) there is a lack of sensor data terminology,
514 (iv) manufacturers/developers provide platform-specific support, (vi) there is lack of support for
515 multi-language/protocol development frameworks, (vii) automated is actually not automated.

516 *Lack of complete multi-protocol edge anomaly detection frameworks.* The two tasks that
517 are executed within anomaly detection pipelines are: (i) data gathering/ingestion, (ii) anomaly
518 detection model development. While these tasks are different by nature, they are essential to
519 implement a complete IoT anomaly detection mechanism. During the evaluations, we see that
520 the required tools are very diverse and no ML framework offers a complete solution. For data
521 ingestion, flexibility and scalability are the most desired features. Hence, microcontrollers are
522 the primary elements that are deployed in IoT environments to gather data. Due to the resource
523 constraints of these environments, we need low power-consuming communication protocols.
524 This is one of the main reasons that why BLE [58] is invented. In the edge, during the evaluations
525 we see that there are three conditions needed for data acquiring: (i) sensors should be compatible
526 with the microcontroller, (ii) communication protocol libraries should be available to establish
527 a network both for the edge and the fog, and (iii) configurable platform that provides visual
528 assistance is needed. As no framework ensures all three conditions, we had to use tools that
529 require different programming languages (e.g., Arduino, JavaScript, Python).

530 *The lack of control over converted ML models.* TensorFlow allows us to convert the neural
531 network models that are generated via the main API to more lightweight models that can be

Table 6: Evaluation Results of WADI Dataset

Model Details			Inference Time (ms)		AUC		Accuracy		Recall		Precision		F1-Score		Scaling/Reduction Time (s)		Model Size (KB)	
Algorithm	SR	API	Fog	Cloud	Fog	Cloud	Fog	Cloud	Fog	Cloud	Fog	Cloud	Fog	Cloud	Fog	Cloud	Fog	Cloud
CNN-AE	Average	TF	1.01413	23.3021	0.418136	0.418136	0.323962	0.323962	0.311666	0.311666	0.91451	0.91451	0.464895	0.464895	6.88913	1.416631	34.64844	291.4189
	Kurtosis	TF	0.995147	24.0796	0.519019	0.519019	0.728089	0.728089	0.755389	0.755389	0.945001	0.945001	0.839623	0.839623	67.18984	23.71106	34.60938	291.0166
	MAD	TF	1.025978	24.03	0.448848	0.448848	0.7072	0.7072	0.740935	0.740935	0.934799	0.934799	0.826653	0.826653	72.07326	21.63906	34.64844	291.4189
	MM	TF	440.2933	96.6556	0.498538	0.498538	0.18035	0.18035	0.138802	0.138802	0.941108	0.941108	0.241923	0.241923	1.433161	0.481541	397203.7	1191798
	NS	TF	440	103.3094	0.5	0.50022	0.059669	0.059669	0.002144	0.002144	0.942253	0.953552	0.004278	0.004278	NA	NA	397194.2	1191799
	Skew	TF	1.001088	24.6464	0.361076	0.361073	0.261203	0.261197	0.248162	0.248156	0.885031	0.885028	0.387632	0.387624	69.34095	34.05004	34.60938	291.0166
	SS	TF	440	106.2	0.699918	0.699918	0.758349	0.758349	0.765979	0.765979	0.971539	0.971539	0.856599	0.856599	0.814415	0.790764	397203.6	1191796
	StDev	TF	1.051563	25.3799	0.497397	0.497397	0.929901	0.929901	0.986375	0.986375	0.941966	0.941966	0.963659	0.963659	20.04371	3.039728	34.64844	291.4189
RNN	Average	TF	8.164231	29.9278	0.431887	0.431887	0.174428	0.174428	0.140811	0.140811	0.892397	0.892397	0.243241	0.243241	6.911079	1.273132	797.8203	4045.537
	Kurtosis	TF	8.020618	30.07291	0.515931	0.515931	0.433345	0.433345	0.422561	0.422561	0.946374	0.946374	0.584251	0.584251	67.16076	22.99305	797.8047	4042.888
	MAD	TF	8.014902	29.02	0.421892	0.421892	0.69089	0.69089	0.726014	0.726014	0.930689	0.930689	0.815708	0.815708	71.88432	19.5469	797.8164	4044.854
	Skew	TF	7.621182	29.82979	0.448924	0.448965	0.544221	0.544209	0.556664	0.556646	0.932372	0.93238	0.697119	0.697107	69.21164	31.53431	797.8125	4044.202
	StDev	TF	8.09226	30.83805	0.420087	0.42009	0.436152	0.436158	0.43825	0.438256	0.922818	0.922819	0.594276	0.594282	20.04073	3.640271	797.8203	4045.537
RNN-AE	MM	TF	18.32667	40.9852	0.470241	0.470231	0.166829	0.166811	0.127211	0.127192	0.917464	0.917453	0.22344	0.223411	1.330966	0.328342	1414.535	10386.43
	NS	TF	17.90922	41.1471	0.5	0.500252	0.942253	0.05964	1	0.002107	0.942253	0.955432	0.970268	0.004205	NA	NA	1410.465	10386.67
	SS	TF	18.08283	49.55683	0.699284	0.699284	0.671158	0.671158	0.667486	0.667486	0.975904	0.975904	0.792754	0.792754	3.304777	1.005414	1414.512	10378.48
OC-SVM	Average	SK	0.894195	0.342938	0.491203	0.491203	0.7086	0.7086	0.801859	0.801859	0.84711	0.84711	0.823864	0.823864	6.88913	1.416631	56.62305	56.62305
	Kurtosis	SK	0.855004	0.341614	0.481406	0.481406	0.8127	0.8127	0.954818	0.954818	0.84496	0.84496	0.896536	0.896536	67.18984	23.71106	50.1582	50.1582
	MAD	SK	0.512389	0.207312	0.500747	0.500747	0.7141	0.7141	0.805624	0.805624	0.850137	0.850137	0.827282	0.827282	72.07326	21.63906	8.673828	8.673828
	MM	SK	0.687799	0.367881	0.5	0.5	0.1501	0.1501	0	0	1	1	0	0	1.433161	0.481541	407.5889	407.5889
	NS	SK	0.680684	0.335778	0.5	0.5	0.1501	0.1501	0	0	1	1	0	0	NA	NA	395.542	395.542
	SS	SK	1.314422	0.576975	0.496702	0.496702	0.8028	0.8028	0.93411	0.93411	0.849	0.849	0.889524	0.889524	0.814415	0.790764	1442.616	1442.616
	Skew	SK	0.545252	0.189985	0.493176	0.493176	0.8383	0.8383	0.986351	0.986351	0.848138	0.848138	0.912038	0.912038	69.34095	34.05004	10.2168	10.2168
	StDev	SK	1.13791	0.469978	0.501338	0.501338	0.7897	0.7897	0.913402	0.913402	0.850274	0.850274	0.880708	0.880708	20.04371	3.039728	92.62793	92.62793
IF	Average	SK	0.1884	0.2868	0.5891	0.5891	0.7152	0.7152	0.7179	0.7179	0.9921	0.9921	0.8330	0.8330	6.88913	1.416631	962.10	962.10
	Kurtosis	SK	0.1884	0.2868	0.6433	0.6433	0.7205	0.7205	0.7221	0.7221	0.9936	0.9936	0.8364	0.8364	67.18984	23.71106	775.20	775.20
	MAD	SK	0.1884	0.2868	0.4325	0.4325	0.7123	0.7123	0.7183	0.7183	0.9876	0.9876	0.8317	0.8317	72.07326	21.63906	634.68	634.68
	MM	SK	0.1978	0.3081	0.7155	0.7155	0.8399	0.8399	0.8426	0.8426	0.9948	0.9948	0.9124	0.9124	1.433161	0.481541	952.58	952.58
	NS	SK	0.1884	0.2868	0.6951	0.6951	0.8444	0.8220	0.8247	0.8247	0.9944	0.9944	0.9016	0.9016	133.64	133.64	NA	NA
	SS	SK	0.2119	0.3103	0.6951	0.6951	0.8444	0.8220	0.8247	0.8247	0.9944	0.9944	0.9016	0.9016	0.814415	0.790764	944.20	944.20
	Skew	SK	0.1884	0.2868	0.6415	0.6415	0.7036	0.7036	0.7049	0.7049	0.9937	0.9937	0.8247	0.8247	69.34095	34.05004	761.20	761.20
	StDev	SK	0.1884	0.2868	0.5519	0.5519	0.7034	0.7034	0.7066	0.7066	0.9910	0.9910	0.8250	0.8250	20.04371	3.039728	765.32	765.32

NA: Not applied. MM: MinMax scaler. NS: No scaler applied. SS: Standard Scaler. MAD: Median Absolute Deviation. StDev: Standard Deviation. IF: Isolation Forest. AE: Autoencoder. TF: TensorFlow. SK: scikit-learn.

Table 7: Evaluation Results of AnoML Dataset

Model Details		Inference Time (ms)			Accuracy			F1 Score		AUC			Recall		Precision		Scale Time (s)		Model Size (KB)						
Model	SR	API	Edge	Fog	Cloud	Edge	Fog	Cloud	Edge	Fog	Cloud	Edge	Fog	Cloud	Edge	Fog	Cloud	Edge	Fog	Cloud					
CNN	SR	TF	174.24	1.95	32.22	0.258	0.251	0.251	0.000	0	0	0.4991	0.485	0.485	0.000	0	0	0.000	0	0	0.716	0.12	19.62	3.172	80
	StDev	TF	172.05	1.13	32.63	0.769	0.871	0.871	0.862	0.91	0.91	0.5738	0.867	0.867	0.979	0.876	0.876	0.770	0.946	0.946	0.241	0.04	19.50	3.152	80
	Average	TF	NA	8E-15	32.20	NA	0.475	0.475	NA	0.57	0.57	NA	0.474	0.474	NA	0.476	0.476	NA	0.72	0.72	2.487	0.8	19.52	3.156	80
	Kurtosis	TF	NA	2E-07	32.04	NA	0.259	0.259	NA	0	0	NA	0.5	0.5	NA	0	0	NA	1	1	2.425	0.57	19.50	3.152	80
	MAD	TF	NA	1.81	32.16	NA	0.576	0.576	NA	0.60	0.60	NA	0.702	0.702	NA	0.441	0.441	NA	0.972	0.972	2.558	0.88	19.52	3.156	80

NA: Not applied. MAD: Median Absolute Deviation. StDev: Standard Deviation. IF: Isolation Forest. TF: TensorFlow.

Table 8: Data Format

Categories	Sensor Type					Protocols				Others		
Features	Temperature	Humidity	Air Quality	Light	Sound	WiFi	Bluetooth Classic	BLE	Zigbee	Location ID	Sensor ID	Microcontroller ID
Identifiers	TH	HU	AQ	LI	SO	WF	BC	BL	ZB	Numbers	Numbers	Numbers
Example Data	24.45	44.31	75	255	644	NA	NA	NA	NA	01	001	1

NA: Not Applicable.

532 deployed to edge and fog. However it offers zero control over the model conversion. During
 533 the evaluations, we realized that some lite models perform better than their main peers which
 534 were unexpected as there is a known trade-off between model accuracy and power consumption.
 535 Due to not having any control over this conversion, the only way to compare these models is
 536 by running inference. Also, scikit-learn models do not have a lite version. Pickle [82] library is
 537 used to pack and deploy the same model to another platform. Hence deploying these models to
 538 microcontrollers is currently not possible. This is why during the evaluations we did not observe
 539 any differences rather than inference time when running scikit-learn models.

540 *Lack of unified IoT sensor data terminology.* In an IoT environment, the sensors generate data
 541 in a similar format. Showing such similarities might be confusing if the data are not identified
 542 with certain terms/identifiers especially where big data are present. Hence, we apply further pro-
 543 cessing to identify the data context. However, there is no sensor terminology that explicitly states
 544 the sensor data format. Thus, we designed ECCG in a way that it generates a code that provides
 545 all required information such as the location of the sensor, sensor type, and the microcontroller
 546 model. Also, another important point to consider regarding IoT sensor data terminology is the
 547 data with floating points. While WiFi, Zigbee, and Bluetooth Classic have no issues when send-
 548 ing floating-point data, BLE by nature designed to send buffer (byte array) as it provides power
 549 optimization. Hence, when sending data over BLE we might require further indicators such as
 550 "F" for floats or "I" for integers to correctly identify the data. Table 8 presents the data for-
 551 mat that is generated via ECCG. We set the number IDs in different lengths to prevent possible
 552 confusion.

553 We evaluated our AnoML pipeline with two datasets, WADI and AnoML. WADI dataset has
 554 data from 127 different sensors, which is not possible to be replayed on edge devices due to
 555 limited computing power. WADI dataset was only evaluated on fog and cloud. We used a pri-
 556 vate machine with NVidia Tesla GPU and TensorFlow-GPU library version 2.4.1 for ML model
 557 training. AnoML dataset has data from five environmental sensors but two of them (temperature
 558 and humidity) were used to build and evaluate models on edge, fog, and cloud. We used Google
 559 Colab with GPU support to build models for the AnoML dataset with TensorFlow-GPU library
 560 version 2.1.1. AnoML dataset was transformed using reduction techniques only, to convert it
 561 into univariate.

562 In terms of accuracy, F1-score, precision, recall, and Area Under Curve (AUC), there was no
 563 significant difference recorded when comparing fog and cloud inferences in all types of models
 564 for the WADI Dataset. It was also observed that batch processing (all at once) in TensorFlow on
 565 the cloud does not affect the efficiency of the pipeline. Figure 11 shows the comparison of each

566 metric for all models for the WADI Dataset. Regarding fog and cloud, we observed a similar
567 trend for the AnoML Dataset. However, we noticed that edge results were comparatively less
568 efficient.

569 We noticed that inference on fog takes distinctly less time as compared to inference on the
570 cloud. But, we observed the opposite results for CNN-AE when doing multivariate analysis
571 where we recorded that cloud inference was faster than fog. Figure 11a provides a visual com-
572 parison between inference times of fog and cloud. We also experienced that when using batch
573 prediction in the cloud (all rows at once), the average prediction time was extremely faster. But,
574 this is not applicable in the real-world scenario where data is being streamed. Both scaling and
575 reduction techniques were faster on cloud (using batch-prediction) when compared to fog as seen
576 in Figure 11g. Inference time trends were also identical in ML models for the AnoML dataset
577 as fog inference time was less than cloud inference. The fog was more than 10 times faster than
578 cloud inference. Inference time on edge was extremely slow when compared to fog and cloud, it
579 was more than 100 times slower than the fog and almost 5 times slower than the cloud.

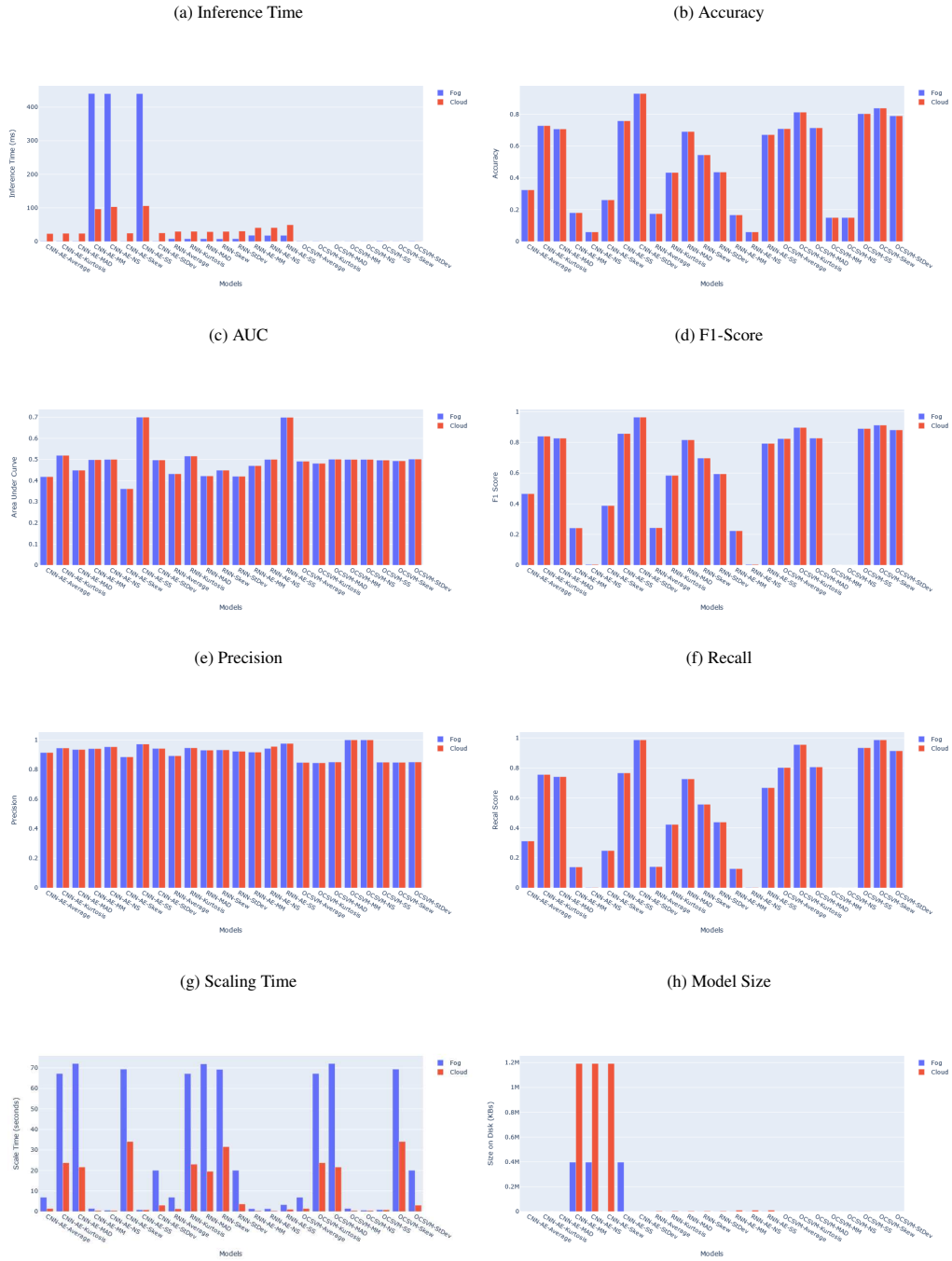
580 WADI Dataset is based on two sub-datasets recorded at different times, normal and attack,
581 as we discussed that normal dataset was used for ML model training thus we used the attack
582 dataset for testing. The attack dataset consists of 172800 rows which were then converted into
583 time-series data. As a result, the dimension of the data became (172770, 30, 127) for scale-based
584 and (172770, 30, 1) for reduction-based models. Reduction-based models took exceptionally
585 less time for training because of the univariate nature of the data. The size-on-disk of CNN-AE
586 scaled-based models is too big as compared to CNN-AE reduction-based models. For RNN-AE
587 models, the size-on-disk of reduction-based models are not significantly different from scale-
588 based RNN models. The main reason for this was that we used LSTM layers for RNN. The
589 size-on-disk variation trend was identical in fog models as seen in Figure 11h. We also learned
590 that scikit-learn models maintain consistency over fog and cloud. The time-related results in
591 the fog were always slower due to the difference in computing power, but there was no change
592 observed in accuracy-related metrics. The obvious reason was that there was no platform/format
593 conversion done for scikit-learn models.

594 Models for the AnoML dataset were lightweight but there was a notable difference for the
595 fog models in cloud model size-on-disk. We also observed that the size of edge models was
596 significantly greater than fog models, even though these models were converted from them. The
597 reason was that edge models were in plain-text (hex-dump) as they were a C-array, but fog
598 models were flat-buffered-binary format.

599 7. Conclusion

600 The proposed system offers support to a data scientist with minimal IoT knowledge to deploy
601 a reconfigurable IoT anomaly detection infrastructure that utilizes a data science pipeline. The
602 proposed framework contains: edge nodes consist of microcontrollers, fog nodes consist of single
603 board computers and virtual cloud nodes. The communication between nodes is not limited
604 to edge node types but limited with protocol specifications (e.g., Bluetooth only allows seven
605 slaves/servers), hence allowing the implementation of different network topologies (e.g., bus,
606 tree, and star). The system also explicitly supports the four major phases of data processing:
607 gathering, training, deployment, and inference. We evaluated several combinations to measure
608 the scalability that is offered by the proposed framework. We observed that the proposed anomaly
609 detection pipeline reduces the required labor for building an IoT anomaly detection system. We

Figure 11: WADI Dataset Evaluation



610 identified the drawbacks of the proposed system including the reasons behind them. We believe
611 our work encourages the future studies that aim to build open-source ML-based pipelines.

612 As future work, we are seeking to improve the ECGG web application by including more
613 edge nodes, sensors, and protocol types. In addition, we envision converting required manual
614 processes (e.g., deploying ML model to edge, uploading training data, deploying edge code)
615 within the AnoML-IoT to automated.

616 Acknowledgment

617 This work has been supported by GCHQ National Resilience Fellowship, Cardiff University
618 EPSRC Capital Award support for Early Career Researchers, and the PETRAS National Centre
619 of Excellence for IoT Systems Cybersecurity, which has been funded by the UK EPSRC under
620 grant number EP/S035362/1.

621 References

- 622 [1] H. Lasi, P. Fettke, H.-G. Kemper, T. Feld, M. Hoffmann, *Industry 4.0, Business & information systems engineering*
623 6 (2014) 239–242.
- 624 [2] M. Mohammadi, A. Al-Fuqaha, S. Sorour, M. Guizani, *Deep learning for IoT big data and streaming ana-*
625 *lytics: A survey, IEEE Commun. Surv. Tutorials* 20 (2018) 2923–2960. doi:10.1109/COMST.2018.2844341.
626 *arXiv:1712.04301*.
- 627 [3] Y. Jiang, C. Li, *Convolutional neural networks for image-based high-throughput plant phenotyping: a review, Plant*
628 *Phenomics* 2020 (2020).
- 629 [4] J. Goh, S. Adepu, M. Tan, Z. S. Lee, *Anomaly detection in cyber physical systems using recurrent neural networks,*
630 *in: 2017 IEEE 18th International Symposium on High Assurance Systems Engineering (HASE), IEEE, 2017, pp.*
631 *140–145.*
- 632 [5] F. T. Liu, K. M. Ting, Z.-H. Zhou, *Isolation forest, in: 2008 eighth IEEE international conference on data mining,*
633 *IEEE, 2008, pp. 413–422.*
- 634 [6] J. Ma, S. Perkins, *Time-series novelty detection using one-class support vector machines, in: Proceedings of the*
635 *International Joint Conference on Neural Networks, 2003., volume 3, IEEE, 2003, pp. 1741–1745.*
- 636 [7] H. Kayan, *AnoML-IoT, 2021. URL: <https://kaggle.com/hkayan/anomliot>.*
- 637 [8] C. M. Ahmed, V. R. Palleti, A. P. Mathur, *Wadi: a water distribution testbed for research in the design of secure*
638 *cyber physical systems, in: Proceedings of the 3rd International Workshop on Cyber-Physical Systems for Smart*
639 *Water Networks, 2017, pp. 25–28.*
- 640 [9] Y. Liu, Z. Pang, M. Karlsson, S. Gong, *Anomaly detection based on machine learning in IoT-based vertical plant*
641 *wall for indoor climate control, Build. Environ.* 183 (2020). doi:10.1016/j.buildenv.2020.107212.
- 642 [10] R. Chalapathy, S. Chawla, *Deep learning for anomaly detection: A survey, arXiv* (2019) 1–50.
643 *arXiv:1901.03407*.
- 644 [11] A. Blázquez-García, A. Conde, U. Mori, J. A. Lozano, *A review on outlier/anomaly detection in time series data,*
645 *arXiv* (2020). *arXiv:2002.04236*.
- 646 [12] J. Song, H. Takakura, Y. Okabe, K. Nakao, *Toward a more practical unsupervised anomaly detection system,*
647 *Information Sciences* 231 (2013) 4–14.
- 648 [13] Y. Zhu, N. M. Nayak, A. K. Roy-Chowdhury, *Context-aware activity recognition and anomaly detection in video,*
649 *IEEE Journal of Selected Topics in Signal Processing* 7 (2012) 91–101.
- 650 [14] L. Martí, N. Sanchez-Pi, J. M. Molina, A. C. B. Garcia, *Anomaly detection based on sensor data in petroleum*
651 *industry applications, Sensors* 15 (2015) 2774–2797.
- 652 [15] N. Jazdi, *Cyber physical systems in the context of industry 4.0, in: 2014 IEEE international conference on*
653 *automation, quality and testing, robotics, IEEE, 2014, pp. 1–4.*
- 654 [16] C. S. Raghavendra, K. M. Sivalingam, T. Znati, *Wireless sensor networks, Springer, 2006.*
- 655 [17] J. Pang, D. Liu, Y. Peng, X. Peng, *Anomaly detection based on uncertainty fusion for univariate monitoring series,*
656 *Measurement* 95 (2017) 280–292.
- 657 [18] Y. Su, Y. Zhao, C. Niu, R. Liu, W. Sun, D. Pei, *Robust anomaly detection for multivariate time series through*
658 *stochastic recurrent neural network, in: Proceedings of the 25th ACM SIGKDD International Conference on*
659 *Knowledge Discovery & Data Mining, 2019, pp. 2828–2837.*

- 660 [19] M. Teng, Anomaly detection on time series, Proc. 2010 IEEE Int. Conf. Prog. Informatics Comput. PIC 2010 1
661 (2010) 603–608. doi:10.1109/PIC.2010.5687485.
- 662 [20] H. S. Wu, A survey of research on anomaly detection for time series, 2016 13th Int. Comput. Conf. Wavelet Act.
663 Media Technol. Inf. Process. ICCWAMTIP 2017 (2017) 426–431. doi:10.1109/ICCWAMTIP.2016.8079887.
- 664 [21] V. Chandola, A. Banerjee, V. Kumar, Anomaly detection: A survey, ACM computing surveys (CSUR) 41 (2009)
665 1–58.
- 666 [22] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss,
667 V. Dubourg, et al., Scikit-learn: Machine learning in python, the Journal of machine Learning research 12 (2011)
668 2825–2830.
- 669 [23] M. Abadi, Tensorflow: learning functions at scale, in: Proceedings of the 21st ACM SIGPLAN International
670 Conference on Functional Programming, 2016, pp. 1–1.
- 671 [24] M. Munir, S. A. Siddiqui, A. Dengel, S. Ahmed, Deepant: A deep learning approach for unsupervised anomaly
672 detection in time series, IEEE Access 7 (2018) 1991–2005.
- 673 [25] T. Wen, R. Keyes, Time series anomaly detection using convolutional neural networks and transfer learning, arXiv
674 preprint arXiv:1905.13628 (2019).
- 675 [26] D. Li, D. Chen, B. Jin, L. Shi, J. Goh, S.-K. Ng, Mad-gan: Multivariate anomaly detection for time series data with
676 generative adversarial networks, in: International Conference on Artificial Neural Networks, Springer, 2019, pp.
677 703–716.
- 678 [27] TensorFlow, Introducing TensorFlow Decision Forests, 2021. URL: <https://blog.tensorflow.org/>.
- 679 [28] N. Ketkar, Introduction to pytorch, in: Deep learning with python, Springer, 2017, pp. 195–208.
- 680 [29] QuinnRadich, Automatic code generation with mlgen, 2021. URL: <https://docs.microsoft.com/>.
- 681 [30] MATLAB, Deep Learning Code Generation - MATLAB & Simulink - MathWorks United Kingdom, 2021. URL:
682 <https://uk.mathworks.com/help/deeplearning/deep-learning-code-generation.html>.
- 683 [31] H. Ren, B. Xu, Y. Wang, C. Yi, C. Huang, X. Kou, T. Xing, M. Yang, J. Tong, Q. Zhang, Time-series anomaly
684 detection service at microsoft, in: Proceedings of the 25th ACM SIGKDD International Conference on Knowledge
685 Discovery & Data Mining, 2019, pp. 3009–3017.
- 686 [32] NilsPohlmann, Create and run ML pipelines - Azure Machine Learning, 2021. URL: <https://docs.microsoft.com/en-us/azure/machine-learning/>.
- 687 [33] Microsoft, Cognitive Services – APIs for AI Developers | Microsoft Azure, 2021. URL: <https://azure.microsoft.com/en-gb/services/cognitive-services/>.
- 688 [34] Microsoft, Anomaly Detector - Anomaly Detection System | Microsoft Azure, 2021. URL: <https://azure.microsoft.com/en-us/services/cognitive-services/anomaly-detector/>.
- 689 [35] H. Zhao, Y. Wang, J. Duan, C. Huang, D. Cao, Y. Tong, B. Xu, J. Bai, J. Tong, Q. Zhang, Multivariate time-series
690 anomaly detection via graph attention network, arXiv preprint arXiv:2009.02040 (2020).
- 691 [36] AWS, Build your own Anomaly Detection ML Pipeline (2021) 1.
- 692 [37] E. Liberty, Z. Karnin, B. Xiang, L. Rousesnel, B. Coskun, R. Nallapati, J. Delgado, A. Sadoughi, Y. Astashonok,
693 P. Das, et al., Elastic machine learning algorithms in amazon sagemaker, in: Proceedings of the 2020 ACM
694 SIGMOD International Conference on Management of Data, 2020, pp. 731–737.
- 695 [38] S. Guha, N. Mishra, G. Roy, O. Schrijvers, Robust random cut forest based anomaly detection on streams, in:
696 International conference on machine learning, PMLR, 2016, pp. 2712–2721.
- 697 [39] M. D. Prado, J. Su, R. Saeed, L. Keller, N. Vallez, A. Anderson, D. Gregg, L. Benini, T. Llewellynn, N. Ouerhani,
698 et al., Bonseyes ai pipeline—bringing ai to you: End-to-end integration of data, algorithms, and deployment tools,
699 ACM Transactions on Internet of Things 1 (2020) 1–25.
- 700 [40] BAIR, Caffe | Deep Learning Framework, 2021. URL: <http://caffe.berkeleyvision.org/>.
- 701 [41] T. F. Foundation, The Open Source platform for our smart digital future, FIWARE. URL: <https://www.fiware.org/>.
- 702 [42] L. L. Fernández, M. P. Díaz, R. B. Mejías, F. J. López, J. A. Santos, Kurento: a media server technology for con-
703 vergent www/mobile real-time multimedia communications supporting webrtc, in: 2013 IEEE 14th International
704 Symposium on "A World of Wireless, Mobile and Multimedia Networks"(WoWMoM), IEEE, 2013, pp. 1–6.
- 705 [43] TensorFlow, TensorFlow Lite | ML for Mobile and Edge Devices, 2021. URL: <https://www.tensorflow.org/lite>.
- 706 [44] I. Drori, Y. Krishnamurthy, R. Rampin, R. Lourenço, J. One, K. Cho, C. Silva, J. Freire, Alphad3m: Machine
707 learning pipeline synthesis, in: AutoML Workshop at ICML, 2018.
- 708 [45] C. B. Browne, E. Powley, D. Whitehouse, S. M. Lucas, P. I. Cowling, P. Rohlfshagen, S. Tavener, D. Perez,
709 S. Samothrakakis, S. Colton, A survey of monte carlo tree search methods, IEEE Transactions on Computational
710 Intelligence and AI in games 4 (2012) 1–43.
- 711 [46] L. Bottou, Large-scale machine learning with stochastic gradient descent, in: Proceedings of COMPSTAT'2010,
712 Springer, 2010, pp. 177–186.
- 713 [47] J. R. Sutton, R. Mahajan, O. Akbilgic, R. Kamaleswaran, Physonline: an open source machine learning pipeline

- 719 for real-time analysis of streaming physiological waveform, *IEEE journal of biomedical and health informatics* 23
720 (2018) 59–65.
- 721 [48] X. Meng, J. Bradley, B. Yavuz, E. Sparks, S. Venkataraman, D. Liu, J. Freeman, D. Tsai, M. Amde, S. Owen, et al.,
722 Mllib: Machine learning in apache spark, *The Journal of Machine Learning Research* 17 (2016) 1235–1241.
- 723 [49] M. Nitsche, S. Halbritter, Development of an end-to-end deep learning pipeline, *Hochschule für Angewandte*
724 *Wissenschaften Hamburg* (2019).
- 725 [50] S. Shaikh, H. Vishwakarma, S. Mehta, K. R. Varshney, K. N. Ramamurthy, D. Wei, An end-to-end machine
726 learning pipeline that ensures fairness policies, *arXiv preprint arXiv:1710.06876* (2017).
- 727 [51] S. Boovaraghavan, A. Maravi, P. Mallela, Y. Agarwal, Mliot: An end-to-end machine learning system for the
728 internet-of-things, in: *Proceedings of the International Conference on Internet-of-Things Design and Implementa-*
729 *tion*, 2021, pp. 169–181.
- 730 [52] M. Molinara, M. Ferdinandi, G. Cerro, L. Ferrigno, E. Massera, An end to end indoor air monitoring system based
731 on machine learning and sensiplus platform, *IEEE Access* 8 (2020) 72204–72215.
- 732 [53] B. Vinzamuri, E. Khabiri, A. Bhamidipaty, G. Mckim, B. Gandhi, An end-to-end context aware anomaly detection
733 system, in: *2020 IEEE International Conference on Big Data (Big Data)*, IEEE, 2020, pp. 1689–1698.
- 734 [54] Y. Li, D. Zha, P. Venugopal, N. Zou, X. Hu, Pyodds: An end-to-end outlier detection system with automated
735 machine learning, in: *Companion Proceedings of the Web Conference 2020*, 2020, pp. 153–157.
- 736 [55] M. Fezari, A. Al Dahoud, Integrated development environment “ide” for arduino, *WSN applications* (2018) 1–12.
- 737 [56] G. R. Hiertz, D. Denteneer, L. Stibor, Y. Zang, X. P. Costa, B. Walke, The ieee 802.11 universe, *IEEE Communi-*
738 *cations Magazine* 48 (2010) 62–70.
- 739 [57] K.-H. Chang, Bluetooth: a viable solution for iot?[industry perspectives], *IEEE Wireless Communications* 21
740 (2014) 6–7.
- 741 [58] R. Heydon, N. Hunn, Bluetooth low energy, *CSR Presentation, Bluetooth SIG* <https://www.bluetooth.org/DocMan/handlers/DownloadDoc.ashx> (2012).
- 742 [59] S. C. Ergen, Zigbee/ieee 802.15. 4 summary, *UC Berkeley*, September 10 (2004) 11.
- 743 [60] J. Song, S. Han, A. Mok, D. Chen, M. Lucas, M. Nixon, W. Pratt, Wirelesshart: Applying wireless technology
744 in real-time industrial process control, in: *2008 IEEE Real-Time and Embedded Technology and Applications*
745 *Symposium*, IEEE, 2008, pp. 377–386.
- 746 [61] P. Warden, D. Situnayake, Tinyml: Machine learning with tensorflow lite on arduino and ultra-low-power micro-
747 controllers, ” O’Reilly Media, Inc.”, 2019.
- 748 [62] Node-RED, Node-RED, 2021. URL: <https://nodered.org/>.
- 749 [63] S. Tilkov, S. Vinoski, Node.js: Using javascript to build high-performance network programs, *IEEE Internet*
750 *Computing* 14 (2010) 80–83.
- 751 [64] Arduino, Arduino Nano 33 BLE Sense | Arduino Official Store, 2021. URL: <https://store.arduino.cc/arduino-nano-33-ble-sense>.
- 752 [65] N. Semiconductor, nRF52840 - Nordic Semiconductor, 2021. URL: nordicsemi.com/Products/nRF52840.
- 753 [66] IEEE, Ieee standard for floating-point arithmetic, *IEEE Std 754-2019 (Revision of IEEE 754-2008)* (2019) 1–84.
754 doi:10.1109/IEEESTD.2019.8766229.
- 755 [67] J. Yao, S. Warren, Applying the iso/ieee 11073 standards to wearable home health monitoring systems, *Journal of*
756 *clinical monitoring and computing* 19 (2005) 427–436.
- 757 [68] Y. A. Badamasi, The working principle of an arduino, in: *2014 11th international conference on electronics,*
758 *computer and computation (ICECCO)*, IEEE, 2014, pp. 1–4.
- 759 [69] A. Due, A. Core, Arduino due, Retrieved 9 (2017) 2019.
- 760 [70] Grove, Sensors - Seeed Studio Electronics, 2021. URL: <https://www.seeedstudio.com/category>.
- 761 [71] D. Quick, K.-K. R. Choo, Google drive: Forensic analysis of data remnants, *Journal of Network and Computer*
762 *Applications* 40 (2014) 179–193.
- 763 [72] T. Carneiro, R. V. M. Da Nóbrega, T. Nepomuceno, G.-B. Bian, V. H. C. De Albuquerque, P. P. Reboucas Filho,
764 Performance analysis of google colaboratory as a tool for accelerating deep learning applications, *IEEE Access* 6
765 (2018) 61677–61685.
- 766 [73] J. Giraldo, D. Urbina, A. Cardenas, J. Valente, M. Faisal, J. Ruths, N. O. Tippenhauer, H. Sandberg, R. Candell, A
767 survey of physics-based attack detection in cyber-physical systems, *ACM Computing Surveys (CSUR)* 51 (2018)
768 1–36.
- 769 [74] S. F. Barrett, D. J. Pack, Atmel avr microcontroller primer: Programming and interfacing, *Synthesis Lectures on*
770 *Digital Circuits and Systems* 7 (2012) 1–244.
- 771 [75] Y. Bai, Practical microcontroller engineering with ARM technology, John Wiley & Sons, 2015.
- 772 [76] M. Hasan, M. M. Islam, M. I. I. Zarif, M. Hashem, Attack and anomaly detection in iot sensors in iot sites using
773 machine learning approaches, *Internet of Things* 7 (2019) 100059.
- 774 [77] J. Davis, M. Goadrich, The relationship between precision-recall and roc curves, in: *Proceedings of the 23rd*
775 *international conference on Machine learning*, 2006, pp. 233–240.

- 778 [78] A. Gulli, S. Pal, Deep learning with Keras, Packt Publishing Ltd, 2017.
- 779 [79] Google, TensorFlow Lite Micro, 2021. URL: <https://www.tensorflow.org/lite/microcontrollers>.
- 780 [80] TensorFlow, TensorFlow Lite converter, 2021. URL: <https://www.tensorflow.org/lite/convert>.
- 781 [81] C. R. Banbury, V. J. Reddi, M. Lam, W. Fu, A. Fazel, J. Holleman, X. Huang, R. Hurtado, D. Kanter, A. Lokhmotov,
782 et al., Benchmarking tinymml systems: Challenges and direction, arXiv preprint arXiv:2003.04821 (2020).
- 783 [82] Python, pickle — Python object serialization — Python 3.9.6 documentation, pickle. URL: [https://docs.](https://docs.python.org/3/library/pickle.html)
784 [python.org/3/library/pickle.html](https://docs.python.org/3/library/pickle.html).