



Identifying Internet of Things software activities using deep learning-based electromagnetic side-channel analysis

Quan Le ^a, Luis Miralles-Pechuán ^b, Asanka Sayakkara ^{c,d}, Nhien-An Le-Khac ^d, Mark Scanlon ^{a,d,*}

^a Centre for Applied Data Analytics Research (CeADAR), University College Dublin, Dublin, Ireland

^b School of Computer Science, Technological University Dublin, Dublin, Ireland

^c School of Computing (UCSC), University of Colombo, Sri Lanka

^d Forensics and Security Research Group, University College Dublin, Ireland

ARTICLE INFO

Article history:

Received 30 May 2020

Received in revised form

10 April 2021

Accepted 26 September 2021

Available online xxx

Keywords:

Electromagnetic side-channel analysis

Internet of Things

Deep learning

Random forests

ABSTRACT

Internet of Things (IoT) is becoming the new frontier in digital forensics due to the abundance of IoT devices appearing in day-to-day life. The diversity and complexity of IoT ecosystems pose a considerable challenge to digital investigators that demand novel approaches. Electromagnetic side-channel analysis (EM-SCA) has been proposed as a promising window to gather forensically useful information from IoT devices. Machine Learning (ML) techniques are instrumental when performing EM-SCA on IoT devices. Our work aims to investigate how machine learning can be applied to accurately identify complex activities on IoT devices from their generated electromagnetic noises. To this end, a range of classification models were created, including deep learning models, to predict the activity from the electromagnetic noise emitted while the device performed the activities. A dataset was generated by using ten different well-known sorting algorithms with diverse computational time complexities and running them on an Arduino Leonardo device to represent a low-powered IoT device. The algorithms were continually sorting arrays of 100 elements randomly generated in ascending order. Experiments were conducted to identify which ML methods performed better with the generated data sets. Furthermore, more experiments were conducted to identify how the methods perform depending on the window size of raw samples and the number of examples against which they are trained. From the experimental results, it is possible to predict which activity is being executed with a high level of accuracy (99.6%) with a convolutional neural network (CNN). It was also found that Random Forests (RF) and Deep Learning (DL) are suitable ML models for making predictions with EM-SCA.

© 2021 The Author(s). Published by Elsevier Ltd. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

1. Introduction

Since the emergence of computers, Moore's law enabled engineers to increase their capabilities with time. Besides making computers more powerful and resource-rich, this trend has also enabled the production of smaller and energy-efficient computing devices such as mobile devices, that have less computing power but can last longer on battery power. The ultimate result of those advances is the emergence of the Internet of Things (IoT); small

computing devices with sensors and networking capability that has started to embed into our everyday lives including sports, health-care, and smart-home systems (Lin et al., 2017).

The integration of IoT devices into everyday life makes it inevitable for law enforcement to encounter them in digital forensic investigations (Du et al., 2020). In modern society, most people carry a smartphone in their pocket and many can have some kind of wearable device such as a smart wristwatch or a personal health tracking device. Therefore, an increasing number of IoT devices are likely to contain important pieces of information related to digital investigation. As a result, IoT forensics has gained significant interest among the law-enforcement community and has become recently a hot topic in the research community (Watson and Dehghantanha, 2016; Yaqoob et al., 2019).

* Corresponding author. Centre for Applied Data Analytics Research (CeADAR), University College Dublin, Dublin, Ireland.

E-mail addresses: quan.le@ucd.ie (Q. Le), luis.miralles@tudublin.ie (L. Miralles-Pechuán), asa@ucsc.cmb.ac.lk (A. Sayakkara), an.lekhac@ucd.ie (N.-A. Le-Khac), mark.scanlon@ucd.ie (M. Scanlon).

The acquisition of digital evidence from IoT devices is a challenging task compared to traditional digital forensics for several reasons. For example, traditional digital forensics involves devices that have easily accessible non-volatile data storage such as hard disks, solid-state drives (SSD), or removable media, e.g., secure digital (SD) cards. When required, a forensic image of the storage can be taken for analysis. In contrast, most IoT devices consist of on-board flash storage that is difficult to access for forensic image acquisition (Sayakkara et al., 2019a). This is due to the manufacturer-specific custom hardware architectures used for IoT devices. This problem already exists for smartphone forensics. However, the high diversity of IoT device hardware architectures has made it significantly more complex to perform IoT forensics. As a result, it is highly necessary to find alternative approaches to acquire forensic insights from IoT devices (Meffert et al., 2017; Sayakkara et al., 2019a).

Electromagnetic side-channel analysis (EM-SCA) is the discipline of gathering information about the internal operations of digital electronics by analysing the electromagnetic noise emissions they produce. These emissions are a by-product of the high-frequency switching operations of semiconductor electronics (Zankl et al., 2018). EM-SCA has been recently proposed as a method to acquire forensically useful insights from IoT devices (Sayakkara et al., 2018, 2020a). When an IoT device that is powered on is seized by law-enforcement, EM-SCA procedures can potentially be performed immediately on the device. It has been proposed that various information can be detected with EM-SCA-based approaches such as tampered firmware, software behaviour, and cryptographic algorithms (Sayakkara et al., 2019a). Due to the non-invasive nature of acquiring and analysing EM emission signals, there is no need to physically tamper with the device being investigated.

The emissions that are produced by IoT devices are weak signals that are not detectable beyond a few centimetres from the device unless powerful signal amplifiers are used. However, unlike unauthorised eavesdropping of devices through EM-SCA, a forensic inspection of a device does not have a requirement to make observations from a distance. The side-channel EM emissions from IoT devices are wide-band signals that are highly dimensional and can be captured with fast sampling rates. Therefore, when captured and stored for analysis, EM data files, called *EM traces*, usually take up multiple gigabytes. When analysing such EM traces for identifying known patterns, applications of artificial intelligence such as Machine Learning (ML) have been explored. ML algorithms are capable of identifying patterns in large and high dimensional data that are not easily recognisable by humans (Sayakkara et al., 2019b).

The application of ML techniques for IoT software behaviour identification still faces many challenges that need to be addressed. From a classification point of view, an IoT device's firmware can be in a wide variety of internal states that can result in a large number of potential classes. Even when considering the same known software activity, slight changes in the variables used within the software activity can affect its emission patterns. Therefore, it is necessary to further evaluate ML-based methodologies for classifying complex software behaviours under varying conditions.

1.1. Contribution of this work

Our work aims to investigate how machine learning can be applied to accurately identify complex activities on IoT devices from its generated EM noise and assess the influence of EM sample parameters. Towards this goal, this work explored the performance of ML-based models for classifying complex software programs running on a representative IoT device with randomised internal variables. In the evaluation, a dataset was generated using

implementations of multiple well-known sorting algorithms running on the target device that use randomly generated input data arrays. The contribution of this work can be summarised as follows:

- Demonstration that machine learning models trained on past data of complex software activities recognise these activities in future recordings.
- Experimentation showing that when using EM signals for ML-based classification purposes, the use of data in the frequency domain is more effective than in the time domain and other feature extraction methods.
- Among various ML methods evaluated, deep learning approaches outperform all the other algorithms when used for classifying complex software activities using EM side-channel data.

2. Background

This section covers the background information related to EM-SCA and the application of ML in its context. Subsection 2.1 provides a comprehensive overview of EM-SCA including methods of data acquisition and tools available for this purpose. Subsection 2.2 provides an overview of ML techniques. Subsections 2.2.1 and 2.2.2 provide background information on two ML architectures, namely; *Random Forest* and *Deep Learning* methods. Finally, Subsection 2.3 illustrates the current status of applying EM-SCA for digital forensics.

2.1. Electromagnetic side-channel analysis

The earliest work on electromagnetic information leakage was published in 1985. Van Eck (1985) demonstrated that cathode-ray tube (CRT) video displays emitted sufficient EM emissions that could be used remotely to reconstruct the content shown on the display. Later, research on eavesdropping computing devices through their unintentional EM emissions started gaining interest among researchers with various objectives such as spying on computer monitor content (Elibol et al., 2012), retrieving cryptographic keys (Kocher et al., 1999), and identifying abnormal deviations of software execution due to malware (Nazari et al., 2017). IoT devices pose more opportunities for EM-SCA due to the abundance of them within the reach of an attacker in contrast to desktop and laptop computers that tend to stay under the possession of their respective owners.

When machine code instructions are being executed on the processor of an IoT device, sequences of binary digits are moved around the processor's registers rapidly. This storing and resetting operations of values in registers impact the draw of electrical current by the processor. The higher the hamming weight of the value being stored, the higher the impact it has on the current draw of the processor. This varying electrical current draw results in varying EM noise to be generated from the circuitry of the processor (Maxwell, 1865). Such weak EM noise can get modulated into any powerful radio frequency signal that is already present in the vicinity. There are two potential radio frequency signal sources on IoT devices; the processor clock and on-board radio transceivers in the case of system-on-chip (SoC) devices (Camurati et al., 2018).

Observation of EM emissions from IoT devices requires specialised equipment. In the traditional form, oscilloscopes with compatible EM probes can be used to capture the signal (Peeters et al., 2007). The identification of the exact frequency of the leakage signal requires a certain amount of trial-and-error procedure. In the worst case, the attacker has to sweep through the EM spectrum looking for suspicious signals. However, the most

common practice is to use the clock frequency of the device being observed as the frequency of the EM emission signal. When capturing EM signals, several important factors need to be met at the hardware level to ensure successful data acquisition. Firstly, the frequency range of the observation device has to cover the clock frequency of the device being attacked. Secondly, the sampling rate of the device has to be fast enough. The sampling rate is the number of samples taken at the tuned frequency per second. Finally, the equipment has to support a wide enough bandwidth in order to capture as many signals around the leakage signal as possible.

As a much cheaper and more flexible alternative to oscilloscopes, software-defined radios (SDRs) devices can be used. These devices are popular among wireless hackers, hobbyists, and security enthusiasts (Tuttlebee, 2003). An SDR device makes use of a bare minimum hardware layer that can be tuned to a specific frequency and digitally sample EM signals with the help of a fast analogue-to-digital converter (ADC). Such digital samples are later processed with software using various software-defined radio packages – among which GNURadio is the most popular open-source package (Cass, 2013; Ossmann, 2016; Ettus and Braun, 2015; Blossom, 2004). Each sample taken from an SDR device represents the amplitude of the tuned frequency at a specific instance of time. In this work, an SDR is used to sample EM emissions from a target IoT device due to the flexibility it offers.

2.2. Relevant machine learning methods

In the last decades, thanks to the increase of computational power and to the development of artificial intelligence, researchers can easily implement ML methods and run them on their computers to conduct experiments (Marsland, 2015). Arthur Samuel, one of the pioneers of ML, defined it as the “field of study that gives computers the ability to learn without being explicitly programmed” (Samuel, 1959). Supervised Learning is the more developed ML learning branch (Marsland, 2015) and it is the branch used in this work – specifically Random Forests and Deep Learning.

2.2.1. Random forests

Random Forests (RF) select those attributes that are able to distinguish elements that belong to different groups and at the same time, divided groups that belong to the same class. RF try to create different trees by choosing different attributes from the database. This supervised method was introduced in 2001 by Breiman (2001), and has become one of the most important supervised methods in ML. Its main advantages are its great predictive power and its rapid execution, and for these reasons, it is an ideal algorithm for the selection of variables. The RF model can be used for both regression and classification, but they are mainly used for classification. RF combines several binary trees that it generates from training data by randomly selecting a subset of variables (Genuer et al., 2010). The two main RF parameters are: *mtry* (number of randomly selected variables in each split) and *ntree* (number of trees that make up the forest). RF are very suitable for feature selection because they have both high levels of accuracy and also a very quick execution time (Genuer et al., 2010) and easy to parallelize (Santner et al., 2009).

2.2.2. Deep learning methods

Deep Learning (DL) is a group of algorithms based on Artificial Neural Networks (ANN) (LeCun et al., 2015). DL models are composed of multiple layers - each includes a number of nodes. Multiple DL layer architectures are developed to capture different correlation patterns in the input; the popular ones are the convolutional layer, the long short term memory layer, and the attention layer (Goodfellow et al., 2016). For convolutional neural

networks (CNN), DL performs hierarchical feature extraction across its multiple layers, where the first layers deal with low-level abstraction features and the later layers compose these low-level features to capture high-level abstraction features. DL algorithms have excelled in many domains such as image and audio recognition, and object detection (Deng et al., 2014).

DL learns by backpropagation (Goodfellow et al., 2016). For an input example, the back propagation algorithm (a gradient descent optimization technique) adjusts the weights of the model (weights determine if a signal is relevant to go to the next layer) to minimize the mismatch between the predicted output of the DL model with respect to its true label. After a number of epochs, the weights are tuned progressively and the training stops when the difference is minimized, the trained model could then be used to predict labels of new inputs (e.g. data in the testing set).

One of the main advantages of DL is that it can learn feature extraction automatically from the labeled data without using human expert knowledge. DL algorithms are well-known for being able to learn from big data sets. Most ML methods' performance plateaus when given a big enough number of training examples, whereas a DL model performance keeps improving and absorbing information with data sets of millions of examples. The downside is that training a deep learning model is generally more time-consuming than the rest of the methods and also, due to its internal complexity of layers and nodes, it is more difficult to understand the internal logic of its predictions.

2.3. Electromagnetic side-channel analysis for forensics

Recent work has suggested to use EM-SCA in digital investigations where classical methods fail (Sayakkara et al., 2018, 2019b). Most of the potential uses of EM-SCA in digital investigations have already been realized in the information security arena. Detection of the specific make and model of a device (Yilmaz et al., 2020), the detection of its internal states of software and hardware (Yilmaz et al., 2019; Chawla et al., 2019), and the retrieval of cryptographic keys from them (Ronen et al., 2017) in realistic scenarios have been performed in the recent past. Therefore, these techniques need to be leveraged into practical digital investigations by providing the necessary facilities to the investigators. Development of tools that are easy-to-use for digital investigators with a minimum knowledge in EM-SCA is such a timely requirement. *EMvidence* is a tool that is currently under constant development towards this goal (Sayakkara et al., 2020b,c). A workshop was delivered at Digital Forensic Research Workshop (DFRWS EU) 2020 describing *EMvidence*'s design and implementation. An explanation for its use, alongside the content of the workshop is available on GitHub.¹

3. Electromagnetic side-channel raw sample dataset

The EM noise emission patterns from IoT devices can depend on a multitude of factors such as the type of the processor, the layout of the printed circuit board (PCB), and the characteristics of the software instructions running on the processor. As a result, it is necessary to experiment with a dataset that is generated on a hardware and software platform easily to reproduce. To this end, the first step was to select a representative IoT device platform that is sufficiently representative of a wider variety of IoT devices and also easy to experiment with.

As the target IoT device, an *Arduino Leonardo* board was select.

¹ <https://github.com/asanka-code/DFRWS-EU2020-Workshop-InsightsfromWaves>.

This is a general-purpose prototype board that has a variety of input/output pins available to be used in a variety of applications. The device consists of an ATmega32u4 micro-controller at its heart that operates at 16 MHz clock frequency with the help of an external oscillator on-board. The device is powered by 5 V power supply through a USB port, which is also used for the purposes of reprogramming the device. Due to the limited random-access memory (RAM) and read-only memory (ROM) capacity (32 KB of flash memory for storing programs and 2.5 KB of RAM), the device cannot run with an operating system. A manufacturer-provided C-like programming language is used to directly write applications to the device that can run without multi-tasking. The small resource profile of the device resembles many battery-operated IoT devices in common use.

When the same software is running on a processor multiple times, it is unlikely that the order of machine code instructions being executed will be the same each time. This is due to the difference in the input data provided to the software in each case. In consequence, the EM emission pattern can get affected and become slightly different from each case as well. For the purpose of simulating such an extreme case, this work uses 10 sorting algorithm implementations each sorting a 100-element long randomly generated integer array. The sorting algorithms are chosen in a way to include a diverse set of computational time complexities in the dataset (see Table 1). The randomness of the input data to the sorting algorithms ensures that their instructions sequences are different at each iteration. Fig. 1 illustrates the average execution time of each sorting algorithm implementation when they are sorting 100 randomly generated integers.

The EM noise of the processor is emitted from various regions of the processor chip across various frequency channels with varying amplitudes. Therefore, when an H-probe antenna is used closer to the processor chip to capture EM noise, the exact location of the antenna over the processor affects the signal being captured. Fig. 2 depicts the procedure of placing an H-probe antenna over the processor of the Arduino Leonardo device. When the antenna is centred directly over the processor, a weaker signal is observed with lower amplitudes and fewer frequency features. By slightly adjusting the position of the antenna over the processor and observing the emission signal, it was empirically identified that when the H-loop antenna is located exactly over the index corner of the chip, the strongest EM emission signal was registered. The necessity for close proximity of the antenna to the chip is due to the low signal amplification of the inexpensive SDR device used in this work. However, with more powerful signal acquisition equipment, it would be possible to capture EM radiation from an increased distance. Fig. 3 (a) and (b) illustrates the emission signal when the H-loop antenna is placed at the centre and the index corner of the processor respectively. Accordingly, the index corner of the processor was chosen as the default antenna position throughout the

rest of the data collection procedure for this work.

When generating the raw sample dataset for the purpose of using it for ML experiments, the following data collection procedure was followed. The data collection was performed in two phases with a time gap of 1 h between them in order to produce two sets of independently taken data for each sorting algorithm. During that time gap, the Arduino Leonardo device was turned off and left unplugged from the power supply. During each phase, the Arduino device was programmed with each sorting program, one at a time. The basic skeleton of each Arduino program is depicted in Algorithm 1. Once programmed with a specific sorting algorithm, the EM emissions of the device are recorded for T time period, where T was set in the two phases to 8 s and 4 s respectively. The rationale being to use them separately as training and testing data for the ML models.

Algorithm 1. Structure of each Arduino program

Algorithm 1: Structure of each Arduino program

Select the sorting algorithm ;

Loop

$array \leftarrow generate\ 100\ random\ integers\ ;$

$result \leftarrow sort\ the\ array\ with\ a\ sorting\ algorithm\ ;$

 reset array and result ;

EndLoop

For data capture, a HackRF SDR² was used with the centre frequency set to 288 MHz and sampling rate set to 20 MHz (Ossmann, 2016). Even though the clock frequency of Arduino Leonardo is 16 MHz, it was empirically identified that certain harmonics of the frequency provide much stronger leakage signals, hence the choice of 18th harmonic, which is 288 MHz. Each EM trace was saved into a NumPy array file (very frequent in python) in order to ease the processing of them later. In phase-1, the EM trace file for each algorithm was approximately 1.4 GB in size causing the first dataset to be approximately 14 GB in size. In phase-2, an EM trace file was approximately 550 MB in size resulting in a dataset of approximately 5.5 GB in size. Fig. 4 illustrates the power spectral density of the recorded signals for 9 of the sorting programs in phase-2. Note that the signal of bubble sort program is shown in Fig. 3).

4. Experimentation and results

In this section, the experiments conducted are outlined to see how well the prediction of the software activity a device is executing performs by just listening to the electromagnetic noise it produces.

4.1. Methodology

For each of our experiments, the methodology outlined in Fig. 5 was followed. It includes five main steps: 1) dataset generation, 2) for a chosen window size slicing time series of raw samples to windows of consecutive raw samples, 3) applying feature engineering, 4) ML model generation, and 5) ML model evaluation.

4.1.1. Preparation of the raw sample dataset

As described in Section 3, the first step involved collecting raw samples for two data sets: the training set and the testing set. For each of the ten sorting algorithms, a time series of approximately

Table 1

Time complexity of the different sorting algorithms used.

Algorithm	Time Complexity		
	Best	Average	Worst
Bubble	$\Omega(n)$	$\Theta(n^2)$	$O(n^2)$
Gnome	$\Omega(n)$	$\Theta(n^2)$	$O(n^2)$
Heap	$\Omega(n \log(n))$	$\Theta(n \log(n))$	$O(n \log(n))$
Insertion	$\Omega(n)$	$\Theta(n^2)$	$O(n^2)$
Merge	$\Omega(n \log(n))$	$\Theta(n \log(n))$	$O(n \log(n))$
Quick	$\Omega(n \log(n))$	$\Theta(n \log(n))$	$O(n^2)$
Radix	$\Omega(nk)$	$\Theta(nk)$	$O(nk)$
Selection	$\Omega(n^2)$	$\Theta(n^2)$	$O(n^2)$
Shell	$\Omega(n \log(n))$	–	$O(n^2)$
Stooge	$\Omega(n^{\log_3 \log 1.5})$	$\Theta(n^{\log_3 \log 1.5})$	$O(n^{\log_3 \log 1.5})$

² <https://greatscottgadgets.com/hackrf/>.

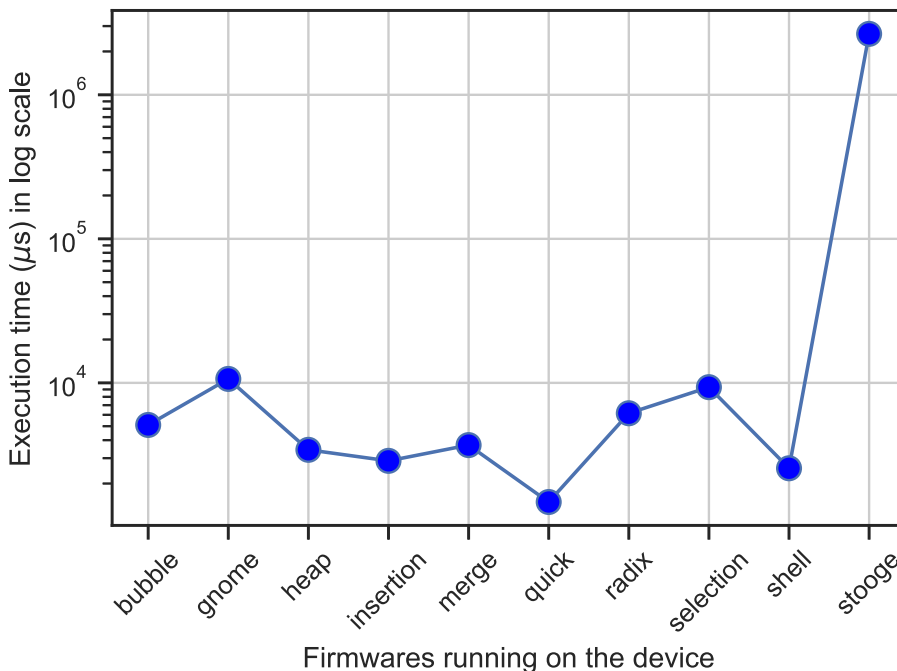


Fig. 1. Average execution time for the 10 sorting programs running on Arduino Leonardo device. The average values were calculated by logging 200 consecutive executions of each program with randomly generated 100-elements long arrays as inputs.

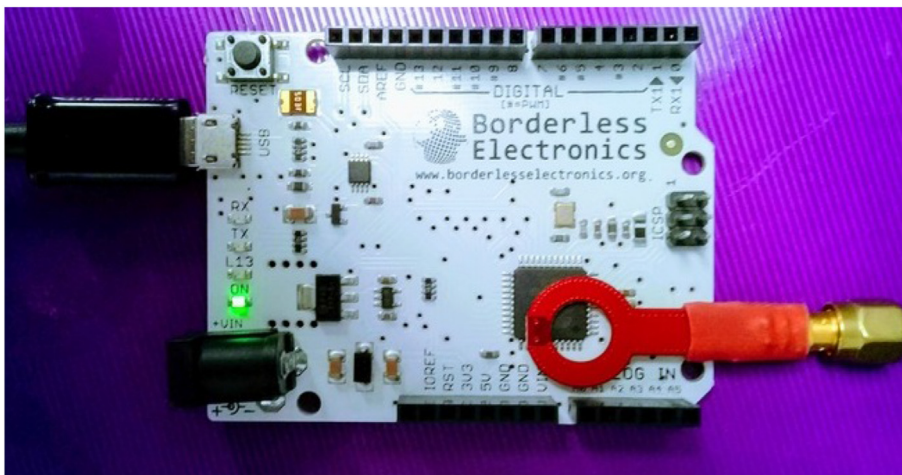


Fig. 2. Positioning of the H-loop antenna over the processor of Arduino Leonardo device to maximize the reception of leakage signal.

175 million raw samples was generated in the training session, and a time series of approximately 76 million raw samples were generated in the testing session. The raw training data and testing data were gathered with a time difference of 1 h ensuring that the antenna was in the same position. Next each time series of raw samples was sliced into windows of consecutive time steps using a chosen time window, each window of raw samples would then be transformed into an input example through feature engineering. Different time windows were applied to see the performance with each size. There is a trade-off between the accuracy and the need to process as little information as possible. The larger the window, the higher the accuracy of the model, but the more expensive it is making the predictions. All the raw samples were not used to generate training examples and testing examples for a chosen window size. Instead, we configured the number of training examples and the number of testing examples specifically for each of

our experiments. A typical dataset configuration for our experiments would include 25,000 examples in the training set and 25,000 examples in the testing set.

4.1.2. Feature engineering

Each window of a chosen window size of consecutive raw samples was converted to an example, where the number of features in the example would depend on our feature engineering approach. The first approach is directly working with samples from the time domain, so the number of features in each example is equal to the chosen window size. It has the advantage that there is no extra cost of running any transformation algorithm. The second option is to transform each time window into the frequency domain by applying the Fast Fourier transform (FFT), which calculates the Discrete Fourier Transform (DFT) of a given time window, again the number of frequency features in each example is

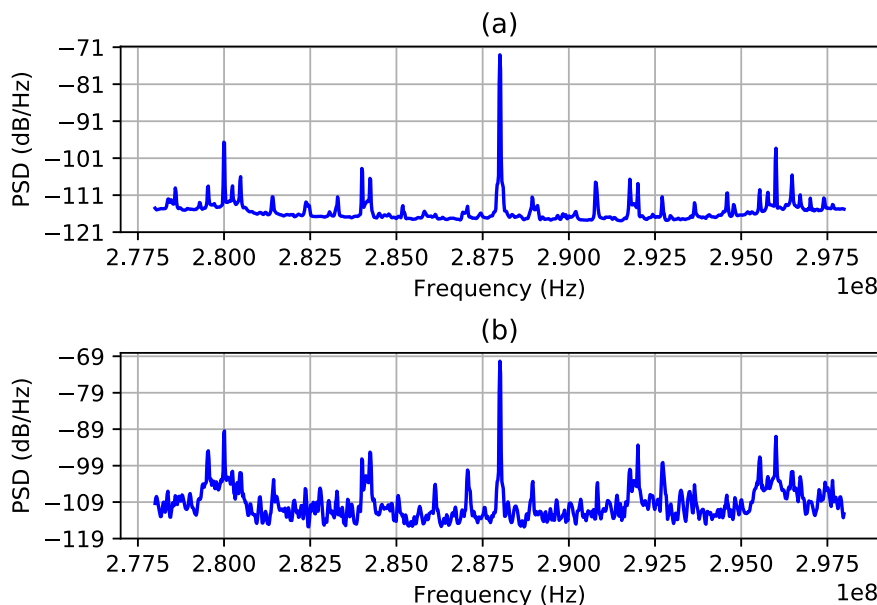


Fig. 3. Power spectral density (PSD) of the EM emission when running Bubble sort algorithm with two different antenna positions.

equal to the chosen window size, as shown in Fig. 6.

Lastly, it is also quite effective to extract features in the time domain and in the frequency domain. This drastically reduces the information taken by the classifier (the input of the model), but at the same time, it increases the computational costs since all the extracted features have to be calculated.

As with other work in the area related to time series domain (Ponce et al., 2016), the features extracted in the time domain are: mean of the signal, standard deviation, root mean square, variance, median, number of the zero-crossing points, skewness, kurtosis, first and third quartile, interquartile range, and correlation. And the features extracted in the frequency domain are: average, median, variance, root mean square, entropy of the signal, Shannon spectral entropy, spectral centroid, flatness, mode frequency, peak frequency, kurtosis, and skewness.

4.1.3. Machine learning models

In this work, a wide range of machine learning models were considered. Support Vector Machines, Decision Trees, K-Nearest Neighbour, and XGBoost were evaluated with their default settings. A Random Forest Classifier was used with 500 trees and a maximum number of features to be the square root of the number of columns (number of input features), together with bootstrapping activated. Each the above classifiers were trained on input data after being preprocessed by the MinMaxScaler transformation, i.e., each input feature was scaled and translated on the training set to be in the range between 0 and 1; and the same transformation was applied on the data in the test set.

For the deep learning model, a Convolutional Neural Network was used with four layers. The first two layers were two convolutional layers - using 32 filters and 64 filters respectively. The outputs of the second convolutional layer were flattened and connected to a fully connected layer with 64 nodes, which was then connected to the final output layer to predict the activity label. As the convolutional layers in the model were used to extract features across a fixed window of consecutive features, the kernel size for both the convolutional layers was set to 8 and large strides were used (2 and 4) - together with a max pooling value of 4 to increase the depth of field of each node in the first dense layer and to reduce

the model complexity. The CNN model was trained with the standard Adam optimization algorithm.

4.1.4. Model training and evaluation

For each experiment and for a chosen machine learning model, we trained it on the training set, and evaluated its performance on the testing set using classification accuracy as the performance metric.

Once the ML model is built, if a real-world deployment was needed, there would only be three steps: applying the time window, transforming into the frequency domain with the FFT, and predicting the activity with the generated ML model.

4.2. Results

In the following subsection we present an explanation and the results for the different experiments.

4.2.1. Experiment I: selecting between time domain, frequency domain, and feature extraction

Our first experiment was designed to select the appropriate feature engineering to apply on each window of raw sample. We considered three feature engineering approaches, as discussed earlier: 1) using directly a window of raw samples from the time domain, 2) converting each window of raw samples from the time domain to the frequency domain using FFT, and 3) using feature selection. The experiment was conducted using 25,000 examples for the training set and 25,000 examples for the testing set. A Random Forest Classifier was employed with the hyper-parameters as described in 4.1. As shown in Fig. 7, according to the first experiment, predicting from the amplitudes of the frequency domain is much more accurate than the other options. In Fig. 8, the confusion matrix of the best model is shown; random forest with a time window of 1,000 samples in the frequency domain and an accuracy of 93.56%.

4.2.2. Experiment II: changing the antenna position

One interesting question in this research is how the position of the antenna on the device affects the performance of the classifiers.

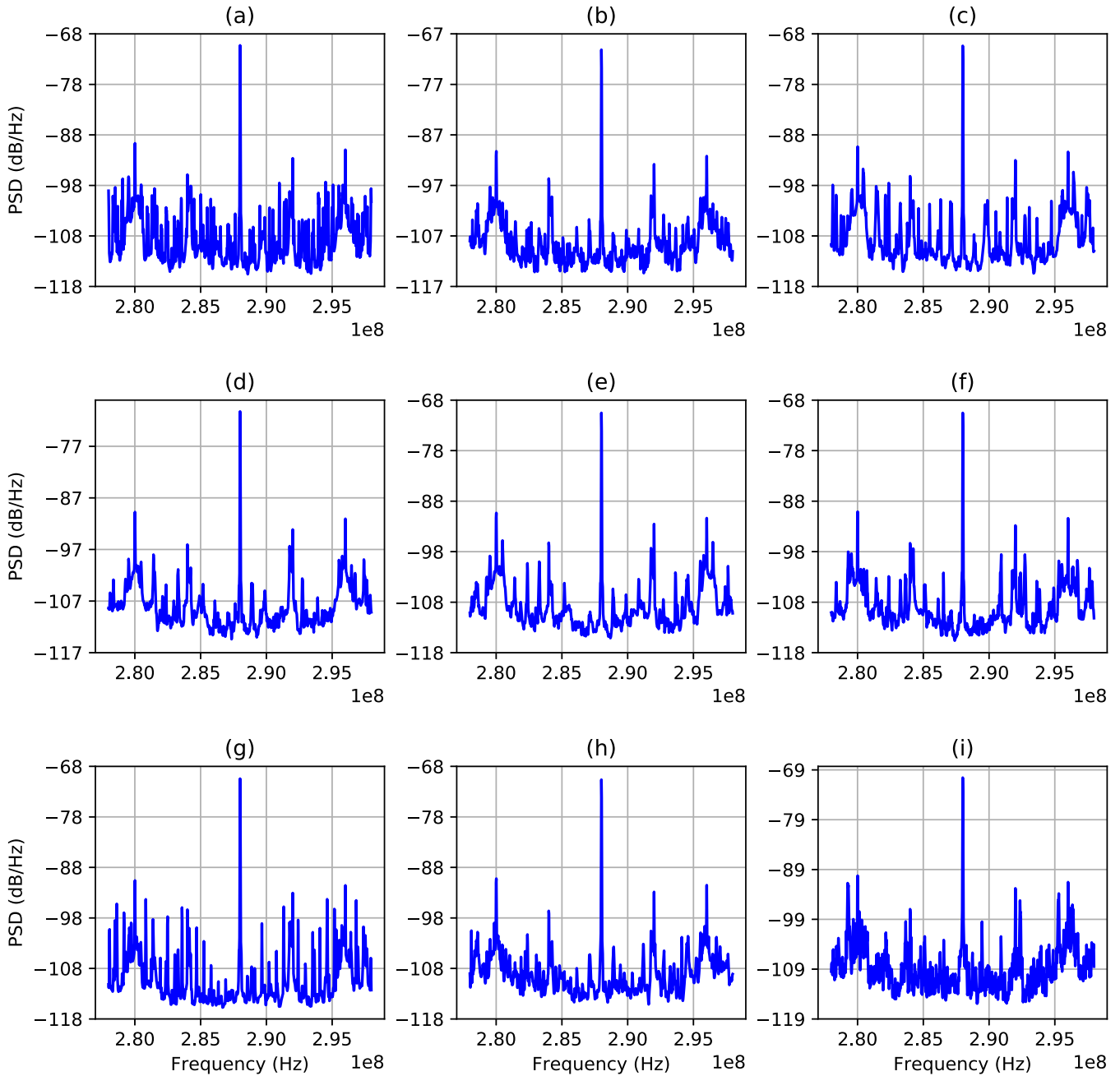


Fig. 4. Power spectral density (PSD) of the EM emission when running 9 sorting algorithms. (a) Gnome sort, (b) Heap sort, (c) Insertion sort, (d) Merge sort, (e) Quick sort, (f) Radix sort, (g) Selection sort, (h) Shell sort, and (i) Stoooge sort. Bubble sort is not depicted here.



Fig. 5. Steps used in the methodology to predict the sorting algorithm executed by the Arduino.

First, the EM traces were created with samples of the ten sorting algorithms with arrays of 100 elements. Second, a new set of EM traces were created with the same algorithms with the same number of elements but, with the position of the antenna changed. After conducting this experiment, the accuracy of the model was

12.75%. As can be seen in Fig. 9, the accuracy of the classes is very low.

4.2.3. Experiment III: selecting a machine learning algorithm

Experiment 1 demonstrates that machine learning models

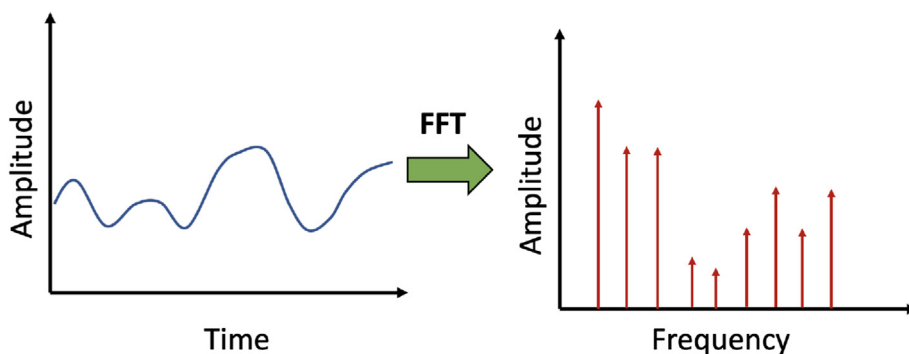


Fig. 6. Transforming a signal from the time domain to the frequency domain with the Fast Fourier Transform.

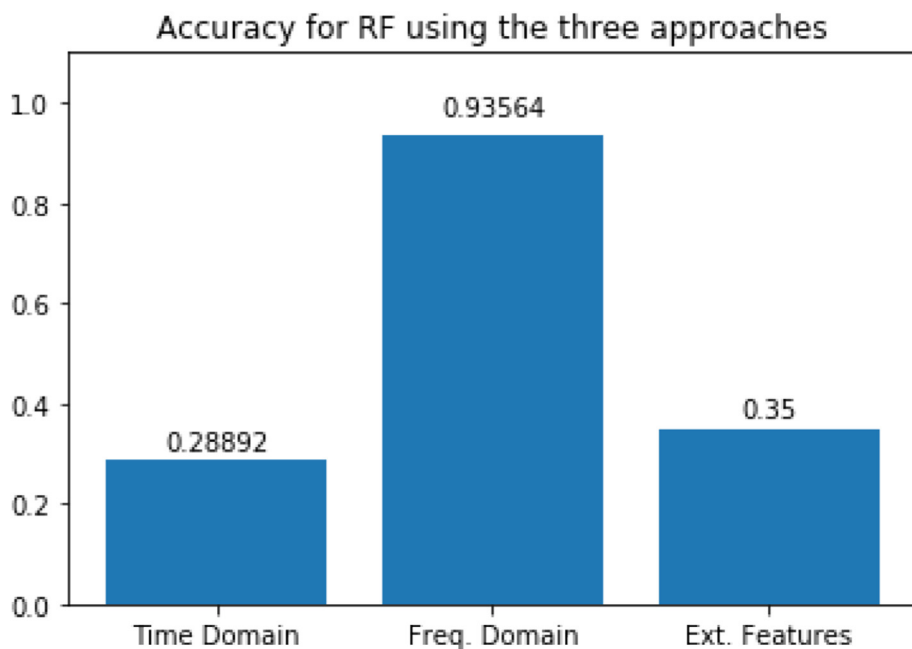


Fig. 7. Accuracy of the predicting in time domain, frequency domain, and with extracted features in time and frequency.

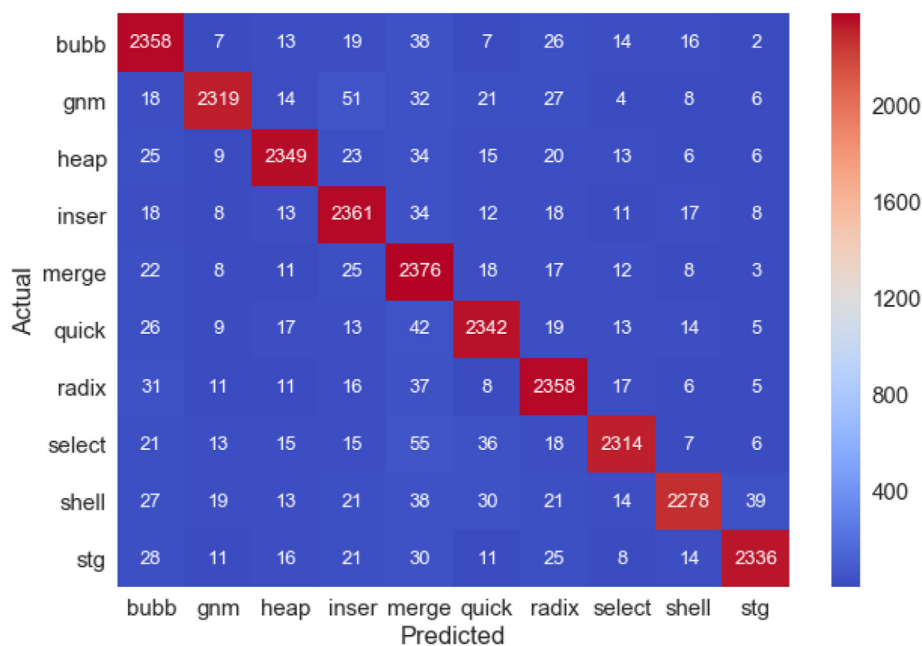


Fig. 8. Confusion matrix of the first experiment.

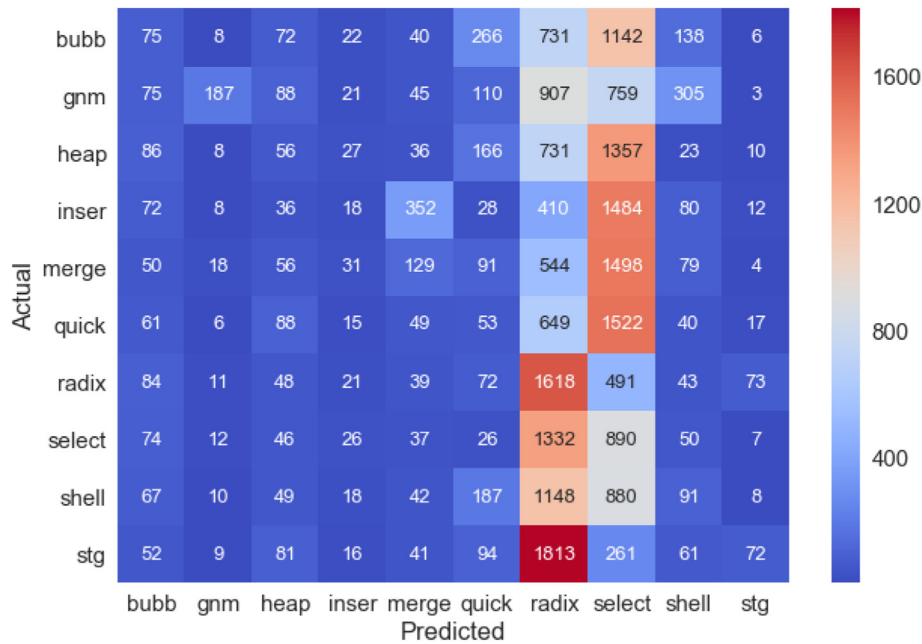


Fig. 9. Confusion matrix of the first experiment.

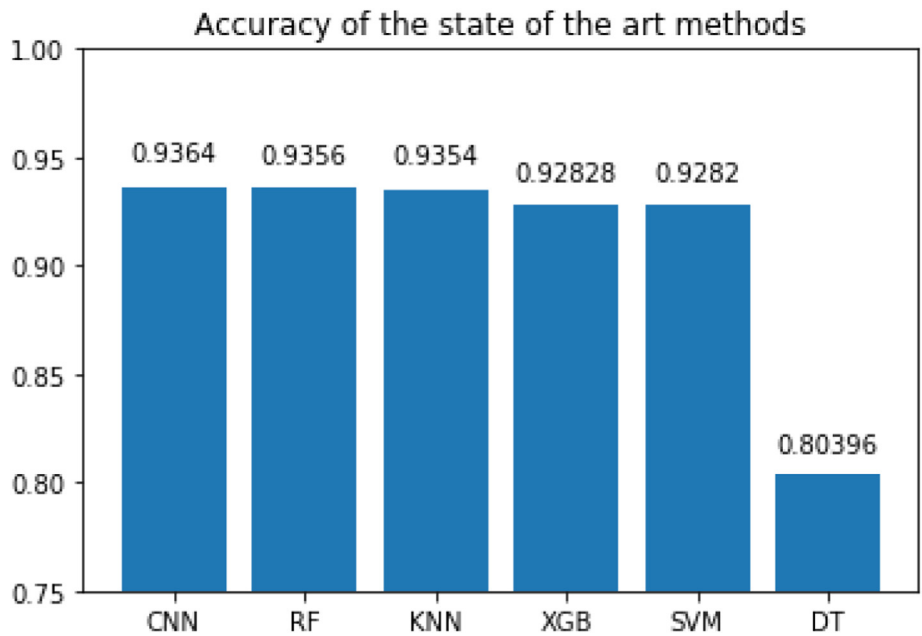


Fig. 10. Performance of different ML methods trained with 25,000 samples with the ten different sorting algorithms.

performed best on data in the frequency domain. In this experiment, different machine learning algorithms were evaluated to see which ones performed best when being trained on data in the frequency domain: Support Vector Machines, Random Forests, K Nearest Neighbour, XGBoost, Decision Trees, Random Forests, and CNN. As for the dataset, the same dataset in the frequency domain was used as in Section 4.2.1 with a training set of 25,000 examples, a test set of 25,000 examples, and the window size of 1,000 time steps.

The performance of all the models is shown in Fig. 10 and the time for running the experiments with each method is shown in Table 2. From the results, the two strongest models were selected, i.e., CNN and Random Forests, to carry out the following

experiments.

4.2.4. Experiment IV: impact of sample window size

In this experiment, the effect of choosing different window sizes for the signal. For the CNN and RF models, a list of window sizes of 500, 1,000, 2,000, 5,000, 10,000, and 15,000 time steps were chosen. Each time series of a chosen window size is converted into an input example of the same size with an equal number of features in the frequency domain using Discrete Fourier Transform (DFT).

For each window size, a training set of 25,000 examples was used from the training recording session, and a test set of 25,000 examples was used from the test recording session.

The results of the experiments are shown in Fig. 11. For a small

Table 2
Time in seconds to build the ML model.

XGB	SVM	DT	K-NN	RF	CNN
402	1027	50	1240	229	250

window size, i.e., the window size of equal or less than 500, RF performed better than the CNN model. For all the window sizes equal or bigger than 1,000, the CNN performed better than RF, even though the margin of improvement is small. The CNN model trained on input examples with a window size of 15,000 time steps achieved an accuracy of 99.6% on the test data.

4.2.5. Experiment V: impact of the number of examples

These experiments were designed to study the effect of the number of training examples on the performance of the two models. For all the experiments in this Section, the window size of one input was fixed to be the default value of 1,000: each 1,000 time step signal time series was converted to an input sample of size 1,000 in the frequency domain using DFT.

For each of the number training examples n_{train} in the list of {5,000, 10,000, 15,000, 25,000, 40,000, 80,000}, a training set of n_{train} examples was generated, while the test set was fixed to be of 25,000 examples. As shown in Fig. 12, as the number of training examples varied from 5,000 to 80,000, the accuracy of the CNN model on the test set varied in the range 92.9%–93.8%; while the CNN model performed better than RF for all values of n_{train} , with a bigger difference with a smaller number of training examples.

4.2.6. Performance of the best model on the whole test recording session

The model that performed best in the experiments in Sections 4.2.4 and 4.2.5 was the CNN model with its input computed from a window of 15,000 time steps; it was trained with 25,000 examples and achieved an accuracy of 99.6% on the test set of 25,000 examples. To verify the performance of this trained model, its accuracy was computed on all the examples of the test recording session. Fig. 13 shows the confusion matrix of this test, where the model achieved an accuracy of 99.8%.

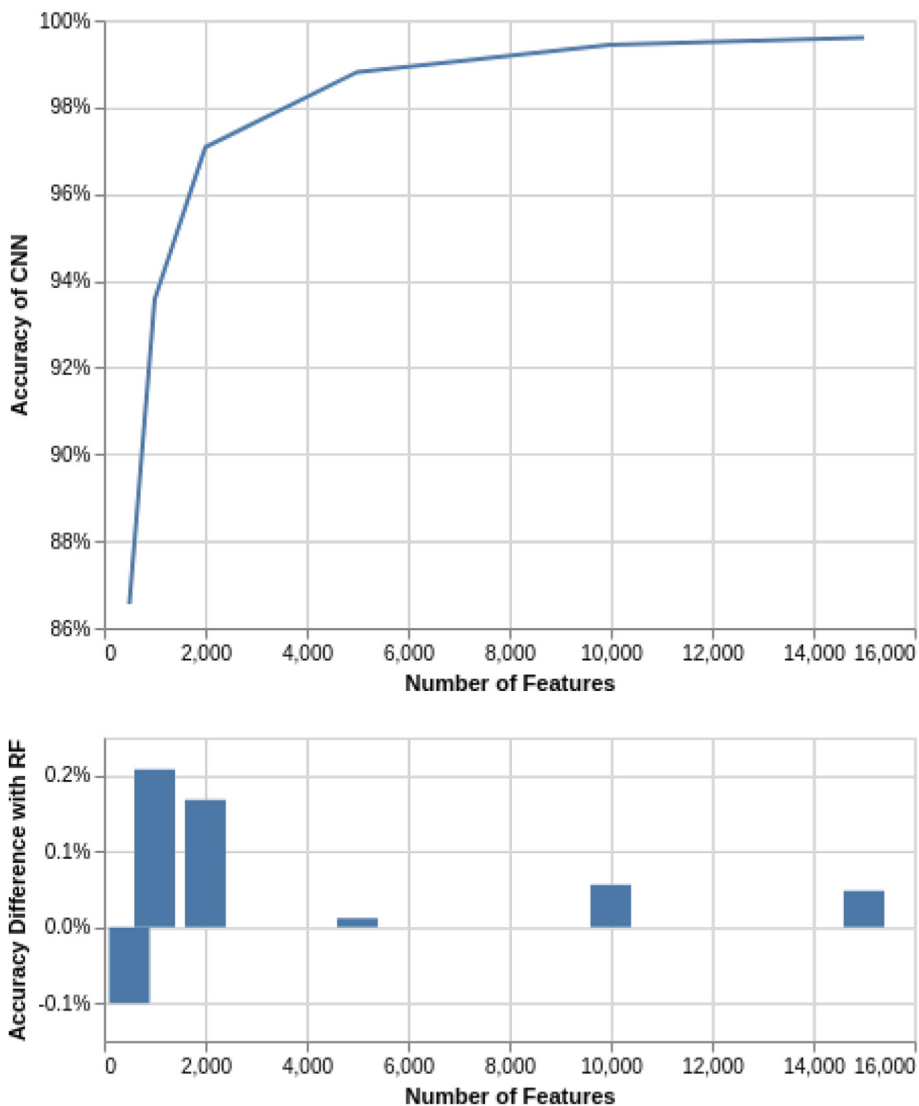


Fig. 11. Experiment of varying the window size (number of features).

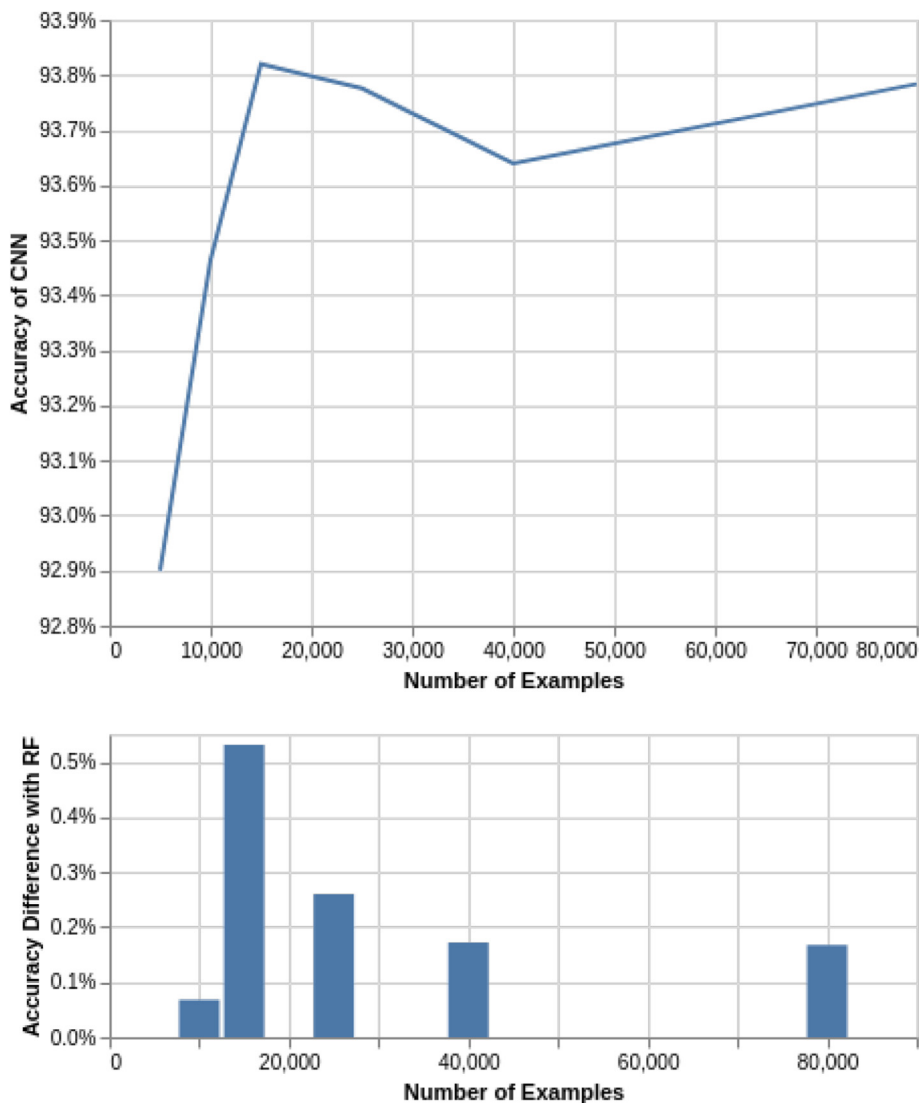


Fig. 12. Experiment of varying the number of training examples.

4.3. Discussion

The experiments outlined in Section 4.2.1 demonstrate that ML models perform best when working on data in the frequency domain. One of the contributions of this work is that a consistent data recording protocol is necessary for machine learning models trained on past data to perform well on data recorded in the future. In this work, this generalisation performance was achieved by fixing the location of the antenna with respect to the Arduino device. In contrast, when the position of the antenna was moved between the two recording sessions, the model trained on the training set performed very poorly on the test data, reaching almost random performance.

The experiments on varying the window size to calculate input show that the performance of the deep learning model improved strongly when the window size was increased – going from 86.6% for a window size of 500 time steps to 99.6% for a window size of 15,000 time steps, as shown in Fig. 11. On the other hand, when the window size was fixed to 1,000 and varied the number of training examples from 5,000 to 80,000 the performance of the CNN model only improved from 93.0% to 93.8%, as shown in Fig. 12). For comparison, while the amount of training data required for a CNN

model trained on 80,000 examples of a window size 1,000 time steps is bigger, its performance of 93.8% accuracy is inferior to the performance of 97.0% accuracy for a CNN model trained on 25,000 training examples of a window size of 2,000 time steps. It can be concluded that choosing the right window size is more impactful than increasing the number of training examples on the performance of the ML model. For instance, a window size of 2,000 time steps (corresponding to an accuracy of 97.0%) is a sweet spot of achieving a high predictive performance of the trained model and having a low processing time.

This work has shown that deep learning models can be used to precisely identify subtle details of running programs on IoT devices opening up the opportunity to acquire forensic insights from IoT devices during the triage examination phase of a digital forensic investigation. The trained models can be integrated with tools, such as EMvidence (Sayakkara et al., 2020b) for the application in practical forensic investigations. As it has been revealed from the experiments, the exact antenna position during data acquisition to build deep learning models should be used consistently in an investigation scenario to acquire EM data from a target device. This limit could be addressed by a more robust training data collection protocol in the future.

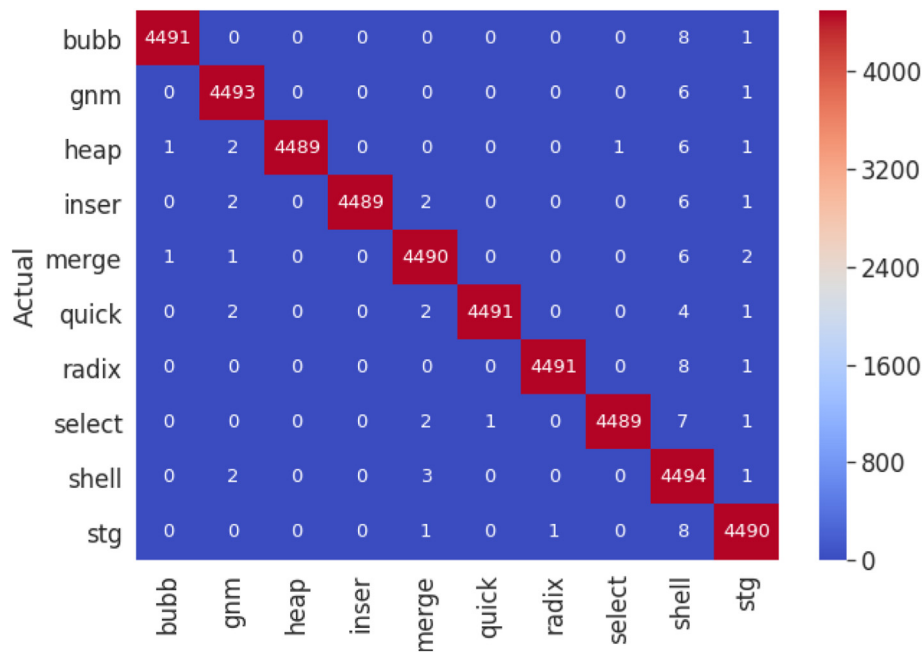


Fig. 13. Confusion matrix of the model with the best performance when testing on the data from the whole test recording session.

The high accuracy achieved with data produced by signal observation for just a few seconds hints that it is sufficient to capture data for a very short period of time in real-world scenarios.

5. Conclusions and future work

In this work, we used a wide variety of software activities as target classes for ML classification. Therefore, the dataset is sufficiently diverse in terms of computational complexity. The objective was to identify machine learning methods that are suitable to classify a wide variety of software activities running on IoT devices. The experimentation proves that it is possible to predict the activity that a device is performing from the generated EM noise with high levels of accuracy. ML classifiers are able to detect the executing algorithm even when the elements are ordered differently, which obviously requires a different order in the execution of the instructions of each sorting algorithm. Experiments have been conducted over well-known sorting algorithms to identify their performance. Among the machine learning algorithms tested, CNNs were found to be most effective in identifying individual software activities.

5.1. Future work

As future work, two interesting avenues for investigation have been identified. First, it is worth exploring different sizes for the arrays. For example, what performance is achievable with an ML model trained with arrays of 100 elements when the size of the array of the testing examples is 10, 50, 500, 1,000, or 10,000 elements. We could also create a more robust deep learning model by training it on data recorded from multiple antenna locations, such a model should be more robust with respect to the antenna location when recording test data. As the new dataset will be more complex, we expect to develop more complex deep learning architectures to fit it, which could include a combination of residual layer, attention layer, or multiple input modules to integrate multiple processing units. The logic behind the predictions of the models was not a focus in this work. It would be interesting to know what features

(frequencies in this case) are more relevant for ML models. Additionally, the application of different explainable machine learning techniques (Molnar, 2019; Montavon et al., 2018) on a trained model to understand why and how the model makes a particular decision.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

References

- Blossom, E., 2004. GNU Radio: Tools for exploring the radio frequency spectrum. *Linux J.* 2004 (122), 4.
- Breiman, L., 2001. Random forests. *Mach. Learn.* 45 (1), 5–32.
- Camurati, G., Poeplau, S., Muench, M., Hayes, T., Francillon, A., 2018. Screaming channels: when electromagnetic side channels meet radio transceivers. In: *Proceedings of the 25th ACM Conference on Computer and Communications Security (CCS)*. CCS '18. ACM.
- Cass, S., 2013. A \$40 software-defined radio. *IEEE Spectrum* 50 (7), 22–23.
- Chawla, N., Singh, A., Kar, M., Mukhopadhyay, S., 2019. Application inference using machine learning based side channel analysis. In: *2019 International Joint Conference on Neural Networks (IJCNN)*. IEEE, pp. 1–8.
- Deng, L., Yu, D., et al., 2014. Deep learning: methods and applications. *Found. Trends® Signal Process.* 7 (3–4), 197–387.
- Du, X., Hargreaves, C., Sheppard, J., Anda, F., Sayakkara, A., Le-Khac, N.A., Scanlon, M., 2020. SoK: Exploring the state of the art and the future potential of artificial intelligence in digital forensic investigation. In: *Proceedings of the 15th International Conference on Availability, Reliability and Security*. ACM. ISBN 9781450388337.
- Elibol, F., Sarac, U., Erer, I., 2012. Realistic eavesdropping attacks on computer displays with low-cost and mobile receiver system. In: *Signal Processing Conference (EUSIPCO), 2012 Proceedings of the 20th European*. IEEE, pp. 1767–1771.
- Ettus, M., Braun, M., 2015. *The Universal Software Radio Peripheral (USRP) Family of Low-Cost SDRs. Opportunistic Spectrum Sharing and White Space Access: The Practical Reality*. Wiley, pp. 3–23.
- Genuer, R., Poggi, J.M., Tuleau-Malot, C., 2010. Variable selection using random forests. *Pattern Recogn. Lett.* 31 (14), 2225–2236.
- Goodfellow, I., Bengio, Y., Courville, A., 2016. *Deep Learning*. MIT press.
- Kocher, P., Jaffe, J., Jun, B., 1999. Differential power analysis. In: *Advances in Cryptology (CRYPTO '99)*. Springer, p. 789, 789.
- LeCun, Y., Bengio, Y., Hinton, G., 2015. Deep learning. *Nature* 521 (7553), 436.
- Lin, J., Yu, W., Zhang, N., Yang, X., Zhang, H., Zhao, W., 2017. A survey on Internet of

- Things: Architecture, enabling technologies, security and privacy, and applications. *IEEE Internet of Things J.* 4 (5), 1125–1142.
- Marsland, S., 2015. *Machine Learning: an Algorithmic Perspective*. CRC press.
- Maxwell, J.C., 1865. A dynamical theory of the electromagnetic field. *Phil. Trans. Roy. Soc. Lond.* 155, 459–512.
- Meffert, C., Clark, D., Baggili, I., Breiting, F., 2017. Forensic State Acquisition from Internet of Things (FSAIoT): a general framework and practical approach for IoT forensics through IoT device state acquisition. In: *Proceedings of the 12th International Conference on Availability, Reliability and Security*. ACM, p. 56.
- Molnar, C., 2019. *Interpretable Machine Learning*. Lulu. com.
- Montavon, G., Samek, W., Müller, K.R., 2018. Methods for interpreting and understanding deep neural networks. *Digit. Signal Process.* 73, 1–15.
- Nazari, A., Sehatbakhsh, N., Alam, M., Zajic, A., Prvulovic, M., 2017. EDDIE: EM-based Detection of deviations in program execution. In: *Proceedings of the 44th Annual International Symposium on Computer Architecture*. ACM, pp. 333–346.
- Ossmann, M., 2016. Software defined radio with HackRF. *Great Scott Gadgets*. <https://greatscottgadgets.com/sdr>.
- Peeters, E., Standaert, F.X., Quisquater, J.J., 2007. Power and electromagnetic analysis: improved model, consequences and comparisons. *Integrat. VLSI J.* 40 (1), 52–60.
- Ponce, H., Miralles-Pechuán, L., Martínez-Villaseñor, M.D.L., 2016. A flexible approach for human activity recognition using artificial hydrocarbon networks. *Sensors* 16 (11), 1715.
- Ronen, E., Shamir, A., Weingarten, A.O., O'Flynn, C., 2017. IoT goes nuclear: creating a ZigBee chain reaction. In: *2017 IEEE Symposium on Security and Privacy (SP)*. IEEE, pp. 195–212.
- Samuel, A.L., 1959. Some studies in machine learning using the game of checkers. *IBM J. Res. Dev.* 3 (3), 210–229.
- Santner, J., Unger, M., Pock, T., Leistner, C., Saffari, A., Bischof, H., 2009. Interactive texture segmentation using random forests and total variation. In: *BMVC*. Citeseer, pp. 1–12.
- Sayakkara, A., Le-Khac, N.A., Scanlon, M., 2018. Electromagnetic side-channel attacks: Potential for progressing hindered digital forensic analysis. In: *Companion Proceedings for the ISSTA/ECOOP 2018 Workshops*. ISSTA '18. Association for Computing Machinery, New York, NY, USA, pp. 138–143. <https://doi.org/10.1145/3236454.3236512>. ISBN 9781450359399 doi:10.1145/3236454.3236512.
- Sayakkara, A., Le-Khac, N.A., Scanlon, M., 2019a. Leveraging electromagnetic side-channel analysis for the investigation of IoT devices. *Digit. Invest.* 29, S94–S103. <https://doi.org/10.1016/j.diin.2019.04.012>. <http://www.sciencedirect.com/science/article/pii/S1742287619301616>.
- Sayakkara, A., Le-Khac, N.A., Scanlon, M., 2019b. A survey of electromagnetic side-channel attacks and discussion on their case-progressing potential for digital forensics. *Digit. Invest.* 29, 43–54. <https://doi.org/10.1016/j.diin.2019.03.002>.
- Sayakkara, A., Miralles-Pechuán, L., Le-Khac, N.A., Scanlon, M., 2020a. Cutting through the emissions: feature selection from electromagnetic side-channel data for activity detection. *Forensic Sci. Int.: Digit. Invest.* 32, 300927. <https://doi.org/10.1016/j.fsidi.2020.300927>. <https://www.sciencedirect.com/science/article/pii/S2666281720300226>.
- Sayakkara, A., Le-Khac, N.A., Scanlon, M., 2020b. EMvidence: a framework for digital evidence acquisition from IoT devices through electromagnetic side-channel analysis. *Forensic Sci. Int.: Digit. Invest.* 32, 300907. <https://doi.org/10.1016/j.fsidi.2020.300907>. <http://www.sciencedirect.com/science/article/pii/S2666281720300020>.
- Sayakkara, A., Le-Khac, N.A., Scanlon, M., 2020c. Facilitating electromagnetic side-channel analysis for IoT investigation: evaluating the EMvidence framework. *Forensic Sci. Int.: Digit. Invest.*
- Tuttlebee, W.H., 2003. *Software Defined Radio: Enabling Technologies*. John Wiley & Sons.
- Van Eck, W., 1985. Electromagnetic radiation from video display units: an eavesdropping risk? *Comput. Secur.* 4 (4), 269–286.
- Watson, S., Dehghantanha, A., 2016. Digital forensics: the missing piece of the internet of things promise. *Comput. Fraud Secur.* 2016 (6), 5–8.
- Yaqoob, I., Hashem, I.A.T., Ahmed, A., Kazmi, S.A., Hong, C.S., 2019. Internet of things forensics: recent advances, taxonomy, requirements, and open challenges. *Future Generat. Comput. Syst.* 92, 265–275.
- Yilmaz, B.B., Ugurlu, E.M., Prvulovic, M., Zajic, A., 2019. Detecting cellphone camera status at distance by exploiting electromagnetic emanations. In: *MILCOM 2019–2019 IEEE Military Communications Conference (MILCOM)*. IEEE, pp. 1–6.
- Yilmaz, B.B., Ugurlu, E.M., Zajic, A., Prvulovic, M., 2020. Cell-phone classification: a convolutional neural network approach exploiting electromagnetic emanations. In: *ICASSP 2020–2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, pp. 2862–2866.
- Zankl, A., Seuschek, H., Irazoqui, G., Gulmezoglu, B., 2018. Side-channel attacks in the internet of things: threats and challenges. In: *Solutions for Cyber-Physical Systems Ubiquity*. IGI Global, pp. 325–357.