

AuthCODE: A Privacy-preserving and Multi-device Continuous Authentication Architecture based on Machine and Deep Learning

Pedro Miguel Sánchez Sánchez^{a,*}, Lorenzo Fernández Maimó^b, Alberto Huertas Celdrán^c and Gregorio Martínez Pérez^a

^aDepartment of Information and Communications Engineering, University of Murcia, Murcia 30100, Spain

^bDepartment of Computer Engineering, University of Murcia, Murcia 30100, Spain

^cCommunication Systems Group (CSG), Department of Informatics (IfI), University of Zurich UZH, 8050 Zürich, Switzerland

ARTICLE INFO

Keywords:

Continuous Authentication
Multi-device Behaviour
Smart Office
Machine Learning
Deep Learning

ABSTRACT

The authentication field is evolving towards mechanisms able to keep users continuously authenticated without the necessity of remembering or possessing authentication credentials. While relevant limitations of continuous authentication systems -high false positives rates (FPR) and difficulty to detect behaviour changes- have been demonstrated in realistic single-device scenarios, the Internet of Things and next generation of mobile networks (5G) are enabling novel multi-device scenarios, such as Smart Offices, that can help to reduce or address the previous challenges. The paper at hand presents an AI-based, privacy-preserving and multi-device continuous authentication architecture called AuthCODE. AuthCODE seeks to improve single-device solutions limitations by considering additional behavioural data coming from heterogeneous devices. AuthCODE proposes a novel set of features that combine the interactions of users with different devices. The features relevance has been demonstrated in a realistic Smart Office scenario with several users that interact with their mobile devices and personal computers. In this context, a set of single- and multi-device datasets have been generated and published to compare the performance of our multi-device solution against single-device approaches. A pool of experiments with machine and deep learning classifiers measured the impact of time in authentication accuracy and improved the results of single-device approaches by considering multi-device behaviour profiles. Specifically, the multi-device approach using XGBoost with 1-minute window of aggregated features, achieved a 69.33%, 59.65% and 89.35% improvement in the FPR when compared to the single-device approach for computer, mobile applications and mobile sensors respectively. Finally, temporal information classified by a Long-Short Term Memory Network, allowed the identification of additional complex behaviour patterns.

1. Introduction

Continuous authentication systems pretend to improve the limitations of traditional mechanisms, which authenticate users from time to time according to credentials such as passwords, codes, or tokens (Almalki et al., 2019). In this context, continuous authentication mechanisms increase the level of security, keeping users authenticated permanently, and enhance the users' quality of experience (QoE), being non-intrusive and minimizing the usage of credentials during the authentication processes (Gonzalez-Manzano et al., 2019).

Existing continuous authentication mechanisms model the user's behaviour when he/she uses a particular device for a given time. We understand by "behaviour" the set of actions that a given user performs somewhat consciously with one or more devices. As an example, for PC devices the user's behaviour could be determined by the mouse movements, keystrokes, or applications opened and closed. Regarding mobile devices, the behaviour might be related to screen interactions, device movements, or patterns to use applications. From a multi-device perspective, the user's behaviour could be characterized by the device used at each particular mo-

ment of the day, the duration that each device is used, or the type of application used in each device. As a result of user monitoring, a behaviour profile is created and stored in a dataset, being the precise selection of dimensions, data and features critical to create accurate behaviour profiles. The next step is to use this profile to train Machine Learning (ML) or Deep Learning (DL) models, which can be classifiers or anomaly detectors depending on the nature of the data and the approached scenario. Finally, the continuous authentication mechanism uses these models to evaluate the similarity between the current device usage profile and the learned user's profile, providing either a user's identification or an anomaly score which can be used to decide if the user is authenticated or not.

Different solutions have implemented the previous steps for single-device scenarios such as personal computers and smartphones (Jorquera Valero et al., 2018). However, they present some important limitations such as the high rate of false positives, which happens when users change some aspect of their behaviour and the system does not recognize them. The lack of privacy considerations and the impossibility of detecting some impersonation attacks are two additional limitations of single-device solutions. In this context, the evolution of information and communications technology such as 5G networks (Huertas Celdrán et al., 2019) or the Internet of Things (IoT) (Firouzi and Farahani, 2020) is influencing the applicability of continuous authentication in multi-device

*Corresponding author. Email address: pedromiguel.sanchez@um.es (P.M.S. Sánchez)

ORCID(s): 0000-0002-6444-2102 (P.M.S. Sánchez);
0000-0003-2027-4239 (L.F. Maimó); 0000-0001-7125-1710 (A.H. Celdrán);
0000-0001-5532-6604 (G.M. Pérez)

scenarios to reduce the limitations of single-device scenarios by including additional information coming from the user's activity on diverse devices used during a time window. On the one hand, multi-device continuous authentication mechanisms could consider data from different devices to decide if the user changed his/her behaviour, which could reduce the false positive rate. On the other hand, multi-device scenarios such as Smart Offices (Qolomany et al., 2019) could also get a benefit from non-invasive and robust multi-device continuous authentication, using them to control the access to sensitive data managed by heterogeneous devices such as IoT devices, tablets, smartphones, or computers. As an example, if an attacker studies the common activities and routines of a user protected by a single-device solution, this attacker can learn his behavior and then mimic him to some extent, performing malicious activities. However, having information from several devices improves the robustness and flexibility of the authentication system since it would be necessary to have access to several devices and mimic the user's behaviour in each of them.

Despite the benefits of existing continuous authentication solutions, their design, implementation and integration in new multi-device scenarios are open and challenging issues. In this sense, we emphasize the following challenges: 1) what features and data are relevant to create collaborative rich users' behaviour profiles in multi-device scenarios; 2) whether multi-device continuous authentication mechanisms can improve the performance of single-device mechanisms in terms of false positive and negatives; 3) how the number or interactions and time affect the accuracy of continuous authentication mechanisms; and 4) how behavioural data is collected and processed in order to guarantee users' privacy.

The main contributions of this paper are the following ones:

- A multi-device continuous authentication architecture, called AuthCODE, that guarantees the privacy of users' sensitive data while improving the authentication performance of single-device solutions. AuthCODE considers a hybrid approach that combines the Mobile Edge Computing (MEC) and Cloud Computing paradigms. Privacy-preserving features are calculated in the MEC and sent to the cloud, where ML/DL models are trained and evaluated to keep the authentication performance.
- A set of privacy-preserving and multi-device features able to model precisely the users' behaviour when they interact, in a collaborative way, with heterogeneous devices during different time windows. In this way, the user's activities can be accurately characterized, but if the resultant dataset falls into the hands of an attacker, it would be harmless as the user's monitored activities do not contain sensitive information such as passwords or other typed text.
- Five datasets, available in (Sánchez Sánchez et al., 2020), which contain single- and multi-device privacy-preserving features modelling the users' behaviour.

The datasets have been created by modelling a realistic Smart Office scenario where five users interact with their computers and mobile devices. Regarding computer interaction, the data aggregate information about keys pressed, mouse actions, as well as application and resource usage statistics. With regard to mobile devices, the data consist of sensor values and application usage statistics.

- A pool of experiments with ML and DL classifiers performed over the previous five datasets that demonstrated how the multi-device behaviour profile of AuthCODE improves the false positive rates (FPR) and f1-score metrics of our single-device approach. Specifically, the f1-score and FPR average reached for multi-device profiles of 1-minute windows were 99.33% and 0.23%, respectively, while the same metrics for single-device authentication were 97.39% and 0.72% in the case of personal computers, 96.70% and 0.57% for mobile app statistics, and 90.36% and 2.16% for mobile sensors. Therefore, the improvement achieved in FPR was 69.33%, 59.65% and 89.35% for computers, mobile app statistics and mobile sensors respectively. Finally, to take advantage of the complex patterns present in multi-device profiles, a set of Long-Sort Term Memory (LSTM) neural network configurations obtained an average f1-score of 89% and FPR of 2.02% for a 5-minute sliding window and +99% and 0.37%, respectively, using +60-minute sliding window.

The remainder of the paper is organized in the following way. Section 2 discusses the related work focused on continuous authentication for multi-device scenarios as well as single-device such as personal computers and mobile devices. A motivating use case is detailed in Section 3. The architectural design of AuthCODE is explained in Section 4. The implementation details of the proposed solution as well as a realistic Smart Office scenario are explained in Section 5. Section 6 measured the performance of AuthCODE in the Smart Office environment. Finally, Section 7 shows the conclusions and future work.

2. Related Work

This section reviews the main continuous authentication solutions found in the literature. A wide variety of continuous authentication proposals focused on single- and multi-device scenarios are identified and analysed, extracting common ideas and potential improvements.

2.1. Continuous authentication in single-device scenarios

In the literature, we can find the next two families of single-device scenarios: mobile devices and personal computers.

In the field of continuous authentication for mobile devices, Jorquera Valero et al. (Jorquera Valero et al., 2018) proposed a continuous and adaptive authentication system

Table 1
Continuous authentication solutions comparison.

Proposal	Device Type	Dimensions	Algorithms	Results / Conclusions
(Jorquera Valero et al., 2018)	Mobile	Sensors and Application Usage Statistics	Isolation Forest	Adaptability, low resource consumption, 92% Recall and 77% Precision, and resilience to adversarial attacks
(Bo et al., 2013)	Mobile	Sensors and touchscreen events	One Class - SVM SVM	72.36% Accuracy and 24.99% FAR
(Patel et al., 2016)	Mobile	Survey	-	Review about current state and challenges
(Li et al., 2018)	Mobile	Sensors with data augmentation	One-Class SVM	7.65% FAR, 9.01% FRR and 8.33% EER.
(Fridman et al., 2016)	Mobile	Text, location, application usage and websites	SVM	95% Accuracy and 5% EER
(Centeno et al., 2017)	Mobile	Sensors	Autoencoder	97.8% Accuracy and 2.2% EER
(Ehatisham-ul Haq et al., 2017)	Mobile	Sensors	Bayes Net and Euclidean distance	87.34–90.78% Accuracy
(Deutschmann and Lindholm, 2013)	Desktop	Mouse, Keyboard and Used applications	Bayes Net	18 seconds to detect an imposter using 15-50 keystrokes and 2.4 mins using 66 mouse interactions
(Fridman et al., 2015)	Desktop	Keyboard and mouse events linked to application.	Naive Bayes and SVM	0.4% FAR and 1% FRR after 30 s 0.1% FAR and 0.2% FRR after 5 minutes
(Aljohani et al., 2018)	Desktop	Keyboard and mouse	Artificial Immune System (AIS)	97.05% average Accuracy (96.6% to 97.74%)
(Mondal and Bours, 2017)	Desktop	Keyboard and mouse	Decision Tree, N. Network, SVM	Many experiments and test performed. 0.04–1% EER
(Lu et al., 2020)	Desktop	Keyboard	CNN, RNN	2.07% and 6.61% FAR, 3.26% and 5.31% FRR, and 2.67% and 5.97% EER , for CNN and RNN, respectively
(Sánchez Sánchez et al., 2019)	IoT	Mouse and keyboard (PC), Application usage (Mobile)	Random Forest	97.43% precision, 96.20% recall, 96.76% F1-Score (Mobile), 96.32 % precision, 90.00% recall, 92.70% F1-Score (PC)
(Ashibani et al., 2019)	Smart Home	User, Device, Network and Environment context	Heuristic analysis and pattern matching	<100 ms authentication time and verified improvement over using only credentials
(Nespoli et al., 2019)	IoT	Location, Person and IoT Devices Ontologies	Ontologies and semantic rules	Execution time less than 4s and more than 78% mean confidence level
AuthCode (Ours)	IoT	Modular architecture (PC and mobile)	MLP, XGBoost, RF and LSTM	Precision: 99.32%, Recall: 99.33%, F1-Score: 99.33% (more results in experiments section)

based on monitoring the application usage statistics and device sensors (gyroscope and accelerometer). The authors used ML-based anomaly detection techniques, concretely Isolation Forest, to identify anomalies in the users' behavioural data. Each user dataset is dynamically updated using his/her current behaviour in order to achieve system temporal adaptability. The solution obtained 92% recall and 77% precision when authenticating different users. In (Bo et al., 2013), Bo et al. identified users thanks to biometrics and typing patterns. The proposed system considered rotation, vibration, and pressure of smartphone touchscreens and sensors.

One-class SVM were used over the user behaviour to classify his profile, achieving 72.36% accuracy and 24.99% FAR (False Acceptance Rate). In (Patel et al., 2016), Patel et al. performed a review about the definition of continuous authentication and the proposals in the field. Authors focused on dimensions and AI techniques that are applied commonly in continuous authentication. They mentioned facial recognition, gestures, application usage and location, and concluded that merging data from different dimensions results in better accuracy and lower error rates. Li et al. (Li et al., 2018) proposed the application of data augmentation techniques

on behavioural information gathered from device sensors (gyroscope and accelerometer) to improve the continuous authentication results. Then, a one-class SVM model was trained and used to evaluate the user, obtaining 7.65% FAR (False Acceptance Rate), 9.01% FRR (False Rejection Rate) and 8.33% EER (Equal Error Rate). Other of the main works in this field is (Fridman et al., 2016), proposed by Fridman et al. This solution uses the typing stylometry, device location, application usage, and accessed websites. Authors achieved 0.4% FAR and 1% FRR after 30 seconds. Centeno et al. (Centeno et al., 2017) applied Autoencoders in their continuous authentication solution. This solution utilises sensor data to extract device holding patterns. Authors obtain better performance and resource consumption by using a cloud platform to perform the authentication process, achieving 2.2% EER. In (Ehatisham-ul Haq et al., 2017), authors utilised motion sensors in their continuous authentication system. Using these data, several device spatial positions were determined. After performing diverse tests, authors found that Bayes Net was the most appropriate algorithm for position recognition. Then, Euclidean distance was used to evaluate a instance compared to its recognised position, obtaining 87.43%-90.78% accuracy (depending on the evaluated position).

In terms of continuous authentication focused on computers and desktop devices, Deutschmann et al. (Deutschmann and Lindholm, 2013) selected keyboard, mouse and usage of applications as representative sources to identify users. Then, classification algorithms were used to sort and filter the gathered information in different categories. The results showed that intruders were detected in 2.4 minutes using the mouse, in 18 seconds using the keyboard, and 1.5 seconds using applications. Lex Fridman et al. (Fridman et al., 2015) utilised keyboard and mouse interactions to identify users. This work proposes to link the keyboard interactions and the application running on the foreground in order to obtain additional information that enables the authentication process. The system used Naive Bayes as classifier for mouse and keyboard events, and SVM for user typing patterns. Using short user interaction periods (30 secs), the system obtained a False Acceptance Rate (FAR) of 0.4% and a False Rejection Rate (FRR) of 1%. This metrics decreased to 0.1% and 0.2% respectively after a 5 minute evaluation. The system proposed by Aljohani et al. in (Aljohani et al., 2018) used an Artificial Immune System (AIS) to perform the continuous authentication task. They used the AIS Negative Selection (NS) algorithm on keyboard and mouse data. After a initialisation period, the system used AIS NS to evaluate keyboard and mouse interaction sets. This solution was evaluated in a group of 24 people, achieving 97.05% average accuracy (from 97.74% to 96.6% in all users). Mondal and Bours (Mondal and Bours, 2017) used keystroke and mouse movement dynamics to build a trust model. Then, this trust model is utilised to perform user continuous authentication. The behaviour of 53 different users was monitored in uncontrolled conditions and then the system was tested using the obtained information. Different Machine Learning classification algorithms were utilised to build a trust model. This trust model utilises a threshold over the dynamic user

authentication score. Last, Lu et al. (Lu et al., 2020) applied Convolutional Neural Networks (CNN) and Recurrent Neural Networks (RNN) to keyboard interactions in computers, achieving 2.07% and 6.61% FAR, 3.26% and 5.31% FRR, and 2.67% and 5.97% EER with CNN and RNN, respectively.

2.2. Continuous authentication in multi-device scenarios

This section analyses continuous authentication research works applied to multi-device scenarios. These solutions are the closest to the scope of this work, but the number of existing works is not very high due to the field novelty.

Sánchez et al. (Sánchez Sánchez et al., 2019) designed a continuous authentication architecture oriented to smart offices. Authors deployed their architecture as a proof of concept on mobile phones and computer devices. Then, the architecture was tested separately on the different devices using Random Forest (RF). In mobile, it achieved 97.43% average precision, 96.20% average recall and 96.76% average f1-score. In computer, it achieved 96.32% average precision, 90.00% average recall and 92.70% average f1-score. However, the authors did not perform any experiment combining the users' behaviour in different devices. Ashibani et al. (Ashibani et al., 2019) proposed a continuous authentication framework designed for Smart Homes. The framework employed contextual information to authenticate users. It authenticates by considering the user's context, the device context, the network context, and the environmental context, obtained from the smart home IoT devices. Nevertheless, unlike our work, this proposal does not consider user behaviour involving several devices. Another multi-device solution was proposed by Nespoli et al. in (Nespoli et al., 2019). This work was based on the application of semantic web techniques such as ontologies and rules for authentication and authorisation purposes in IoT. Specifically, IoT devices were used to gather information about the environment status. This information was used to perform the user's modelling and let him/her utilise certain services. Authors implemented the architecture and evaluated the resource consumption, scalability and authentication process confidence. Results reached an authentication confidence average of 78%. However, as authors claimed, the system performance is highly related to the deployment context. Then, vast training is necessary to accurately model the users' behaviour.

In conclusion, the existing continuous authentication proposals considering single- and multi-device scenarios do not combine the behaviour of users with different devices to infer additional behavioural patterns, as our work does. Table 1 provides a global overview of the key aspects considered by each solution. As it can be appreciated, previous continuous authentication solutions obtained good performance during the authentication process. However, the ones achieving better performance are only focused on single-device scenarios. Therefore, our proposal goes beyond the state-of-the-art and improves some of the limitation found in current solutions.

3. Motivating example: Multi-device profiling

This section seeks to provide an illustrative example of the deficiencies of single-device continuous authentication as well as to motivate how multi-device systems could reduce or even mitigate these deficiencies.

Let's suppose a scenario based on a Smart Office where two employees, Bob and Alice, utilize their computers and smartphones to perform their daily tasks. Each device hosts a continuous authentication application that monitors and permanently authenticates the owner according to his/her behaviour. This behaviour is determined by the usage of the screen, sensors, and applications in the mobile device. Regarding PCs, other dimensions such as the keyboard, mouse, resource usage and application monitoring can be considered.

In such Smart Office, Alice has been observing Bob when he used his smartphone, and she learned how to replicate his behaviour to impersonate Bob. This impersonation is achieved by repeating characteristics of Bob's activities, such as the way he holds the device, the speed at which he interacts with the screen and writes, or applications used and the application change patterns. One day, while Bob is working on his computer, Alice takes advantage of Bob's distraction and gets his smartphone to access and obtain sensitive information. Since Alice replicated Bob's behaviour, the single-device continuous authentication application cannot detect that Alice is using the device, instead of Bob. In this scenario, only a multi-device continuous authentication system could detect that Bob's smartphone and computer are used exactly at the same time, something that is not typical in Bob's behaviour while working.

Now, let's suppose that Bob installed and started using two new applications on his smartphone. It implies a smooth change in his behaviour, affecting the authentication accuracy of the single-device continuous authentication application. After that, the authentication application sometimes fails and does not detect Bob's behaviour as the owner (which means a high rate of false positives). This annoying situation can be reduced or mitigated by having a multi-device continuous authentication system that combines the behaviours of Bob with their devices in different time windows.

Finally, a multi-device continuous authentication system can also detect suspicious behaviours in one device compared to others' behaviour. As an example, multi-device could identify that recreational apps are used in the smartphone while the computer is being used for work. These situations are not detectable when single-device behaviour profiles are considered.

In conclusion, single-device continuous authentication solutions are capable of recognizing the device owner based on its behaviour. However, these solutions are sensitive to false-positive rates when small behaviour changes happen and can suffer impersonation attacks if the attacker knows the device owner's normal behaviour. These drawbacks can be improved by a multi-device authentication system able to combine information from several devices. This combination would improve the authentication robustness and flexibility, since an attacker would need to mimic the user in all the

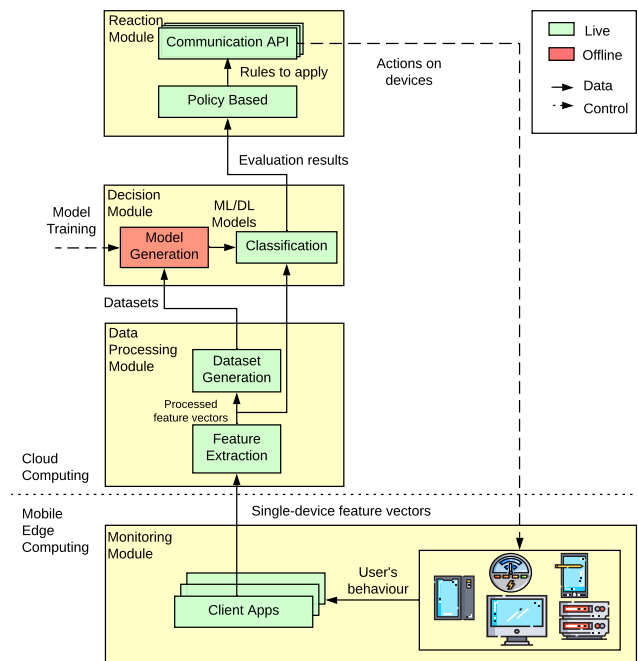


Figure 1: Design of the modules and components making up the AuthCODE architecture.

devices to attack the system integrity.

4. AuthCODE architectural design

This section describes the design details of our AI-based, privacy-preserving and multi-device continuous authentication architecture called AuthCODE. The architecture has been designed following a modular approach to allow flexible and dynamic modifications of its modules. In addition, AuthCODE combines the MEC, where users' sensitive data is managed and features are calculated to protect users' privacy, with the cloud computing, where computationally complex data processing and AI-base techniques are executed. This is a key characteristic, since sensitive data is maintained in the users' devices, and the performance of resource-constrained devices is not affected because AI-based models are trained and evaluated in the cloud. Figure 1 shows modules and components making up the AuthCODE architecture.

The AuthCODE architecture is composed of the following four modules:

- **Reaction.** Provides users with a final authentication score according to the outputs of the Decision module. Furthermore, it exposes interfaces with heterogeneous devices that enable global continuous authentication mechanisms for multi-device scenarios.
- **Decision.** Hosts and executes AI-based techniques able to train and evaluate different models that will authenticate users based on their behaviour with multiple devices.
- **Data Processing.** Filters, aggregates, and processes individual features acquired by the Monitoring module

to generate relevant combined feature vectors making up the single- and multi-device behavioural datasets.

- **Monitoring.** Monitors the data generated by users interacting with their heterogeneous devices and calculates single-device vectors of features that do not contain sensitive data. Once the feature vectors are calculated they are sent to the Data Processing module for further processing.

From bottom to up, *Monitoring* is the lowest module of AuthCODE and it is composed of several applications (Client Apps). Each client app runs on top of a device and monitors the data generated due to the user's actions. To ensure the privacy of users' sensitive data, each app processes and aggregates sensitive data in different windows of time established by the administrator, preventing the user's actions from being reconstructed or sensitive information from being inferred. After the aggregation, the app generates single-device vectors of features that do not contain sensitive data and send them periodically to the Data Processing module, which runs in the cloud computing, as explained in Section 5.

The *Data Processing* module periodically receives single-device feature vectors from each app. The Feature extraction component filters, aggregates and processes single-device feature vectors to calculate processed relevant features (which could belong to one or more devices depending on the scenario) that model the user's behaviour. The aggregation process of this module is also performed periodically and considering different time windows, which are established by the administrator as well. Finally, the Data Processing module generates datasets modelling the user's behaviour with different devices. These datasets can contain single-device or multi-device features, depending on the scenario needs.

The *Decision* module has two main tasks:

- **Off-line training.** The datasets generated by the previous module are used to train a set of models by using different ML and DL algorithms. The scenario requirement will decide if it is needed one model per device, one multi-device model, or both. The training process is performed by the *Model Generation* component only once and during the system bootstrapping.
- **Real-time evaluations.** Periodically and once the models have been trained, the *Classification* component evaluates the real-time features vectors against the models to provide an authentication score per model.

To conclude, the *Reaction* module aggregates the different authentication scores and calculates a global one as well as provides interfaces with the devices enabling a global and non-invasive multi-device continuous authentication. For that, the *Policy-based* component considers rules that establish the user's authentication level. This level decides proper security actions such as unlock a particular device without requiring additional credentials, lock the device, or ask for authentication credentials. Strict rules can be applied over

devices managing sensitive data or performing critical tasks, while more permissive rules can be applied over devices with secondary roles. Finally, the *Communication API* sends the previous security reactions to the different devices as well as configures the time windows sent to the Monitoring module to aggregate sensitive data and calculate features.

5. AuthCODE deployment & Datasets

This section shows the implementation details of the AuthCODE architecture as well as the generation of our five datasets in a realistic multi-device scenario such as a Smart Office. In our Smart Office, five users interacted with their smartphones, tablets, laptops and desktop computers for 60 days. Below we provide the implementation details of the modules making up our architecture.

5.1. Mobile Edge Computing

The Monitoring module and its Client Apps have been deployed close to the end users, in the MEC. They are hosted by Smart Office devices or by third-parties, in case of resource constrained devices. This decision ensures the performance and privacy-preserving capabilities of AuthCODE.

5.1.1. Personal computer devices

We have implemented a client app for Windows, the most used desktop OS, and another app for Linux distributions based on Debian. The apps monitor the following dimensions: 1) mouse movements/events, 2) keyboard events and 3) applications/resources usage statistics. Table 2 shows the selected dimensions and the data acquired from each dimension. We used Python, specifically, the pyinstaller tool to generate both Windows and Linux executable, the pynput library to monitor the mouse and keyboard events, and psutil and pywin32 (only in Windows) to monitor the resources consumption and applications usage (Team; Foundation). Finally, a time window is used to aggregate data and calculate features, which are sent to the Data Processing module using a REST API. In our implementation, the time window is set to 60 seconds. The feature selection process is explained and justified in Section 6.

5.1.2. Mobile devices

We have also deployed a client app for smartphones and tablets running from Android 5.0 OS (API 21). We chose Android since it is the most used operating system on mobile devices. In this case, the monitored dimensions are 1) the application usage statistics, and 2) the device sensors. Table 3 shows the data extracted from the previous dimensions. We have used Android.app.usage class (Developers) to gather the application usage statistics and the Android.hardware.SensorEventListener interface to obtain the gyroscope and accelerometer sensors data. To maintain a low resource consumption, we implemented a short-time service, which is triggered cyclically through the Android.app.AlarmManager. As in the personal computers case, the previous data is aggregated in time windows, and features are calculated and send to the Data Processing module

Table 2

Privacy-preserving features obtained from personal computers.

Dimension	Features
Time (1 feature)	- Vector timestamp.
Keyboard (24 002 features)	- Keystroke & word counter. - Erasing keys percentage. - Pressed keys histogram - Average & standard deviation of time that keys are pressed/released. - Average & standard deviation of consecutive keystroke intervals . - Number of written words & length histogram. - Digraphs typed (two consecutive keys) & mean time to type the digraph.
Mouse (45 features)	- Clicking speed average & standard deviation per mouse button and left double clicking. - Average mouse speed per direction. - Movement length histogram.
Application and resource usage (17 features)	- ID of the last and penultimate application used. - Active application counter average. - Counter of application changes. - CPU/RAM usage average and standard deviation. - Bytes transmitted & received through the network interfaces.

through an REST API. In our implementation, the time window is set to 60 seconds. The feature selection process is explained in Section 6.

5.2. Cloud computing platform

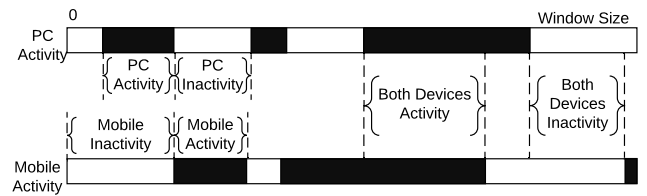
Due to storage and processing requirements, a private cloud platform hosts the Data Processing, Decision, and Re-action modules of AuthCODE. The Data Processing module exposes a REST API to receive single-device feature vectors from the Client apps. Periodically, AuthCODE follows the next steps to create our five datasets: 1) it processes and aggregates single-device feature vectors from the same user to generate multi-device feature vectors for every user using the user ID as label, 2) translates the features domains to make them suitable for ML/DL classifiers. Table 4 shows some of the most relevant multi-device features obtained after aggregating each device vector timestamps in a given time window. Figure 2 illustrates the four different activity combinations that arise when two devices are involved: none of the devices is active, only one of them is active, or both devices are active simultaneously. These features are tested and validated in Section 6. After following the previous steps AuthCODE generated the following five datasets (Sánchez Sánchez et al., 2020):

- *Dataset 1.* Single-device behaviour profile obtained from the personal computer, comprising aggregated

Table 3

Privacy-preserving features obtained from mobile devices.

Dimension	Features
Time (1 feature)	- Vector timestamp.
Application usage statistics (13 features)	- Foreground application counters (number of different and total apps) for the last minute and day. - Most common app ID and number of usages in the last minute and day. - ID of the currently active app - ID of the last active app prior the current one. - ID of the application most frequently utilised prior to the current application. - Bytes transmitted & received through the network interfaces.
Sensors (Gyroscope and Accelerometer) (40 features)	- Average, maximum, minimum, variance and peak-to-peak (max-min) of X, Y, Z coordinates. - $Magnitude = \sqrt{X^2 + Y^2 + Z^2}$

**Figure 2:** Explanatory diagram on the derived features from both user devices.

data about keyboard and mouse activity, as well as application usage statistics (see Table 2).

- *Dataset 2.* Single-device behaviour profile obtained from the mobile device, with application usage statistics (see Table 3).
- *Dataset 3.* Single-device behaviour profile with features computed from the sensors of the mobile device (see Table 3).
- *Dataset 4.* Multi-device behaviour profile combining the most relevant features of the mobile device and personal computer.
- *Dataset 5.* Multi-device behaviour profile generated from the active/inactive intervals of both devices (see Table 4).

Once the datasets are generated, the Decision module uses ML and DL algorithms to classify the users. Thus we had to implement a selected set of ML/DL models both for single- and multi-device profiles. To this end we used three well-known Python libraries: Scikit-learn, Keras and Pandas. Scikit-learn provides a wide variety of ML algorithms for both classification and anomaly detection. The Keras framework is widely used to implement and train DL models.

Table 4

Privacy-preserving multi-device features.

Features (32 in total)
Hour when the time window starts.
Weekday when the time window starts.
Total number of vectors from PC.
Total number of vectors from mobile devices.
Number of changes pc-mobile.
Number of changes mobile-pc.
Number of minutes with activity of both devices.
Mean, stdev, max, min of the PC Activity periods duration
Mean, stdev, max, min of the Mobile Activity periods duration
Mean, stdev, max, min of the Both Devices Activity periods duration
Mean, stdev, max, min of the PC Inactivity periods duration
Mean, stdev, max, min of the Mobile Inactivity periods duration
Mean, stdev, max, min of the Both Devices Inactivity periods duration

And finally, the Pandas library was employed to manipulate and process data. Finally, the Reaction module deploys the Policy-based component with management rules, which are implemented in Python. Moreover, the Communication API component implements a REST API to send the user's authentication result and the actions to be performed to the Smart Office devices.

6. Experiments

We measured the AuthCODE performance in terms of authentication accuracy and resource consumption. To accomplish this objective, the interactions of five individuals with their computers and mobile devices were collected for 60 days by means of the client apps detailed in Section 5.

6.1. Comparing single- and multi-device authentication.

This experiment analysed and compared the classification accuracy of AuthCODE considering single- and multiple-device behaviour profiles. Additionally, it justified the list of features selected for both single- and multi-device profiles (see Table 2, 3 and 4). To analyse the previous aspects we considered Dataset 1-4 of Section 5.

Given the huge number of features of some of the datasets (+24 000 features in Dataset 1) it was necessary a preliminary stage of feature selection. Our first step was to preprocess the three datasets discarding all the constant features and encoding each categorical one by using one-hot representation. Next, we chose RF and XGBoost (Chen and Guestrin, 2016) to perform an initial classification process of each single-device dataset for additional feature selection purposes. Both algorithms provide an estimation of the discriminative power of each feature. Additionally, they were chosen due to their ability to manage high number of features and their good performance dealing with imbalanced classes (some users have

Table 5

Classification algorithms and hyperparameters tested.

Algorithm	Hyperparameters
Naive Bayes	No hyperparameter tuning required
k-NN	$k \in [3, 20]$
SVM	$C \in [0.01, 100]$, $\gamma \in [0.001, 10]$ $kernel \in \{'rbf', 'linear', 'sigmoid', 'poly'\}$
XGBoost	$lr \in [0.01, 0.30]$, $max_depth \in [3, 15]$ $min_child_weight \in [1, 7]$, $\gamma \in [0, 0.5]$ $colsample_bytree \in [0.3, 0.7]$
MLP	$layers \in [1, 5]$, $neurons_layer \in [50, 1000]$
Random Forest	$Number_of_trees \in [50, 1000]$

more activity than others). Each dataset was partitioned in 10-minute segments and 10% segments were randomly chosen to create the test set. The previous and next segments for each selected segment were discarded to prevent data leakage in the training set due to correlation issues. Both algorithms were trained using each of the first three datasets, and the estimated discriminative of each feature was used to select a subset of features that comprised 95% of the total importance. Finally, Dataset 4 was created by combining these three resulting datasets to include multi-device information. The four final datasets were then used to train a set of candidate ML algorithms besides RF and XGBoost: Naive Bayes, k Nearest Neighbours (k-NN), Support Vector Machine (SVM) and Multi-Layer Perceptron (MLP). Each training was carried out using a validation set randomly selected similarly as the test set. The corresponding validation set of each dataset was used to perform proper optimization of hyperparameters for each ML algorithm. The list of hyperparameters per ML algorithm is detailed in Table 5. The performance metrics used to evaluate the models were the following (FPR: False Positive Rate, FRR: False Rejection Rate) :

$$Precision = \frac{TP}{TP + FP} \quad (1)$$

$$Recall = \frac{TP}{TP + FN} \quad (2)$$

$$F1 - Score = \frac{2 \times precision \times recall}{precision + recall} \quad (3)$$

$$FPR = \frac{FP}{FP + TN} \quad (4)$$

6.1.1. Single-device classification on personal computer data

The preliminary preprocessing of Dataset 1 reduced its number of features from +24 000 to 12 160. Additionally, the timestamp was replaced by just the time of day. With this

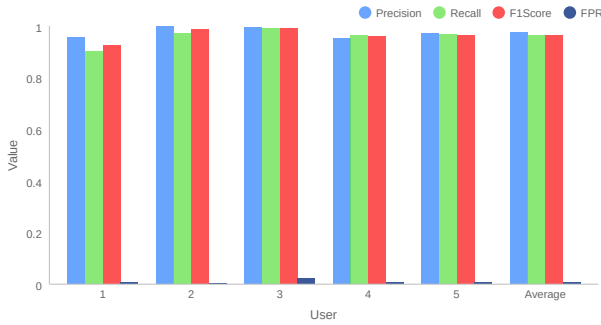


Figure 3: Classification performance achieved by MLP on Dataset 1 (PC data).

transformed dataset, XGBoost reached a slightly better performance than RF. Therefore, we used the features importance provided by XGBoost to selected the 150 most relevant ones, training with them every candidate classification model. This decision drastically reduced the dimensionality (from 12 160 to 150 features) for future testing without loss of classification performance.

As it can be seen in Table 6, MLP with one single hidden layer of 500 neurons obtained the best average classification performance, reaching a f1-score of 97.39%. Furthermore, Figure 3 shows the results of identifying each of the five users from their personal computer usage. Additionally, MLP reached for every single user a precision, recall, and f1-score higher than 95%, 90% and 95%, respectively. Regarding FPR, its values were smaller than 3% for every user.

6.1.2. Single-device classification on mobile device data

We evaluated the classification performance of the RF and XGBoost models on application statistics data (Dataset 2) and sensor data (Dataset 3) separately. RF achieved better classification performance than XGBoost on both datasets. Based on the discriminating importance of each feature provided by RF, the 50 most important features of Dataset 2 and the 40 most important features of Dataset 3 were selected, reducing the number of features from 818 to 50, and from 88 to 40, respectively.

Once selected the most discriminating features, Table 6 shows that RF, with *Number_of_trees*: 500, obtained the best performance with Dataset 2, and XGBoost, with *lr*: 0.25, *max_depth*: 10, *min_child_weight*: 3, *gamma*: 0.5 and *colsample_bytree*: 0.5, was the best model for Dataset 3. Additionally, Figure 4 shows the classification results for each user by considering the previous models for each dataset.

6.1.3. Multi-device classification

We used Dataset 4 to evaluate whether the combination of the most discriminating features of personal computers (150) and mobile devices (50 for apps statistics and 40 for sensors) could improve the authentication results obtained in the previous two experiments. With that goal in mind, the previous features were grouped in time windows of one minute. In this way, we created feature vectors representing

Table 6

Comparison of classification algorithms for single- and multiple-device behaviour profiles.

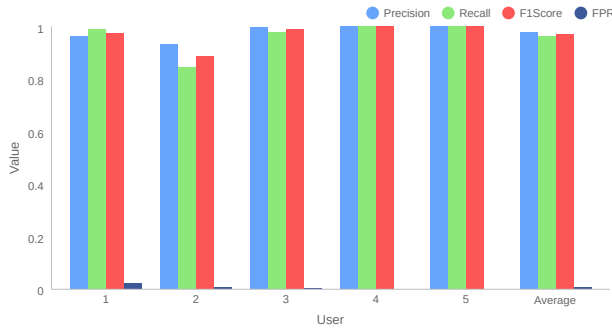
Model	Naive Bayes	K-NN	SVM	XG Boost	MLP	RF
Dataset 1: Personal computer						
Avg. Precis.	0.6977	0.9280	0.9593	0.9708	0.9752	0.9710
Avg. Recall	0.6425	0.9187	0.9532	0.9525	0.9727	0.9279
Avg. F1-Scr.	0.6054	0.9221	0.9561	0.9610	0.9739	0.9459
Avg. FPR	0.0951	0.0125	0.0092	0.0096	0.0072	0.0085
Dataset 2: Applications usage statistics						
Avg. Precis.	0.8281	0.9473	0.9318	0.9466	0.9505	0.9775
Avg. Recall	0.8005	0.9368	0.9057	0.9390	0.9487	0.9534
Avg. F1-Scr.	0.7286	0.9418	0.9174	0.9421	0.9504	0.9670
Avg. FPR	0.0523	0.0082	0.0112	0.0118	0.0085	0.0057
Dataset 3: Sensors						
Avg. Precis.	0.2915	0.7701	0.8682	0.9242	0.8439	0.9113
Avg. Recall	0.2751	0.7327	0.8171	0.8871	0.6502	0.8547
Avg. F1-Scr.	0.2360	0.7483	0.8385	0.9036	0.7328	0.8783
Avg. FPR	0.1933	0.0618	0.0355	0.0216	0.0882	0.0271
Dataset 4: Multi-device						
Avg. Precis.	0.7425	0.9342	0.9606	0.9932	0.9712	0.9799
Avg. Recall	0.7577	0.9362	0.9671	0.9933	0.9476	0.9694
Avg. F1-Scr.	0.6862	0.9351	0.9637	0.9933	0.9591	0.9743
Avg. FPR	0.0544	0.0115	0.0062	0.0023	0.0064	0.0045

the activities of each user interacting with the two devices in the same minute. Figure 5 depicts a scheme of the morphology of the vector generated. If any of the devices (PC or mobile) has no activity in that minute, their features are established to 0. Only vectors with activity are generated.

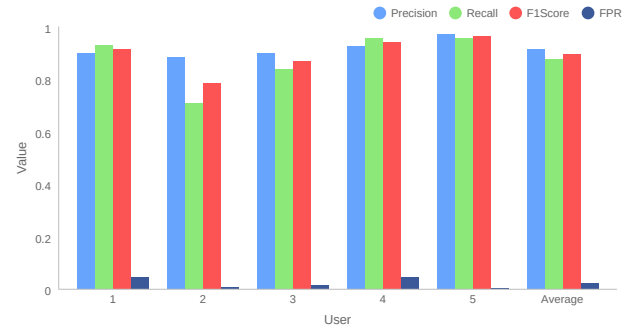
The relatively low number of features (240) and the better class balance due to the combination of the other three datasets, made the multi-device Dataset 4 suitable to be used without changes to train the set of ML algorithms, including their hyperparameter optimization. Table 6 shows the best classification results for each algorithm on Dataset 4. The winner was XGBoost, with the following hyperparameters: *lr*: 0.25, *max_depth*: 10, *min_child_weight*: 5, *gamma*: 0.1 and *colsample_bytree*: 0.7.

The results of the previous three subsections demonstrated how single-device classification results can be improved by combining single-device features to create multi-device profiles. Table 6 compares the results of each one of the experiments of this section, demonstrating that the best results were obtained with multi-device classification and XGBoost with the following hyperparameters: *lr*: 0.25, *max_depth*: 10, *min_child_weight*: 5, *gamma*: 0.1 and *colsample_bytree*: 0.7.

Analysing the previous results, two main conclusions are obtained. On the one hand, AuthCODE classifies and



(a) Classification performance achieved by RF on Dataset 2 (Application usage statistics).



(b) Classification performance achieved by XGBoost on Dataset 3 (Sensors data).

Figure 4: Users' behaviour classification in mobile devices.

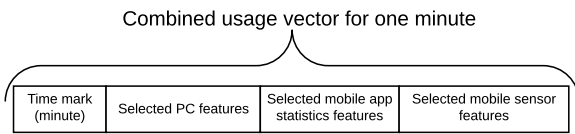


Figure 5: Structure of multi-device features vector making up Dataset 4.

authenticates users in single-device scenarios such as PC and mobile devices with 97.39%, 96.70% and 90.36% average f1-score and 0.72%, 0.57% and 2.16% average FPR for computer, mobile application statistics and mobile sensors, respectively. On the other hand, the proposed multi-device behaviour profile dataset improves the authentication results of single-device datasets, obtaining 99.33% average f1-score and 0.23% average FPR. It is interesting to note that the FPR with multi-device achieved a 69.33%, 59.65% and 89.35% improvement for computer, mobile applications and mobile sensors respectively.

6.2. Measuring the time impact in multi-device authentication.

In the previous section, the user identification improvement obtained by using multi-device behaviour profiles collected in 1-minute time windows was evaluated. In this experiment we analysed the impact of using a sequence of such 1-minute snapshot vectors, seen as the evolution of a user's activity over the time. The temporal information carried in the sequence can be represented in different ways. We studied two approaches: authentication by means of derived temporal activity features and by means of time series of vectors.

6.2.1. Derived features classification

In this experiment we evaluated whether it was possible to identify users according to the time they spent interacting with their devices. For this purpose, we labelled each 1-minute vector from Datasets 1-3 with its associated device (computer or mobile). Subsequently, the resulting vectors were sorted by their timestamps and grouped in a variety of time windows (1 hour, 3 hours, 6 hours, 12 hours and 24

Table 7

Derived usage features classification results using RF (Dataset 5).

Time window	1 h	3 h	6 h	12 h	24 h
Average Precision	0.7298	0.7893	0.8005	0.8224	0.9310
Average Recall	0.7291	0.7883	0.8005	0.8175	0.9253
Average F1-Score	0.7259	0.7884	0.7980	0.8193	0.9245
Average FPR	0.0621	0.0498	0.0480	0.0491	0.0215

hours). For each time windows size, a dataset was created with features about the usage periods of the devices present in the window and the device changes made by the user (see Table 4). Figure 5 illustrates a window containing a group of vectors belonging to two devices. It can be noticed how the user switches between the two devices and the periods where there is activity on just one device, both devices or even no activity at all.

Our set of ML algorithms were trained with this dataset, including a hyperparameter optimization procedure. RF was the model that achieved the best performance (number of trees = 200). Table 7 lists the results for each selected time windows. The classification results improve as window size increases, ranging from 72.59% f1-score and 6.21% FPR with 1-hour window to 92.45% f1-score and 2.15% FPR with 24-hour window. These results make sense, since the larger the window, the more interactions contains, providing more information to differentiate clearly each user's behaviour.

Although the performance achieved did not improve the results obtained in Section 6.1, this experiment demonstrated the potential of using derived features to identify different users based on their device usage routines.

6.2.2. Time window processing using LSTM

In this experiment we leveraged the ability of DL, and more specifically Long-Short Term Memory Neural Net-

works (LSTM), to learn complex patterns in sequences of vectors obtained from the activity data from the users' devices. LSTM is effective at capturing long-term temporal dependencies; therefore, our aim was to evaluate whether it was able to use the temporal information to improve the level of authentication reached in the experiments of Section 6.1. In those experiments we created Dataset 4 by aggregating 1-minute time windows of multi-device activity data in vectors of 240 features, obtaining an average f1-score of 99.33% and FPR of 0.23% with XGBoost. The same Dataset 4 was subsequently processed to be used as input of a selection of LSTM architectures in a variety of configurations. The steps we carried out were the following:

- Dataset 4 includes 60 days of monitoring data but a some of the first and last days do not contain data from all users. Therefore, we selected a continuous period of 40 days in which the five users had activity data.
- The feature vectors were sorted by the minute in which they were generated. To make explicit the lack of activity in the sequence, every minute without activity had -1 in the rest of the vector features.
- We used the numpy to work with array views and generate a dataset of sequences of a given size. The numpy function `stride_tricks.as_strided` is particularly useful, allowing us to obtain a sliding window view of our dataset. With this function we obtained 10 separate views of the dataset for a range of sliding window sizes (2, 5, 10, 20, 30, 60, 90, 120, 240 and 360 minutes).

Our main architecture was composed of an optional 1D-convolutional layer (Conv1D), as suggested by some works in the literature (Eapen et al., 2019; Zhong et al., 2019), followed by a number of LSTM layers (from 1 to 4) with different number of nodes (from 16 to 256) and a 5-node fully connected softmax output layer. All these configuration values together with the use of batch normalization and the dropout ratio in each layer, were hyperparameters to be optimized. Every dataset view was split into training/validation/test subsets. Each view was used as input during the hyperparameter optimization procedure. The validation and test datasets were composed of the sequences belonging to 7 days each: days 26 to 33 and 34 to 40 respectively. The model configuration that achieved the higher performance was a 2-layer LSTM with 64 and 32 nodes, without Conv1D layer, no batch normalization and 20% dropout. Figure 6 shows a scheme of this network.

Table 8 lists the average precision, recall and f1-score when classifying behaviour data using different time window sizes. The table shows how as the time window size was increased, in general, the classification results were better. The lowest average f1-score and FPR (2 minute window) were 82.38% and 3.81% respectively, and they improved with window size until they reached 99% and 0.37% when the time window is 60 minutes. For greater windows the metrics stabilised in similar values, even when the window is six times larger (360 minutes).

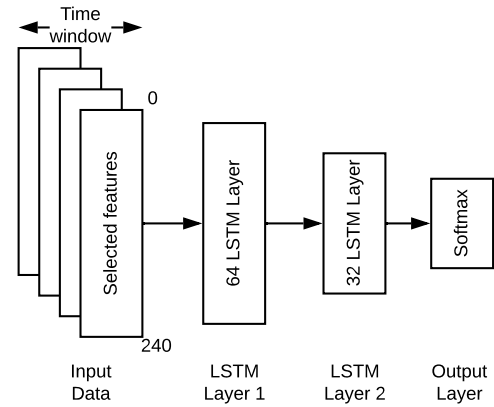


Figure 6: LSTM network architecture.

From these results some conclusions can be drawn. As it can be noticed in Table 7 and Table 8, LSTM-based classification using Dataset 4 view as a dataset of sufficient large vector sequences, improved the results obtained when using Dataset 5 and RF. Also in Table 8, we can observe how performance improves as the time window increases. After the 60 minute window, the results stabilised around 99% for f1-score and 0.30% for FPR. So, it can be a good trade of between window size and performance. These results can be applied as a complement to the one-minute multi-device vector classification, giving an additional temporal information to the authentication process.

6.3. Resource Consumption

This experiment measured the resource consumption of client apps, running on smartphones and computers, as well as the ML/DL models of AuthCODE running on our server. For each testing device (see Table 9), relevant resources such as battery, memory, storage and processing, were monitored to check the impact of AuthCODE. Resources whose consumption is dynamic such as battery, memory, or CPU were monitored during 10 days, averaging the measurements of that period to calculate final results. Resources whose consumption is static, such as storage, were measured once the AuthCODE components were deployed.

6.3.1. Client Apps Consumption

Battery, memory, storage and processing are the most critical resources of constrained resource devices such as laptops and smartphones. This experiment aimed to measure the impact of our client apps in the hosting devices resources.

- **Battery.** In average, the client app consumed 167 mAh ($\approx 4\%$) and 201 mAh ($\approx 6\%$) of the Xiaomi and Huawei battery devices, respectively. In terms of laptops, ≈ 4.75 mAh ($< 1\%$) and ≈ 6.54 mAh ($< 1\%$) of the HP and Acer batteries respectively were consumed.
- **Memory.** In both laptops (HP and Acer), the client app consumed ≈ 25 MB. In contrast, for both mobile devices (Xiaomi and Huawei), the client app used ≈ 104 MB.

Table 8

Temporal windows classification results using LSTM network (Dataset 6).

Time Window (minutes)	2	5	10	20	30	60	90	120	180	360
Average Precision	0.9070	0.9532	0.9738	0.9794	0.9883	0.9971	0.9948	0.9984	0.9989	0.9988
Average Recall	0.7584	0.8380	0.8453	0.9390	0.9538	0.9840	0.9777	0.9890	0.9932	0.9907
Average F1-Score	0.8238	0.8898	0.9029	0.9579	0.9698	0.9902	0.9856	0.9934	0.9959	0.9945
Average FPR	0.0381	0.0202	0.0199	0.0086	0.0065	0.0037	0.0038	0.0034	0.0028	0.0026

Table 9

Device specification of resource consumption tests.

Device	Processor	Mem. (GB)	Battery (mAh)
Laptop: HP 15-bs00x	Intel i7-7500U 4 Cores @ 2.7 GHz	8	2850
Laptop: Acer Nitro 5 AN51	Intel i7-7700HQ 4 Cores @ 2.8 GHz	8	3270
Smartphone: Xiaomi Pocophone	Snapdragon 845 8 Cores @ 2.8 GHz	6	4000
Smartphone: Huawei P10 Lite	Kirin 658 8 Cores @ 1.7 GHz	4	3000
Cloud Server	Intel Xeon E5-2697v4 18 Cores @ 2.3 GHz	64	-

- *Storage.* The client app executable file occupied 6.50 MB in the personal computers, and 15 MB in the smartphones (regardless the device model). The amount of raw data needed to compute a single-device feature vector is temporarily stored in the devices, having a size of ≈ 50 KB. Once the single-device feature vector is sent to the server, this storage is released.
- *Processing.* In both laptops (HP and Acer), the average daily CPU usage ranged between 2% and 5%. In both smartphones (Xiaomi and Huawei), daily CPU usage remained under 1%.

Based on the previous results, it can be concluded that neither computer client nor mobile client apps have a significant impact on the device resource consumption. Therefore, our client apps are suitable even for resource constrained devices.

6.3.2. ML/DL Model Consumption

This experiment measured the resource and time consumption of the AuthCODE modules deployed in our server to train and evaluate our ML and DL models.

- *Time.* The average time needed to process each feature vector and evaluate it varied from 1.5 to 2 seconds. The measured process included the filtering and selection of features, the grouping of vectors in a 60-minute window, and the vector evaluation using XGBoost and the 60-minute LSTM network.

- *Memory.* Once trained the models, they were loaded in memory and utilised to evaluate live users' vectors in real time. The memory usage of XGBoost and LSTM was 4.75 MB and 49.50 MB.
- *Storage.* The Python scripts implementing Data Processing, Decision and Reaction modules of AuthCODE had a size of ≈ 6 KB. The generated models had a size of 2.05 MB for the XGBoost classifier (.pickle format) and 1.1MB for the LSTM models (.h5 format). In total, 3.15 MB of the server storage were used.
- *Processing.* When the Data Processing module received features to evaluate the model and authenticate users, AuthCODE utilised $\approx 3\%$ in average.

In conclusion, we have proven that the AuthCODE modules makes efficient use of the processing, memory, storage, and battery resources of heterogeneous devices such as mobile devices, computers, and servers. In addition we have shown that the time required to evaluate and authenticate a given user remains in about 2 seconds, which is acceptable for our continuous authentication scenario prototype. Finally, it is important to keep in mind that the resource consumption optimization task is not the main objective of this work.

7. Conclusions and Future Work

This paper presents AuthCODE, a multi-device continuous authentication architecture, deployed in the MEC and cloud infrastructures, that utilises ML and DL techniques to authenticate users according to their behaviour. AuthCODE proposes a list of privacy-preserving and multi-device features that combine the user's interactions with different devices. The relevance of the previous features as well as the improvement of multi-device profiles compared to single-ones have been validated in a Smart Office scenario, where we generated and published five behavioural datasets. Several experiments over the previous datasets demonstrated that multi-device profiles improved the authentication accuracy and FPR of solutions based on single-device profiles, reaching an f1-score of 99.33% and a FPR of 0.23% with XGBoost. Therefore, our multi-device approach improves by 69.33%, 59.65% and 89.35% the FPR obtained by computer, mobile applications and mobile sensors separately.

Additionally, the inclusion of temporal information in the form of vector sequences provides a further improvement in the authentication performance of the single-vector models, allowing the identification of complex behaviour patterns associated to each user. With this approach, an LSTM achieved

an f1-score of 99.02% and a FPR of 0.37% with a 60-minute sequence of vectors. To conclude, several experiments also demonstrated the suitability of the proposed solution in terms of resource consumption at the mobile edge and cloud computing paradigms.

However, some limitations of the current proposal should be also commented. First, as user's behaviour is evaluated every minute, there is a short period of time in which an attacker could make use of the device without being detected. Besides, as ML/DL classification algorithms have been applied, user identification performance depends on how different each user's behaviour is from other users, and there can be model scalability problems if the number of users increases too much. Then, anomaly detection algorithms will be tested in the future. Finally, processing and evaluation is carried out in a centralized server, which can lead to scalability, availability or security issues. These limitations motivate to continue researching in multi-device continuous authentication solutions.

As future work, we plan to evaluate the authentication accuracy of AuthCODE with more users and new experiments aiming to detect common behaviours by using new ML/DL algorithms and filtering actions per type of application. Furthermore, we will extend the use case implementation by considering IoT devices, other operating systems, and new dimensions such as writing patterns or network traffic statistics. It will allow us to generate and release novel datasets, useful for researchers and scientist to keep improving the multi-device continuous authentication challenge.

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.

CRedit authorship contribution statement

Pedro Miguel Sánchez Sánchez. Methodology, Writing - original draft, Data curation, Software. **Lorenzo Fernández Maimó.** Methodology, Writing - Review & Editing. Formal analysis. **Alberto Huertas Celdrán.** Methodology, Conceptualization, Writing - Review & Editing. **Gregorio Martínez Pérez:** Supervision, Project administration, Funding acquisition.

Acknowledgment

This work has been partially supported by Armasuisse S+T with project CYD-C-2020003, by the University of Zürich UZH, and by the European Union Horizon 2020 Research and Innovation Program under grant agreement No. 830927, namely the H2020 Concordia Project. Special thanks to all those voluntaries who installed the client applications: Oscar Fernández, Pedro A. Sánchez, Francisco J. Sánchez, Pantaleone Nespole, Mattia Zago, Sergio López, Eduardo López, Manuel Gil, José M. Jorquera, Javier Pastor and Gregorio Martínez.

References

- Aljohani, O., Aljohani, N., Bours, P., Alsolami, F., 2018. Continuous authentication on pcs using artificial immune system, in: 2018 1st International Conference on Computer Applications & Information Security (ICCAIS), IEEE. pp. 1–6.
- Almalki, S., Chatterjee, P., Roy, K., 2019. Continuous authentication using mouse clickstream data analysis, in: International Conference on Security, Privacy and Anonymity in Computation, Communication and Storage, Springer. pp. 76–85.
- Ashibani, Y., Kauling, D., Mahmoud, Q.H., 2019. Design and implementation of a contextual-based continuous authentication framework for smart homes. *Applied System Innovation* 2, 4.
- Bo, C., Zhang, L., Li, X.Y., Huang, Q., Wang, Y., 2013. Silentsense: silent user identification via touch and movement behavioral biometrics, in: Proceedings of the 19th annual international conference on Mobile computing & networking, pp. 187–190.
- Centeno, M.P., van Moorsel, A., Castruccio, S., 2017. Smartphone continuous authentication using deep learning autoencoders, in: 2017 15th Annual Conference on Privacy, Security and Trust (PST), IEEE. pp. 147–1478.
- Chen, T., Guestrin, C., 2016. Xgboost: A scalable tree boosting system, in: Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining, pp. 785–794.
- Deutschmann, I., Lindholm, J., 2013. Behavioral biometrics for darpa's active authentication program, in: 2013 International Conference of the BIOSIG Special Interest Group (BIOSIG), IEEE. pp. 1–8.
- Developers, A., . Android library. URL: <https://developer.android.com/>. accessed on: April 15, 2020.
- Eapen, J., Bein, D., Verma, A., 2019. Novel deep learning model with cnn and bi-directional lstm for improved stock market index prediction, in: 2019 IEEE 9th annual computing and communication workshop and conference (CCWC), IEEE. pp. 0264–0270.
- Firouzi, F., Farahani, B., 2020. Architecting IoT Cloud. Springer International Publishing. doi:10.1007/978-3-030-30367-9_4.
- Fridman, L., Stolerman, A., Acharya, S., Brennan, P., Juola, P., Greenstadt, R., Kam, M., Gomez, F., 2015. Multi-modal decision fusion for continuous authentication. *Computers and Electrical Engineering* 41, 142–156. doi:10.1016/j.compeleceng.2014.10.018.
- Fridman, L., Weber, S., Greenstadt, R., Kam, M., 2016. Active authentication on mobile devices via stylometry, application usage, web browsing, and gps location. *IEEE Systems Journal* 11, 513–521.
- Fundation, P.S., . The python package index (pypi). URL: <https://pypi.org/>. accessed on: April 15, 2020.
- Gonzalez-Manzano, L., Fuentes, J.M.D., Ribagorda, A., 2019. Leveraging user-related internet of things for continuous authentication: A survey. *ACM Computing Surveys (CSUR)* 52. doi:10.1145/3314023.
- Ehatisham-ul Haq, M., Azam, M.A., Loo, J., Shuang, K., Islam, S., Naeem, U., Amin, Y., 2017. Authentication of smartphone users based on activity recognition and mobile sensing. *Sensors* 17, 2043.
- Huertas Celdrán, A., Gil Pérez, M., García Clemente, F.J., Martínez Pérez, G., 2019. Towards the autonomous provision of self-protection capabilities in 5g networks. *Journal of Ambient Intelligence and Humanized Computing* 10, 4707–4720. doi:10.1007/s12652-018-0848-6.
- Jorquera Valero, J.M., Sánchez Sánchez, P.M., Fernández Maimó, L., Huertas Celdrán, A., Arjona Fernández, M., De Los Santos Vélchez, S., Martínez Pérez, G., 2018. Improving the security and qoe in mobile devices through an intelligent and adaptive continuous authentication system. *Sensors* 18, 3769.
- Li, Y., Hu, H., Zhou, G., 2018. Using data augmentation in continuous authentication on smartphones. *IEEE Internet of Things Journal* 6, 628–640.
- Lu, X., Zhang, S., Hui, P., Lio, P., 2020. Continuous authentication by free-text keystroke based on cnn and rnn. *Computers & Security* 96, 101861. URL: <http://www.sciencedirect.com/science/article/pii/S0167404820301334>, doi:https://doi.org/10.1016/j.cose.2020.101861.
- Mondal, S., Bours, P., 2017. A study on continuous authentication using a combination of keystroke and mouse biometrics. *Neurocomputing* 230, 1–22.

- Nespoli, P., Zago, M., Huertas Celdrán, A., Gil Pérez, M., Gómez Mármol, F., Clemente, G., Félix, J., 2019. Palot: profiling and authenticating users leveraging internet of things. *Sensors* 19, 2832.
- Patel, V.M., Chellappa, R., Chandra, D., Barbelo, B., 2016. Continuous user authentication on mobile devices: Recent progress and remaining challenges. *IEEE Signal Processing Magazine* 33, 49–61.
- Qolomany, B., Al-Fuqaha, A., Gupta, A., Benhaddou, D., Alwajidi, S., Qadir, J., Fong, A.C., 2019. Leveraging machine learning and big data for smart buildings: A comprehensive survey. *IEEE Access* 7, 90316–90356.
- Sánchez Sánchez, P.M., Huertas Celdrán, A., Fernández Maimó, L., Pérez, G.M., Wang, G., 2019. Securing smart offices through an intelligent and multi-device continuous authentication system, in: *International Conference on Smart City and Informatization*, Springer. pp. 73–85.
- Sánchez Sánchez, P.M., Fernández Maimó, L., Huertas Celdrán, A., Martínez Pérez, G., 2020. Authcode - dataset. URL: <http://dx.doi.org/10.21227/ttcs-ak23>, doi:10.21227/ttcs-ak23.
- Team, P.D., . Pyinstaller. pyinstaller bundles python applications. URL: <https://www.pyinstaller.org/>. accessed on: April 15, 2020.
- Zhong, L., Hu, L., Zhou, H., 2019. Deep learning based multi-temporal crop classification. *Remote sensing of environment* 221, 430–443.



Pedro Miguel Sánchez Sánchez received the M.Sc. degree in computer science from the University of Murcia. He is currently pursuing his PhD in computer science at University of Murcia. His research interests are focused on continuous authentication, networks, 5G, cybersecurity and the application of machine learning and deep learning to the previous fields.



Lorenzo Fernández Maimó received the M.Sc. and Ph.D. degrees in computer science from the University of Murcia. He is currently an Associate Professor with the Department of Computer Engineering, University of Murcia. His research interests primarily focus on machine learning and deep learning applied to cybersecurity, and computer vision.



Alberto Huertas Celdrán received the M.Sc. and Ph.D. degrees in computer science from the University of Murcia, Spain. He is currently a postdoctoral fellow associated with the Communication Systems Group (CSG) at the University of Zurich UZH. His scientific interests include medical cyber-physical systems (MCPS), brain-computer interfaces (BCI), cybersecurity, data privacy, continuous authentication, semantic technology, context-aware systems, and computer networks.



Gregorio Martínez Pérez is Full Professor in the Department of Information and Communications Engineering of the University of Murcia, Spain. His scientific activity is mainly devoted to cybersecurity and networking, also working on the design and autonomic monitoring of real-time and critical applications and systems. He is working on different national (14 in the last decade) and European IST research projects (11 in the last decade) related to these topics, being Principal Investigator in most of them. He has published 160+ papers in national and international conference proceedings, magazines and journals.