

Applying IoTDevID to a New Dataset: the CIC-IoT-2022 Case Study

Kahraman Kostas, Mike Just, and Michael A. Lones

Abstract—In this paper, we have examined under various headings the aspects to be considered in device identification studies using machine learning methods, common mistakes that may occur and how to avoid them. Our paper briefly touched upon the following topics: identification methods and their pros and cons, available data types and properties, common mistakes made during feature extraction and their solutions, what to consider about the use of machine learning methods, and how to choose appropriate evaluation methods.

Index Terms—IoT security, IoT fingerprinting, machine learning, device identification

dataset

One of the biggest problems in device identification studies is the lack of adequate datasets. The simulations used in many networking studies cannot be used in device identification studies, and the fact that the dataset can only be created with real device data is the most prominent of the difficulties. To create a proper device identification dataset, many types of IoT devices are needed, and the supply of these devices is a serious financial burden. In addition, normal data collection requires a long time, labour, and specialised space. In our baseline study, IoTDevID, we used two datasets, the Altoo University and UNSW datasets, which were produced for device identification studies. In this study, we used the Altoo dataset to develop our method and the second dataset to validate our results. In 2022, a new device identification dataset, CIC-IoT-22, was made public. This dataset is very interesting as it contains many types and numbers of devices, contains the state of the devices under different conditions, and contains attack data in addition to benign data. This dataset is very useful to demonstrate the usefulness, robustness, and generalisability of our method.

In this dataset, data was collected in 6 different situations. These situations can be summarised as follows. In the **Power** state, each device is isolated from other devices and rebooted and the network packets related to this device are collected. In the **Interactions** state, the device is interacted with by buttons, applications or voice commands and the network packets generated during this process are collected. In **Scenarios**, the network data of these devices are collected in situations such as entering the house, leaving the house, unauthorised entry to the house at night and day or user error. In case of **attack** state, data is collected by applying Flood attacks and RTSP Brute Force attacks to the devices. the **Idle** state consists of recording every 8-hour period for 30 days in the evening hours when the devices are working but not actively used. The **active**

state contains the data of the devices being used during the day for 30 days. This data is generated by people entering the lab and using the devices.

Some important points about the dataset: In this study, the most important sections for us are IDLE and Active. In these two sections, enough data has been collected from almost all devices. Although it is stated on the paper that 60 devices were used in this process, according to our own experiments and the information provided in the dataset, these sections contain 40 devices. These 40 devices consist only of lan WIFI devices, they do not include Zigbee and z-wave devices. Zigbee and Z-Wave devices have data isolated from other devices, including power interaction and hede hodo, but these data are both very limited and do not contain normal usage data.

FE

Python, Scapy and WireShark were used for feature extraction. Only individual package-based features are used for feature extraction. Many of these features are derived from packet headers, but there are also payload-based features such as payload entropy and payload bytes. Although the feature exruction system created about 100 features in total (features and their descriptions can be found in the table), very few of these features, only the sub-features selected during the feature selection phase of the IoTDevID study, were used in the experiment.

Labelling: Labelling was performed using the list of device names/MAC addresses couples in the dataset. In each fingerprint extracted, the source MAC address part was replaced with the given name and the MAC addresses not given in this list (5 MAC addresses that we believe belong to the hub, switch or the computer where the data is collected) were ignored.

Each of the pcap files we use for feature extraction contains network traffic recorded on a day, and is named with the date it was recorded. For example, data recorded on 24.11.2021 is labelled A211124 if Active and I211124 if Idle. In this context, 30 IDLE and 24 active sessions were recorded. as a preliminary study,, we aimed to test the performance of all these sessions by comparing them with each other. In order to compare the sessions with each other, they should contain similar devices. Unfortunately, data was not collected from every device in every session, and in some sessions some devices did not generate any data at all. Table 12 shows how much data was generated by each device in each session in terms of network packets. Therefore, we only compare sessions that contain the same devices with each other. For this comparison, we create a session ID. In this ID, each device is represented by a binary digit. If the session has that device, it is indicated with 1, if not, it is indicated with 0. For example, if sessin1 contains devices A, and C but not device B, then the

K. Kostas, M. Just, and M. A. Lones are with the Department of Computer Science, Heriot-Watt University, Edinburgh EH14 4AS, UK, e-mail: kk97, m.just, m.lones@hw.ac.uk

Kahraman Kostas supported by Republic of Turkey - Ministry of National Education

ID number is 101(ABC). Sessin1 can be compared to other sessions with the same ID number without any problem. In this context, we have created a 40-digit ID for each session according to totalling 40 devices.

The results we obtained by using devices with the same ID as training and test data are given in Fig. 1. We used the F1 score to present these results for roughly two reasons. Firstly, unlike accuracy, f1 score gives reliable results on unbalanced data sets. Secondly, the F1 score does not only give overall results, but also allows us to analyse the results by class. When the results are analysed in this context, it is seen that the f1 score varies between 40%-88% in pairwise session comparisons. Another point we would like to draw attention to here is that this process is a multiple classification process with approximately 40% classes. In this context, even 40 F1 is a much better result than chance/random success.

I. INTRODUCTION

An IoT device is any item that has a unique identity that can connect to other devices and perform control commands [1]. It is a bridge that connects our cyber world and our physical world [2]. Using IoT, we can manage our physical world from our cyber world.

This rapid progress has enabled many of us to include IoT devices in our lives. However, as a result of this rapid progress, many IoT devices were launched by many different manufacturers in a very short time. This rapid progress brings to mind the issue of whether devices are safe or how to ensure their safety.

It is unlikely to apply the solutions applied to classical computers to IoT devices. Because most of these devices have very limited possibilities in terms of battery, processor and storage. Also, due to the heterogeneous nature of IoT, there is no standard among these devices. Many use a unique operating system and hardware. In addition, IoT devices do not have the standard interfaces as conventional computers. This situation limits user-device interaction. In today's world, where there are a wide variety of devices, users can't take the security measures that every device needs. IoT device identification is a method that aims to find the device identity, such as brand and model, by analyzing the device behaviour.

Thanks to this method, devices on the network can be detected, and the necessary security measures can be taken for these devices. For example, the software of the needed ones can be updated, the behaviour of the devices and the addresses they will connect to can be restricted, or these devices can be isolated from the rest of the network.

In this study, we have divided the device identification process into four steps (see Figure 1), and included possible mistakes to be made in each step and how to deal with these mistakes.

Parallel to these steps, the structure of the paper is as follows. In Section II, device identification methods are discussed. Section III examines data types and considerations when choosing the appropriate dataset. Section IV highlights the feature extraction step and its key points. Section V discusses how to decide on the appropriate machine learning

method. In Section VI , evaluation methods and their pros and cons are explained.

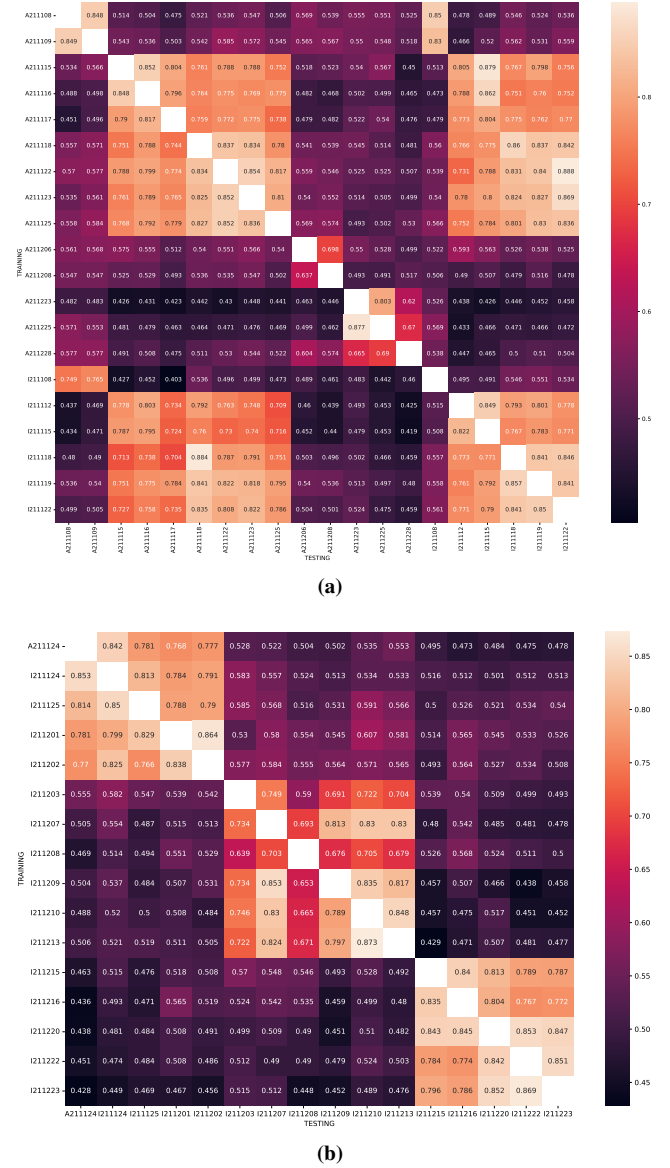


Fig. 1: F1 scores of sessions pairs assigned as training and test sets.

II. FAMILIARIZING WITH METHODS

In this section, different device identification methods are introduced and their advantages and disadvantages are discussed.

A. Define the limits of the scope of the IoT concept you will use in your study.

The Internet of Things (IoT) can be defined as all physical objects that can share data with other devices, have embedded sensors, software and various technologies [3]. However, this definition is too general. For your own work, you must decide what to consider as IoT and define the boundaries of your definitions. For example, if your data will contain, will you collect classic computers or mobile devices such as tablets and

phones or routers and switches under this definition or will you consider them as non-IoT devices? Will you accept Wireless Sensor Network (WSN) devices as IoT or will you create a separate subgroup for them? When classifying IoT devices, will you focus on a certain sub-area (home type, smart cities, agriculture, army, etc.).

B. Consider different methods of device identification.

Deciding what kind of definition you will make is very important as it will affect other steps of your work. You can identify according to devices with 3 different approaches [4]: **Unique Identification:** By accepting each of the devices as unique, a separate class is created for each device [5].

Type Identification: Identification is performed according to the device type. If there are multiple devices of the same make and model, they are seen as a single class [6].

Class Identification: Different devices that are not the same but have similar features, e.g., produced by the same brand for similar tasks, are gathered under a single class [7], [8].

Labels			Devices	Labels			Devices
Class	Type	Unique		Class	Type	Unique	
1	1	1	Aria	8	15	17	HomeMaticPlug
2	2	2	D-LinkCam	9	16	18	HueBridge
	3	3	D-LinkDayCam		17	19	HueSwitch
3	4	4	D-LinkDoorSensor	10	18	20	iKettle2
	5	5	D-LinkHomeHub		19	21	SmarterCoffee
	6	6	D-LinkSensor	11	20	22	Lightify
	7	7	D-LinkSiren	21	21	23	MAXGateway
	8	8	D-LinkSwitch	12	22	24	TP-LinkPlugHS100
	9	9	D-LinkWaterSensor		23	25	TP-LinkPlugHS110
	10	10	EdimaxCam 1	13	24	26	WeMoInsightSwitch 1
4	10	11	EdimaxCam 2		27	27	WeMoInsightSwitch 2
5	11	12	EdimaxPlug1101W		29	29	WeMoSwitch 1
	12	13	EdimaxPlug2101W		30	30	WeMoSwitch 2
6	13	14	EdnetCam 1	14	26	28	WeMoLink
	14	15	EdnetCam 2	15	27	31	Withings
7	14	16	EdnetGateway				

Fig. 2: Labelling the Aalto dataset according to 3 identification approaches.

Fig. 2 shows the labelling of the Aalto University IoT Devices Captures dataset [9] with three viewpoints. In the dataset containing 33 devices in total, 33 different labels are formed in the unique method, 27 different labels in the type method, and 15 different labels in the class method. Which of these three perspectives you choose is also effective in determining the network traffic characteristics you will use in your design.

C. Remember that using flow-based or packet-based features affects the generalizability of your models.

In your design, you can use packet-based, flow-based features or both. However, if you use Unique Device Identification, the features you will only get from the packets will not be enough. you will also need to use the flow-based features. On the other hand, using flow features will reduce the generalizability of your model, resulting in your model being specific only to the design you are using.

D. Analyze the strengths and weaknesses of these methods and choose according to your needs.

Type Identification and Class Identification methods allow you to obtain a model with higher generalizability by using only packet-based features. Using the Type Identification method does not promise great success in distinguishing similar devices, because similar devices show similar behavior patterns. Against this disadvantage, Class Identification can be used, which gathers similar devices under the same label. This approach is based on the assumption that similar devices can be grouped under the same group because they have similar hardware and software. However, in this method, expert involvement is required to decide on the similarities and groups of the devices. Before deciding on the method you will use in your design, it will be appropriate to examine the advantages and disadvantages of these methods, and to decide according to your purpose.

III. FAMILIARIZING WITH DATA

In this section, how the data you will use can be obtained, the characteristics and qualities that the data should have are examined.

A. Decide how you will obtain the data.

One of the important steps in device identification studies is how to obtain the data. In this context, you can use simulation or real devices. Although using simulation is fast, cheap, and easy, simulation programs currently lack the ability to simulate the various types of IoT devices because of the extremely heterogeneous nature of IoT devices. Therefore, the use of simulation is mostly seen in studies where many homogeneous devices such as WSN are used.

B. Do not forget about privacy if you are using real data.

In the use of real IoT devices, there are options to create an experiment set or usage of real data. When using test sets it is important that the devices can mimic as much as possible usage patterns. Although using real data may seem advantageous in many ways, the danger of causing privacy disclosures should be considered. If you are collecting data from a real environment, you should consider the ethical aspect, obtain the necessary relevant permissions and censor any information that may disclose confidentiality such as user identity or IP addresses.

C. Use data from previous studies.

Another data acquisition method is the use of other previously collected publicly available data. The use of such data is very advantageous as it does not require extra cost and provides the results to be comparable with the literature. You can scan the literature to find the appropriate data set, especially survey studies on device identification and fingerprinting can be very useful. You can also use websites specialized for datasets, e.g., Dataset Search, Kaggle, and UCI ML Repository

D. Unsure the data is suitable and has the necessary qualifications.

In addition, no matter how you obtain your data, the data you will use for device identification must have the following characteristics:

- It should contain benign data. In device identification, you need to extract the behavior patterns of the devices. For this, you need the normal/benign data of the devices. However, if the benign and malicious data packets are separated, you can also use the normal parts of the data sets prepared for anomaly detection (e.g., UNSW IoT benign and attack traces [10]). Another point that can be overlooked is that the tools used in vulnerability testing are very similar to attack tools. You can classify the data obtained by performing vulnerability tests on them as attack data, not normal data.
- The dataset should contain enough data. It is essential for a robust study that the dataset contains as many and varied devices as possible. In this way, the operation of your model is observed more soundly and the probability of success by chance is reduced.
- Devices must be labelled. In the dataset you will use, you should know which device the packets belong to (which device was produced by).

E. Ensure that the correct labels are assigned to the data.

It is quite common to use MAC and IP addresses to identify devices, but care should be taken when using these addresses. Especially if there are Non-IP devices in the datasets, this information can be misleading. The Aalto dataset is a good example for this. The non-IP HueSwitch device is connected to the gateway where data collection is made, via the HueBridge device. The HueSwitch transmits the network packets to the HueBridge device using ZigBee, which then HueBridge re-encapsulates these network packets and forwards them to the gateway via Ethernet. HueSwitch packets arriving at gateway carry the MAC addresses of the HueBridge device, not their own MAC addresses. Therefore, HueBridge and HueSwitch devices are represented by a single MAC address in the Aalto dataset. So, MAC and IP addresses cannot be used for labelling. That's why the creator of the dataset labelled the devices separately.

F. Remember that IoT devices are highly heterogeneous and this affects data distribution.

IoT devices are tools that contain a wide variety of software and hardware, produced for many different purposes. This is why data from devices can be quite diverse. So, IoT devices have an unbalanced data distribution due to their highly heterogeneous nature. One device in a network can generate a large amount of data, while another device can generate a very small amount of data. Fig. 3 shows the number of packets produced by the devices in the Aalto dataset.

G. Use data augmentation to deal with scarce data.

This unbalanced distribution of data can negatively affect some machine learning methods. You can deal with this

problem by using data augmentation techniques. The most important point during this process is to isolate the test and training datasets from each other **before** data augmentation. In this way, the augmented (post-generated) data in the training dataset is prevented from leaking into the test data.

H. Protect test data from all influence

Another common mistake is to apply augmentation of test data. It is best to keep the test data in its original state. Even the simple process of resampling, such as copying packets, is harmful because it will break the distribution of test data. The distribution of test dataset should also show the true world distribution. Thus, you can evaluate the outputs of your model more realistically and check whether there is a real contribution on data augmentation.

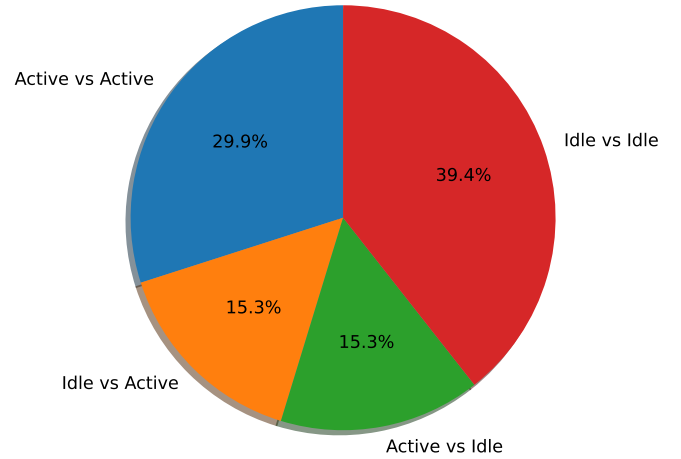


Fig. 3: The number of packets produced by the devices in the Aalto dataset.

IV. FAMILIARIZING WITH FEATURES EXTRACTION

In this section, the feature extraction process and important points to be considered during this process are highlighted.

A. Save time by choosing the right tools for preliminary analysis.

One disadvantage of working with network data is dealing with huge amounts of data. Analysing this data and extracting features can be very time consuming. Therefore, it is very important to choose the right tools for these operations. Recently, python-based libraries such as `dpkt` and `scapy` have been used extensively in the analysis of pcap data and feature extraction. Note that although they are easy to use and supported by a lot of documentation, they are extremely slow. Especially in very large data, using C-based programs such as `Wireshark` for your pioneering analysis will save you time. Even, with `tshark`, which you can automate with bash code, you can do many operations much faster.

B. Avoid features that uniquely identify your device.

Get to know the features you will use during the feature extraction process. Avoid including passive features such as MAC and IP addresses in your feature set. These properties are identifying features. They uniquely identify devices but do not provide any information about device behaviour. Even if you get high results using these features in your experiments, these results cannot be generalized because they suffer from overfitting.

C. Avoid features that uniquely identify a session.

Some session-based features are highly identifying, despite not static. However, these features uniquely identify the session, not the device. For example, various features such as port numbers, IP ID, TCP sequence numbers are determined at the beginning of the session, and are used until the session ends. Because these properties are randomly determined, they are not predictable or generalizable. You can avoid using these features in your study. If you are going to use it, you can do it in a healthy way. For example, isolate the test and train datasets from each other to ensure that they consist of different sessions. Thus, you can prevent these features, which are identifying for each session, from leaking from the training data to the test data.

D. Beware of features that are implicitly identifying.

Although time-related features such as timestamps do not appear to be stand-alone identifying, they can be identifiers as they will uniquely show a device's operating timespan. Header checksum features are another example of this situation. Checksum features are a summary of the header in which they are used. So IP checksum contains information about IP addresses, TCP checksum contains sequence, acknowledgment numbers, UDP checksum contains port numbers. If you find it inconvenient to use any attribute in the header, you should also avoid using the header's checksum feature.

E. Be careful when extracting features from raw data.

An alternative method in the feature extraction stage is the use of raw bytes. This method is based on obtaining feature sets through various size reduction methods by converting network data into raw bytes. A network packet converted to raw bytes is given in Fig. 4. This method is advantageous in terms of both obtaining original features and reducing expert involvement in the feature extraction stage. There is no harm in using this method in the packet payload, but if you include the packet headers, you may leak many features that you would not normally want to use in the feature set. So, when applying this method, it should be noted that raw data contains identifying features. Therefore, before performing feature extraction, censoring the raw data corresponding to the identifying features ensures that more robust features are extracted, and prevents errors that would lead to overfitting.

V. FAMILIARIZING WITH MACHINE LEARNING

In this section, it is focused on what should be considered when using machine learning methods and how the methods can be selected.

A. Consider other options when creating a multiclass model.

Device identification is a multi-class problem. However, there are some disadvantages of using multi-class models in solving this problem. The multi-class approach makes it difficult to extend the model. Because every time a new device is added to your system, you will have to retrain the model. In addition, scalability issues may arise as the number of devices increases. One2all approach can be preferred to deal with this problem. In this approach, a separate binary model is created for each device, and a multi-class result is obtained from the outputs of these models. In case a new device is included in this system, there is no need to recreate the whole models, only the model of this new device is added.

B. Do not underestimate the classical ML approaches

Do not underestimate the classical/older approaches when choosing the ML method you will use. Although deep learning methods have been popular recently, it has been observed that classical methods such as decision trees are much more successful than deep learning methods, especially in evaluating tabular data that lack temporal relationships, such as device identification [11].

C. Consider the interpretability of the models

Interpretability of algorithms can also be a matter of preference. Methods such as liner/logistic regression, DT, kNN have a high level of interpretability, while SVM, Ensemble methods and deep learning have a low level of interpretability (see Fig.5). High interpretability can be very useful when analysing and evaluating features [12].

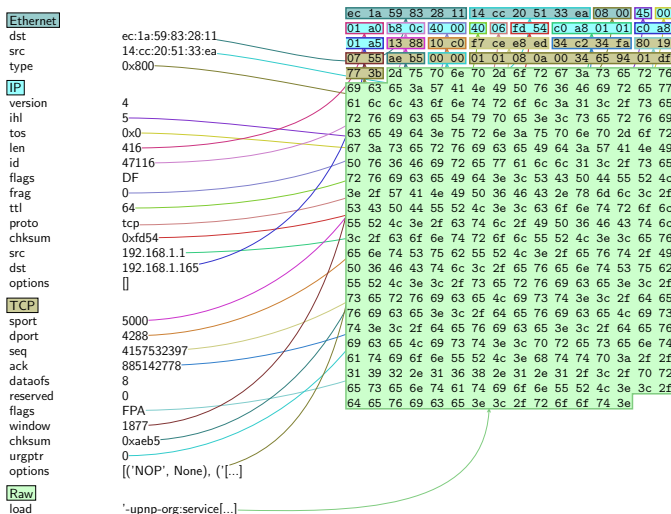


Fig. 4: The fields contained in a network packet and their byte equivalents..

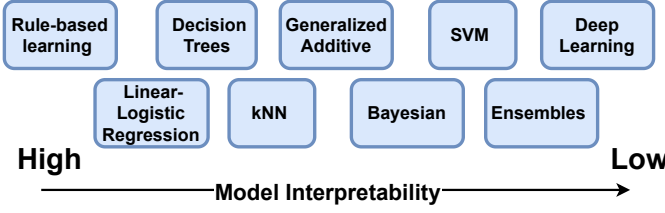


Fig. 5: Interpretability levels of some machine learning methods.

D. Note that there is no free lunch.

There is no perfect machine learning method that works best on every data type [13]. Try to find the appropriate machine learning method for your data. With so many methods available, finding the best method can be costly. You can try different strategies for this. For example, taking a look at the methods used by similar previous studies can be very helpful, especially survey studies. Another method would be to try one of each type of machine learning (for example, tree-based, kernel-based, ANN-based etc.) and experiment more deeply with the type that performs best.

E. Consider the success - inference time trade-off.

Do not accept detection success as the only criterion. In network technology, transactions are done at the micro-milliseconds level. The fast sorting capability of your model, which will be used in the real world and will work in a critical area such as security, is as important as the correct classification capability. Although some algorithms such as SVM and kNN achieve high performance levels, their inference time can be very high. Therefore, it is highly recommended to include the inference time criterion in your evaluation and to avoid algorithms with very high inference times.

VI. FAMILIARIZING WITH EVALUATION

In this section, performance evaluation methods are examined and their pros and cons are discussed.

A. Consider data distribution when choosing the evaluation method.

Various evaluation criteria are used to see how successful the study is. Accuracy is the most popular of these criteria. Since it is used for evaluation in many studies, using accuracy will add comparability to your study. However, in datasets that suffer from unbalanced distribution, such as IoT datasets, using accuracy as the only method can be quite misleading. Devices with very high or very low results with too many samples may give unrealistic results by pulling the result of all data up or down. In addition, the performance of devices with too few samples may be overlooked. Another disadvantage is that accuracy is a holistic method. So, it cannot give results per device/class. This is why many studies use the recall, for the presentation of device-based results.

B. Note that some evaluation concepts have many names.

Nevertheless, naming the recall value per class is also very problematic in the literature. Many names are used to express this criterion, such as identification rate, recognition rate, detection rate, accuracy rate, and individual device classification performance. In order not to add to this confusion, you can use a simpler and more general nomenclature such as overall recall or per-class recall.

C. Remember that some methods can be misleading.

However, using recall alone can be misleading as it does not account for False Positives. As a solution to this, you can use precision with recall to keep the balance between them. Another solution is to use the F1 score, which is the harmonic mean of the recall and precision. This metric alone shows the recall-precision balance, and you can observe both overall and class-based results with it.

VII. CONCLUSION

In this study, the methods used in device identification with machine learning and common mistakes are discussed. In this context, by addressing the positive and negative aspects of the methods used, identification methods, appropriate data types, the feature extraction process, selection and use of machine learning techniques, and evaluation methods are examined in order to provide practical solutions to common mistakes.

REFERENCES

- [1] A. Al-Fuqaha, M. Guizani, M. Mohammadi, M. Aledhari, and M. Ayyash, "Internet of things: A survey on enabling technologies, protocols, and applications," *IEEE communications surveys & tutorials*, vol. 17, no. 4, pp. 2347–2376, 2015.
- [2] C. Patel and N. Doshi, *Internet of Things Security: Challenges, Advances, and Analytics*. CRC Press, 2018.
- [3] F. Hussain, R. Hussain, S. A. Hassan, and E. Hossain, "Machine learning in iot security: Current solutions and future challenges," *IEEE Com. Surveys & Tutorials*, vol. 22, no. 3, pp. 1686–1721, 2020.
- [4] P. Yadav, A. Feraudo, B. Arief, S. F. Shahandashti, and V. G. Vassilakis, "Position paper: A systematic framework for categorising iot device fingerprinting mechanisms," in *Proceedings of the 2nd International Workshop on Challenges in AI and ML for IoT*, 2020, pp. 62–68.
- [5] S. A. Hamad, W. E. Zhang, Q. Z. Sheng, and S. Nepal, "IoT device identification via network-flow based fingerprinting and learning," in *18th IEEE TrustCom*. IEEE, 2019, pp. 103–111.
- [6] M. Miettinen, S. Marchal, I. Hafeez, N. Asokan, A.-R. Sadeghi, and S. Tarkoma, "IoT sentinel: Automated device-type identification for security enforcement in IoT," in *37th Int. Conf. DCS*. IEEE, 2017.
- [7] T. D. Nguyen, S. Marchal, M. Miettinen, H. Fereidooni, N. Asokan, and A.-R. Sadeghi, "Diot: A federated self-learning anomaly detection system for IoT," in *2019 IEEE 39th Int. Conf. on Distributed Computing Systems (ICDCS)*, 2019, pp. 756–767.
- [8] A. Aksoy and M. H. Gunes, "Automated IoT device identification using network traffic," in *ICC 2019-2019 IEEE International Conference on Communications (ICC)*. IEEE, 2019, pp. 1–7.
- [9] S. Marchal, "IoT devices captures, aalto university," 2017, accessed: 2021-08-25. [Online]. Available: <https://research.aalto.fi/en/datasets/iot-devices-captures>
- [10] A. Hamza, H. H. Gharakheili, T. A. Benson, and V. Sivaraman, "IoT benign and attack traces," 2019, accessed: 2021-09-25. [Online]. Available: <https://iotanalytics.unsw.edu.au/attack-data>
- [11] S. M. Lundberg, G. Erion, H. Chen, A. DeGrave, J. M. Prutkin, B. Nair, R. Katz, J. Himmelfarb, N. Bansal, and S.-I. Lee, "From local explanations to global understanding with explainable ai for trees," *Nature Machine Intelligence*, vol. 2, no. 1, pp. 2522–5839, 2020.

- [12] A. B. Arrieta, N. Díaz-Rodríguez, J. Del Ser, A. Bennetot, S. Tabik, A. Barbado, S. García, S. Gil-López, D. Molina, R. Benjamins *et al.*, “Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai,” *Information Fusion*, vol. 58, pp. 82–115, 2020.
- [13] D. H. Wolpert, “The supervised learning no-free-lunch theorems,” *Soft computing and industry*, pp. 25–42, 2002.

TABLE I: The list of individual packet-based features used in device identification and feature descriptions

No	Feature	Description
1	ts	Time Stamp
2	Ether_dst	Destination Media Access Control (MAC) Address
3	Ether_src	Source MAC Address
4	IP_src	Source Internet Protocol (IP) Address
5	IP_dst	Destination IP Address
6	WS_src	WireShark Source Address
7	WS_dst	WireShark Destination Address
8	pck_size	Packet (Frame) Size
9	Ether_type	Ethernet Type
10	LLC_dsap	Logical Link Control - Destination Service Access Point
11	LLC_ssap	Logical Link Control - Source Service Access Point
12	LLC_ctrl	Logical Link Control - Control
13	EAPOL_version	Extensible Authentication Protocol (EAPOL) version
14	EAPOL_type	Extensible Authentication Protocol (EAPOL) type
15	EAPOL_len	Extensible Authentication Protocol (EAPOL) Length
16	IP_version	IP version
17	IP_ihl	IP Internet Header Length
18	IP_tos	IP type of service
19	IP_len	IP Length
20	IP_flags	IP Flags
21	IP_Z	IP Zero
22	IP_MF	IP More Fragments
23	IP_id	IP identifier
24	IP_chksum	IP Checksum
25	IP_DF	IP Don't Fragment
26	IP_frag	IP fragmentation
27	IP_ttl	IP Time To Live
28	IP_proto	IP Protocols
29	IP_options	IP Options
30	ICMP_type	Internet Control Message Protocol (ICMP) Type
31	ICMP_code	ICMP Code
32	ICMP_chksum	ICMP Checksum
33	ICMP_id	ICMP identifier
34	ICMP_seq	ICMP Sequence Number
35	ICMP_ts_ori	ICMP ConditionalField
36	ICMP_ts_rx	ICMP ConditionalField
37	ICMP_ts_tx	ICMP ConditionalField
38	ICMP_ptr	ICMP ConditionalField
39	ICMP_reserved	ICMP ConditionalField
40	ICMP_length	ICMP length
41	ICMP_nexthopmtu	ICMP Next Hop Maximum Transmission Unit (MTU)
42	ICMP_unused	ICMP ConditionalField
43	TCP_seq	TCP Sequence Number
44	TCP_ack	TCP Acknowledgment Number
45	TCP_dataofs	TCP data offset
46	TCP_reserved	TCP Reserved
47	TCP_flags	TCP Flags
48	TCP_FIN	FINished Flag
49	TCP_SYN	Sync Flag
50	TCP_RST	Reset Flag

TABLE I: The list of individual packet-based features used in device identification and feature descriptions

No	Feature	Description
51	TCP_PSH	Push Flag
52	TCP_ACK	Acknowledgment Flag
53	TCP_URG	Urgent Flag
54	TCP_ECE	ECE Flag
55	TCP_CWR	CWR Flag
56	TCP_window	TCP Window Size
57	TCP_chksm	TCP Checksum
58	TCP_urgptr	TCP Urgent Pointer
59	TCP_options	TCP Options
60	UDP_len	User datagram protocol (UDP) Length
61	UDP_chksm	UDP Checksum
62	DHCP_options	Dynamic Host Configuration Protocol (DHCP) Options
63	BOOTP_op	Bootstrap Protocol (BOOTP) Options
64	BOOTP_htype	BOOTP Hardware Len
65	BOOTP_hlen	BOOTP Hardware Length
66	BOOTP_hops	BOOTP Hardware Options
67	BOOTP_xid	BOOTP Transaction Identifier
68	BOOTP_secs	BOOTP Seconds
69	BOOTP_flags	BOOTP Flags
70	BOOTP_sname	BOOTP Server Name
71	BOOTP_file	BOOTP Boot Filename
72	BOOTP_options	BOOTP Options
73	DNS_length	Domain Name System (DNS) Length
74	DNS_id	DNS Identifier
75	DNS_qr	DNS Query-Response
76	DNS_opcode	DNS Operation Code
77	DNS_aa	DNS Authoritative Answer
78	DNS_tc	DNS TrunCation
79	DNS_rd	DNS Recursion Desired
80	DNS_ra	DNS Recursion Available
81	DNS_z	DNS Reserved for future use
82	DNS_ad	DNS Authentic Data
83	DNS_cd	DNS Checking Disabled
84	DNS_rcode	DNS Response Code
85	DNS_qdcount	DNS The unsigned fields query count
86	DNS_ancount	DNS Answer Count
87	DNS_nscount	DNS Authority Count
88	DNS_arcount	DNS Additional Information Count
89	sport_class	Source Port Class (IoTDevID classing)
90	dport_class	Destination Port Class (IoTDevID classing)
91	sport23	Source Port Class (keep wellknown ports between 0-1023)
92	dport23	Destination Port Class (keep wellknown ports between 0-1023)
93	sport_bare	Source Port Number
94	dport_bare	Destination Port Number
95	payload_bytes	Payload size in Bytes
96	entropy	Payload Entropy
97	Protocol	WireShark Protocol
98	sport	Source Port Number
99	dport	Destination Port Number
100	Label	Packet Level Label

TABLE II: Sessions, the total number of packets generated by the devices and devices in the session.

ML	Date	Type	AMCREST WiFi Camera	Amazon Alexa Echo Dot 1	Amazon Alexa Echo Dot 2	Amazon Alexa Echo Spot	Amazon Alexa Echo Studio	Arlo Base Station	Arlo Q Camera	Atom Coffee Maker	Bosch/Stihlman-AI Camera	D-Link Water Sensor	D-Link Mini Camera	Fury Homebase 2	Globe Lamp ESP_B1680C	Google Nest Mini	Gosund ESP_032979 Plug	Gosund ESP_039A9F Socket	Gosund ESP_0C3994 Plug	Gosund ESP_10098F Socket	Gosund ESP_10ACD8 Plug	Gosund ESP_14TF99 Plug	Gosund ESP_1ACEE1 Socket	HeimVision Smart WiFi Camera	HeimVision SmartLife R-Lamp	Home Eye Camera	LG Smart TV	Loche Cam Dog	Next Indoor Camera	Netatmo Camera	Netatmo Weather Station	Philips Hue Bridge	Ring Base Station AC1236	SIMCAM 1S (AMPAKTe)	Smart Board	Sonos One Speaker	Teckin Plug 1	Teckin Plug 2	Yurton Plug 1	Yurton Plug 2	Yurton Roomba	
Train	20211002	ACTIVE	554	691	134	154	130	17	21.88	693	24.42	0.00	1.11	8.07	6.75	21.27	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9		
Train	20211105	ACTIVE	552	710	698	1519	1312	1.48	495	21.54	69.3	26.20	0.00	1.11	8.07	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9		
Train	20211108	ACTIVE	563	1073	1092	1499	1306	1.47	450.88	78.96	69.4	26.05	0.00	1.11	8.07	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9		
Train	20211109	ACTIVE	609	1144	1094	1472	1317	1.47	459.97	47.26	69.3	24.37	0.00	1.09	7.97	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9	
Train	20211110	ACTIVE	674	1407	1094	1498	1495	2.39	22.25	37.38	69.2	24.09	0.00	1.09	8.06	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9	
Train	20211111	ACTIVE	566	1214	1183	1493	1421	1.82	11.10	55.18	69.4	17.26	0.00	1.17	8.15	6.88	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9
Train	20211115	ACTIVE	566	1214	1183	1493	1421	1.82	11.10	55.18	69.4	17.26	0.00	1.17	8.15	6.88	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9
Train	20211116	ACTIVE	566	1214	1183	1493	1421	1.82	11.10	55.18	69.4	17.26	0.00	1.17	8.15	6.88	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9
Train	20211117	ACTIVE	566	1214	1183	1493	1421	1.82	11.10	55.18	69.4	17.26	0.00	1.17	8.15	6.88	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9
Train	20211118	ACTIVE	597	1135	1333	1421	1471	4.81	49.7	69.4	24.18	0.00	1.11	8.07	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9
Test	20211119	ACTIVE	554	1183	1115	1341	1422	1.47	49.1	36.4	69.2	24.12	0.00	1.11	8.07	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9	
Test	20211122	ACTIVE	562	954	935	1491	1349	1.47	15.10	58.72	69.4	24.12	0.00	1.12	8.15	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9	
Test	20211124	ACTIVE	628	1153	1040	1534	1358	1.41	7.82	58.50	71.10	24.76	0.00	1.15	8.12	6.89	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9
Test	20211125	ACTIVE	574	1122	1028	1461	1334	1.47	4.90	17.54	69.3	24.88	0.00	1.10	7.96	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9	
Test	20211125	ACTIVE	576	1137	1025	1600	1408	1.55	5.15	34.63	69.4	24.23	0.00	1.09	8.07	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9	
Test	20211126	ACTIVE	576	1137	1025	1600	1408	1.55	5.15	34.63	69.4	24.23	0.00	1.09	8.07	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9	
Test	20211127	ACTIVE	576	1137	1025	1600	1408	1.55	5.15	34.63	69.4	24.23	0.00	1.09	8.07	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9	
Test	20211207	ACTIVE	560	1106	1068	1559	1576	1.52	44.38	41.61	69.3	27.73	0.00	1.09	8.29	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9	
Test	20211208	ACTIVE	566	1097	1055	1557	1444	1.52	43.8	41.61	69.3	27.73	0.00	1.09	8.29	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9	
Test	20211208	ACTIVE	566	1097	1055	1557	1444	1.52	43.8	41.61	69.3	27.73	0.00	1.09	8.29	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9	
Test	20211223	ACTIVE	616	1066	1059	1534	1363	1.52	13.70	29.37	69.2	24.79	0.00	1.08	7.99	6.72	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9
Train	20211223	ACTIVE	626	1173	1149	1533	1333	1.46	4.93	3.79	69.2	24.79	0.00	1.08	7.99	6.72	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9
Train	20211225	ACTIVE	548	1163	1176	1488	1439	1.46	13.16	23.09	69.3	25.21	0.00	1.08	8.01	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9	
Train	20211228	ACTIVE	548	1156	1193	2039	1507	1.46	8.51	23.50	69.2	24.60	0.00	1.07	8.45	6.72	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9
Train	20211228	ACTIVE	548	1156	1193	2039	1507	1.46	8.51	23.50	69.2	24.60	0.00	1.07	8.45	6.72	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9
Train	20211228	ACTIVE	548	1156	1193	2039	1507	1.46	8.51	23.50	69.2	24.60	0.00	1.07	8.45	6.72	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9
Train	20211228	ACTIVE	548	1156	1193	2039	1507	1.46	8.51	23.50	69.2	24.60	0.00	1.07	8.45	6.72	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9
Train	20211228	ACTIVE	548	1156	1193	2039	1507	1.46	8.51	23.50	69.2	24.60	0.00	1.07	8.45	6.72	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9
Train	20211228	ACTIVE	548	1156	1193	2039	1507	1.46	8.51	23.50	69.2	24.60	0.00	1.07	8.45	6.72	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9
Train	20211228	ACTIVE	548	1156	1193	2039	1507	1.46	8.51	23.50	69.2	24.60	0.00	1.07	8.45	6.72	6.73	6.74	6.75	6.74	6.75	6.74	6.75	6.74	6.75	6.74	14.70	6.74	25.28	690	2.92	24.23	26.15	1.70	12.67	3.26	21.08	0.55	6.62	6.29	6.53	89.9
Train	20211228	ACTIVE	548	1156	1193	2039	1507	1.46	8.51	23.50	69.2	24.60	0.00																													