



D2.1 – USE CASE SPECIFICATION

| | |
|------------------------|--|
| Work Package | WP 2, Use Case Specification |
| Lead Author | Marco Mosconi (MCI), Eleonora Ciceri (MCI), Stefano Galliani (MCI) |
| Contributing Author(s) | Monir Azraoui (ORA), Sébastien Canard (ORA), Dominique Le Hello (ORA), Angel Palomares Perez (ATOS), Melek Önen (EURC) |
| Reviewers | Melek Önen (EURC), Boris Rozenberg (IBM) |
| Due date | 30.04.2019 |
| Date | 30.04.2019 |
| Version | 1.0 |
| Dissemination Level | PU (Public) |



The research leading to these results has received funding from the European Union's Horizon 2020 Research and Innovation Programme, through the PAPAYA project, under Grant Agreement No. 786767. The content and results of this deliverable reflect the view of the consortium only. The Research Executive Agency is not responsible for any use that may be made of the information it contains.



Project No. 786767

D2.1 – Use Case Specification Dissemination Level – PU

Revision History

| Revision | Date | Editor | Notes |
|----------|------------|---|---|
| 0.1 | 21.01.2019 | Marco Mosconi (MCI) | TOC definition |
| 0.2 | 15.02.2019 | Eleonora Ciceri, Marco Mosconi, Stefano Galliani (MCI) | Contribution on healthcare use cases description |
| 0.3 | 15.03.2019 | Monir Azraoui, Sébastien Canard, Dominique Le Hello (ORA) | Contribution on Mobile and phone usage use case description |
| 0.4 | 02.04.2019 | Eleonora Ciceri (MCI), Melek Önen (EURC), Boris Rozenberg (IBM) | First revision |
| 0.5 | 10.04.2019 | Eleonora Ciceri, Marco Mosconi (MCI), Monir Azraoui, Sébastien Canard (ORA) | Updates after first revision |
| 0.6 | 18.04.2019 | Melek Önen (EURC), Boris Rozenberg (IBM) | Second revision |
| 0.7 | 23.04.2019 | Eleonora Ciceri, Marco Mosconi (MCI), Monir Azraoui, Sébastien Canard (ORA) | Updates after second revision |
| 0.8 | 24.04.2019 | Beyza Bozdemir (EURC) | Quality check and dedicated updates |
| 1.0 | 30.04.2019 | Melek Önen (EURC) | Final updates, ready for submission |



Project No. 786767

D2.1 – Use Case Specification
Dissemination Level – PU

Table of Contents

| | |
|---|----|
| Executive Summary | 5 |
| Glossary of Terms..... | 6 |
| 1 Introduction | 7 |
| 1.1 Purpose of the document..... | 7 |
| 1.2 Outline | 8 |
| 2 Healthcare use cases | 9 |
| 2.1 Single source architecture: Privacy-preserving arrhythmia detection | 9 |
| 2.1.1 Introduction and current solution limitations | 9 |
| 2.1.2 Use case definition | 10 |
| 2.1.3 Privacy requirements | 16 |
| 2.2 Multiple source architecture: Privacy-preserving stress management | 16 |
| 2.2.1 Introduction and current solution limitations | 16 |
| 2.2.2 Use case definition | 17 |
| 2.2.3 Privacy requirements | 24 |
| 3 Mobility and phone usage use cases | 26 |
| 3.1 Single source architecture: Privacy-preserving mobility analytics..... | 26 |
| 3.1.1 Introduction and current solution limitations | 26 |
| 3.1.2 Use case definition | 28 |
| 3.1.3 Privacy requirements | 33 |
| 3.2 Multiple source architecture: Privacy-preserving mobile usage analytics | 33 |
| 3.2.1 Introduction and current solution limitations | 33 |
| 3.2.2 Use case definition | 35 |
| 3.2.3 Privacy requirements | 40 |
| 4 Threat detection use case..... | 42 |
| 4.1 Introduction and current solution limitations | 42 |
| 4.2 Use case definition..... | 44 |
| 4.2.1 Story | 44 |
| 4.2.2 Players..... | 45 |



D2.1 – Use Case Specification Dissemination Level – PU

Project No. 786767

| | | |
|-------|-------------------------------------|----|
| 4.2.3 | Involved data | 46 |
| 4.2.4 | Processing and analytics | 46 |
| 4.2.5 | Protocol | 47 |
| 4.3 | Data sensitivity requirements | 49 |
| 5 | Conclusions | 50 |
| 6 | References | 51 |

List of Tables

| | | |
|----------|--|----|
| Table 1 | UC1 Players | 11 |
| Table 2 | Description of UC1 | 13 |
| Table 3 | Privacy Requirements for UC1 | 16 |
| Table 4 | UC2 Players | 19 |
| Table 5 | Description of UC2 - Dataset collection | 20 |
| Table 6 | Description of UC2 - Model training | 21 |
| Table 7 | Description of UC2 - Stress detection and classification | 22 |
| Table 8 | Privacy requirements for UC2 | 24 |
| Table 9 | UC3 Players | 28 |
| Table 10 | Description of UC3 - Audience measurements | 31 |
| Table 11 | Description of UC3 - Trajectories analysis | 32 |
| Table 12 | Privacy requirements for UC3 | 33 |
| Table 13 | UC4 Players | 36 |
| Table 14 | Description of UC4 | 38 |
| Table 15 | UC5 Players | 45 |
| Table 16 | Description of UC5 | 47 |

List of Figures

| | | |
|----------|----------------------------|----|
| Figure 1 | UC5 research phase | 44 |
| Figure 2 | UC5 commercial phase | 45 |



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

Executive Summary

This deliverable reports the work carried out for the execution of the task T2.1 (“*Use case definition*”). The document reports the use cases identified in the two high-level use case umbrellas, namely, the healthcare umbrella and the mobile and phone usage umbrella. In each of these umbrellas, we identified the following use cases, as follows:

- **Arrhythmia detection use case (healthcare umbrella).** In this use case, sensitive health data in the form of ECG are collected by a single source and an external entity that runs the PAPAYA platform helps detect arrhythmias from the ECG signal without having access to its content.
- **Stress detection use case (healthcare umbrella).** In this use case, sensitive health data from IoT sensors are collected by multiple sources and used to train a collaborative model via the PAPAYA platform, with the ultimate goal of automatically detecting stress conditions in workers.
- **Mobility analytics use case (mobile and phone usage umbrella).** In this use case, sensitive data are collected from a single source and an external entity that runs the PAPAYA platform measures the audience in one or several areas, or extracts mobility patterns.
- **Mobile usage analytics use case (mobile and phone usage umbrella).** In this use case, sensitive data on mobile phone application usage are collected from multiple sources and used to extract analytics for statistical purposes.
- **Threat detection use case (mobile and phone usage umbrella).** In this use case business-sensitive data (rather than personal data) from several sources are processed, to detect threats in systems or networks.

Each use case is accompanied with the related players’ description, its privacy requirements and its expected processing steps, and the description explains how the application of PAPAYA in its context would solve the current limitations coming from implementations that do not consider privacy-enhancing technologies.

This document, when complemented with the outcome of tasks T2.2 (“*User requirements (privacy and usability)*”) and T2.3 (“*Platform requirements (functional and utility)*”), characterizes the ground for understanding the expected key features of the PAPAYA platform, and thus serves as a basis for the design and development of the privacy-enhancing technologies (as in WP3), the design and development of the PAPAYA platform (as in WP4) and the validation of the PAPAYA project (as in WP5).



Project No. 786767

Glossary of Terms

| | |
|--------|---|
| BMI | Body Mass Index |
| ECG | Electrocardiogram |
| GDPR | General Data Protection Regulation |
| ID | Identifier |
| IoT | Internet of Things |
| MCI | MediaClinics Italia |
| MSISDN | Mobile Station Integrated Services Digital Network Number |
| NN | Neural Network |
| OMO | Orange Mobile Operator |
| ONU | Orange Network User |
| ORA | Orange |
| ToC | Table of Contents |
| TPC | Third-Party Customer |
| T2.1 | Task 2.1 Use case definition |
| WP | Work Package |



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

1 Introduction

1.1 Purpose of the document

This document reports the results of the work carried out for the execution of task T2.1 (“*Use case definition*”). The aim of this work is to identify the use cases (both in the healthcare and mobile and phone usage umbrellas) where PAPAYA will play a role. This document will thus serve as a ground for the definition of specific functionalities to be integrated in PAPAYA in the following Work Packages and Tasks.

The definition of use cases explores two different architectures where PAPAYA would enable the outsource of analytics to an external processor. In this project, we will address the following use cases:

- **Arrhythmia detection use case (UC1, in the healthcare umbrella).** In this use case (presented in Section 2.1), sensitive health data in the form of ECG are collected by a single source and an external entity that runs the PAPAYA platform detects arrhythmias from the ECG signal using neural networks. This use case maps with the first PAPAYA usage scenario (“Single data owner” - the data owner applies data analytics primitives to his sensitive data) defined in the description of work of the project.
- **Stress detection use case (UC2, in the healthcare umbrella).** In this use case (presented in Section 2.2), sensitive health data from IoT sensors are collected by multiple sources and used to train a collaborative model (using neural networks) via the PAPAYA platform, with the ultimate goal of automatically detecting stress conditions in workers. This use case maps with the second PAPAYA usage scenario (“Multiple data owners, collaborative work” - multiple data owners process a large dataset containing data from all the different data owners and derive relevant information such as global machine learning model) defined in the description of work of the project.
- **Mobility analytics use case (UC3, in the mobile and phone usage umbrella).** In this use case (presented in Section 3.1), sensitive data are collected from a single source and an external entity that runs the PAPAYA platform measures the audience in one or several areas (using counting techniques based on Bloom filters) or extracts mobility patterns (using some clustering algorithms). This use case maps with the third PAPAYA usage scenario (“Single-source data owner, third-party querier” - the data comes from a single source that protects it from being read by a third party; however, the data owner allows a third party to perform analytics tasks over its encrypted data, provided that the third party will only learn the analytics result) defined in the description of work of the project.
- **Mobile usage analytics use case (UC4, in the mobile and phone usage umbrella).** In this use case (presented in Section 3.2), sensitive data on mobile phone application usage are collected by multiple sources and used to extract analytics (such as counting or



D2.1 – Use Case Specification Dissemination Level – PU

Project No. 786767

aggregation) for statistical purposes. This use case maps with the fourth PAPAYA usage scenario (“Multiple source user, third party querier” - The data to be analyzed comes from different sources and is queried by a third party; neither the server nor the querier sees the collected data in clear, but only in an encrypted version, so as to achieve end-to-end privacy) defined in the description of work of the project.

- **Threat detection use case (UC5, in the mobile and phone usage umbrella).** In this use case (presented in Chapter 4), business-sensitive data from several sources (rather than personal data) are processed, to detect threats in systems or networks (using anomaly detection algorithm). This use case was added as an extra usage.

Notice that each use case is tagged with an acronym in the form UCx. These acronyms will be used throughout the upcoming deliverables to refer to the aforementioned use cases as well.

As each of the proposed use case deals with different analytics (e.g., neural networks, clustering, counting, etc.), they all raise different privacy requirements, hence enlarging the number of techniques to be considered when extracting analytics in a privacy-preserving way.

The definition of the identified use cases builds on discussions and reviews with end users and stakeholders involved in the healthcare and mobile and phone usage umbrellas, and describe how the system is supposed to behave in the covered situations. The provided description underlines how the work done on each of the identified use cases solves problems related to the as-is instantiations of the services. Indeed, most of the use cases start from an actual implementation of a service that is actually lacking some privacy-preserving solution to outsource data storage or processing to untrusted environments, and builds up from there, to show how the application of PAPAYA and its privacy-preserving technologies would enable a more secure and scalable processing of data.

Operatively, for each use case, this document presents the use case definition (which includes the description of the involved players, the involved data, the processing description and the underlying protocol(s)) and the related privacy requirements.

1.2 Outline

The document is divided into two main parts: Chapter 2, where we review the two use cases (UC1 and UC2) in the healthcare umbrella, and Chapter 3, where we review UC3 and UC4 in the mobile and phone usage umbrella. Chapter 4 overviews UC5 where business-sensitive data instead of personal data is processed.



Project No. 786767

2 Healthcare use cases

In this chapter we report the use case definition for the healthcare umbrella, namely, the *Privacy-preserving arrhythmia detection* use case (UC1) and the *Privacy-preserving stress management* use case (UC2). The two use cases differ from the type of architecture they are putting in place, and the type of underlying analytics: in the first use case, ECG data coming from a single source (e.g., a hospital) are sent to untrusted premises so as to detect the arrhythmias; in the second use case, stress-related data coming from several sources (e.g., several companies) are used to collaboratively train a neural network model.

2.1 Single source architecture: Privacy-preserving arrhythmia detection

In this section we present the description of the first use case in the healthcare umbrella, i.e., the *privacy-preserving arrhythmia detection* use case (UC1).

2.1.1 Introduction and current solution limitations

The *privacy-preserving arrhythmia detection* use case (UC1) targets patients who need to perform cardiac parameters analyses for some clinical reason, e.g., in the case where the patient suffers from a chronic condition, with the goal of verifying the presence/absence of arrhythmias.

The analysis of patients' cardiac parameters is provided by MCI, and it is offered to data subjects (patients) as a commercial service via pharmacies. With this service, each patient is given a wearable device that collects his/her ECG data during a fixed amount of time (e.g., 24 hours), with the goal of getting them analysed and reported by a healthcare professional. Once the patient returns the device to the pharmacy, ECG data are uploaded to the "MCI internal cloud", so that professionals can access them and report findings about his/her cardiac health.

The current deployment of such a service comes with some limitations. First, the burden of analyzing long streams of ECG data for a large number of patients may be difficult to be handled on premises, where the potential lack of computational resources would limit the performance of the analysis. To overcome this issue and increase the available computational and storage resources, one could think of adding a cloud environment to the deployment setting. In this view, the deployment of the service would be done in a mixed trusted-untrusted environment, where:

- an **on-premises site** (trusted), possibly limited in resources, retains the data as collected by the patient and to be served to the cardiologist;
- a **cloud environment** (untrusted) automatically extracts and analyzes relevant ECG sections to be reviewed by the cardiologist (avoiding the cardiologist to analyze the whole stream of data). This may be currently done via the usage of some untrusted available



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

service, such as CardioCalm¹, which is neither in control of MCI nor in control of the patient.

Nevertheless, the more we move away from trusted premises (as in the case in which we move the computation to the cloud side), the more the data privacy is endangered. This is a critical issue, as the most recent regulations on data protection (such as the GDPR) impose strict analysis constraints for the so-called “special categories of data” (as of Article 9), which include health-related data. In conclusion, we identify two limitations: i) on-premises-only solutions suffer from the lack of computational resources, for the sake of preserving patients’ data protection; ii) mixed deployments that use untrusted cloud environments may lack the appropriate security measures to preserve patients’ privacy.

2.1.2 Use case definition

In this section, we detail the definition of the use case. To tackle the aforementioned limitations, we will hypothesize to have ECG record analysis outsourced to PAPAYA, so that highly demanding computations are tackled in a large cloud environment while still ensuring patients’ privacy.

2.1.2.1 Story

Giacomo Rossi is a patient with a chronic disease that may unfortunately result in cardiac problems over time. To keep his health status in check, his medical doctor prescribed him a periodic cardiac check-up.

Giacomo discovers from a friend that pharmacies are offering patients the access to the CardioPharma service, which may allow them to perform a cardiac health assessment in an easy and convenient way. As this is exactly fitting with his doctor’s prescription, Giacomo goes to the nearest pharmacy to ask for a CardioPharma access.

Paola Bianchi, the pharmacist, explains Giacomo how CardioPharma works. Giacomo would be provided with a wearable device, called MC CardioMonitor, which would collect his ECG data over a fixed period of 24 hours. After one day, he would be asked to return the CardioMonitor to the pharmacy, so that the acquired ECG data would be uploaded to the MCI platform, pre-processed to find significant patterns, analyzed by a cardiologist, and reported in a document that would be forwarded to him. As the service is simple to use and without harm for him, Giacomo decides to apply for CardioPharma and signs the consent form.

Paola fetches her tablet and opens the application. As this is the first time Giacomo is using the CardioPharma service, Paola registers Giacomo’s profile by collecting basic information (among which, e.g., first and last name, gender, date of birth). A basic anamnesis is collected too: it states which are the drug therapies Giacomo is subjected to, registers some anthropometric measures (i.e., height and weight), and indicates the presence/absence of a pacemaker. All these data are

¹ <http://www.cardiocalm.com/site/>



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

indicated as *personal data* by the GDPR and include *special categories of data* (i.e., the health-related ones). Now Paola fetches a CardioMonitor wearable device, pairs it with Giacomo's profile, makes Giacomo wear it and starts the data acquisition. Giacomo is free to go home and perform his daily routine, knowing that he will have to come back to the pharmacy after 24 hours.

A day has passed and Giacomo comes back to Paola's pharmacy to return the CardioMonitor device. Paola collects the device, downloads the ECG data and pushes them on MCI cloud.

At the Sacro Cuore healthcare facility, Carlo Neri (a cardiologist) gets a notification stating that a new reporting activity is required for a chronic patient. The incoming data are Giacomo's ones, but yet Carlo will not be able to see Giacomo's identifiers (e.g., name and surname): the only data he will access will be the anamnesis and the ECG data. Carlo accesses the CardioPharma-Doctor web application, accepts the incoming job and opens the related data. Carlo is happy to see that the CardioPharma platform has already conducted a series of preprocessing activities on the data (via the PAPAYA platform), so that he can access the whole stream of ECG data if he needs to, but still some significant areas containing possible proofs for health problems are highlighted. Carlo starts analyzing the ECG data, and when he is done, he reports his findings in a document. The report, signed by Carlo, contains data that revolve around Giacomo's health profile: his BMI, the presence/absence of arrhythmias, his pulse rate and a brief professional comment.

Paola receives back Carlo's report, and forwards it to Giacomo.

2.1.2.2 Players

In the following, we list the stakeholders playing a role in this use case.

Table 1UC1 Players

| | |
|---------------------|---|
| Patient | The patient (i.e., the 'data subject' in GDPR perspective) uses the CardioPharma service to get his cardiac health in check. Its GDPR role is data subject. |
| Pharmacist | This subject collects patients' data via the CardioMonitor wearable device, registers patients' profiles, receives the cardiologists' reports and forwards them to the patient. His objective is to provide new services to increase his own business, while bringing health services to a new level of capillarity. Its GDPR role is data processor. |
| Cardiologist | This subject visualizes pseudonymized patients' data (anamnesis and ECG data) and writes a cardiology report. The automatic analyses performed on the external cloud environment are crucial to keep the required analysis time as short and effective as possible. The cardiologist trusts the quality of the data collection and the quality of the analysis. |



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

| | |
|--------------------------------|---|
| | <p>Access to data is in pseudonymized form so as to minimize the risk of data disclosure. The doctor indeed belongs to the trusted area of the system, but has no need to know patients' names nor their identifiers to provide the ECG report.</p> <p>His GDPR role is data processor.</p> |
| MCI | <p>MCI retains the patients' data in his trusted cloud environment and is responsible for serving as a contact point for all the other players.</p> <p>Its GDPR role is data controller.</p> |
| External cloud provider | <p>The external cloud environment is used to perform preprocessing activities on ECG data. It collects raw ECG data and produces analyses. The PAPAYA platform could play a big role here</p> <p>Its GDPR role is data processor.</p> |

2.1.2.3 Involved data

In this section, we list the personal data whose processing is involved in the current use case. In this section, we indicate as **identifier** all the information that, alone, can be used to identify a person.

- **Patient's biographical data:** first and last name (identifier); healthcare ID (identifier); birth date; gender; ethnicity; contacts (may be identifier).
- **Anamnesis:** weight; height; drug therapies; usage of a pacemaker.
- **ECG:** ECG data.
- **ECG report:** BMI; arrhythmias; pulse rate; doctor's evaluation; doctor's signature (optional, identifier).

2.1.2.4 Processing and analytics

Artificial Neural Networks (NN) can support pharmacists and doctors to analyse patients' data and quickly diagnose issues or identify trends that would predispose them to a particular disease. Heart arrhythmia is considered as a serious disease. Recent research results show that early classification of arrhythmia types can help prevent fatal death. Heart arrhythmia can be detected by using Electro-Cardiograms (ECG) that record the electrical activities of the heart of a patient. Running a neural network over such data can help check for possible Atrial Fibrillation (AF). These neural network parameters (i.e., its model) will be stored at an untrusted server and should be executed over protected ECG data.

For the demonstration of this use case, the neural network will be trained locally on publicly available data, and the trained model will be used at runtime to classify ECGs.



D2.1 – Use Case Specification Dissemination Level – PU

Project No. 786767

Processing would be subjected to what GDPR dictates. Specifically, implementation of appropriate technical and organisational measures to secure personal data shall be put in place when data are transferred between environments, being them either trusted or untrusted (i.e., *encryption upon transferring*), or stored in untrusted environments (i.e., *encryption at rest*). Also, processing of patient’s personal data (being them personal data or of a special category of data, e.g., health-related) is possible only when consent for manipulation is given by the patient.

2.1.2.5 Protocol

In this section, we present the structure of the protocol within UC1, including preconditions, postconditions and flow.

Table 2 Description of UC1

| | |
|-------------------------|---|
| ID and name | UC-ECG-1 Detect arrhythmias in a patient’s ECG |
| Primary actor | Patient |
| Secondary actors | Pharmacist, MCI, external cloud provider, cardiologist |
| Description | A patient needs to check his ECG (for some clinical reason, e.g., for a chronic condition). He uses the CardioPharma service so as to collect his ECG data for some time, extract data about his arrhythmias and get a report from a cardiologist |
| Preconditions | <p>PRE-1 Pharmacy is registered to the CardioPharma service</p> <p>PRE-2 A consent to treat patient’s data is signed by the patient</p> <p>PRE-3 Neural network is trained to recognize arrhythmias in ECG signals</p> <p>PRE-4 Neural network model is uploaded on the external cloud (i.e., the PAPAYA platform)</p> |
| Postconditions | POST-1 A report for the patient, redacted by the cardiologist, is available for its download |
| Normal flow | <p>1.0 Generation of patient’s arrhythmias report</p> <ol style="list-style-type: none"> 1. Pharmacist enters the patient’s identifier (e.g., healthcare ID) 2. MCI app retrieves patient’s profile and displays it 3. Pharmacist enters the CardioMonitor identifier to pair it with patient 4. Pharmacist sets up the CardioMonitor and gives it to patient 5. The CardioMonitor starts the ECG data recording, which will last for a fixed amount of time (e.g., 24 hours) |



Project No. 786767

D2.1 – Use Case Specification Dissemination Level – PU

| | |
|--------------------------------|---|
| | <ol style="list-style-type: none"> 6. When the monitoring period is over, pharmacist retrieves the CardioMonitor from the patient and downloads the ECG data 7. MCI app uploads the ECG data on the MCI cloud 8. MCI cloud encrypts the ECG data 9. MCI cloud forwards encrypted ECG data to the external cloud 10. The PAPAYA platform on the external cloud analyses the ECG data so as to detect arrhythmias in an oblivious manner 11. The external cloud returns the encrypted results about patient's arrhythmia classification to the MCI cloud 12. MCI cloud decrypts results 13. MCI cloud forwards the (cleartext but pseudonymized) data about the patient's profile, the patient's full ECG and the patient's classified heart beats to cardiologist 14. Cardiologist produces a report for the patient 15. MCI system forwards the produced report to the pharmacy 16. Pharmacist notifies the patient about the availability of the report |
| <p>Alternative flow</p> | <p>1.1 Generation of patient's arrhythmias report when no arrhythmias are found</p> <ol style="list-style-type: none"> 1. Pharmacist enters the patient's identifier (e.g., healthcare ID) 2. MCI app retrieves patient's profile and displays it 3. Pharmacist enters the CardioMonitor identifier to pair it with patient 4. Pharmacist sets up the CardioMonitor and gives it to patient 5. The CardioMonitor starts the ECG data recording, which will last for a fixed amount of time (e.g., 24 hours) 6. When the monitoring period is over, pharmacist retrieves the CardioMonitor from the patient and downloads the ECG data 7. System uploads the ECG data on the MCI cloud 8. MCI cloud protects the ECG data 9. MCI cloud forwards protected ECG data to the external cloud 10. The PAPAYA platform on the external cloud analyses the ECG data so as to detect arrhythmias |



Project No. 786767

D2.1 – Use Case Specification

Dissemination Level – PU

| | |
|---------------------------|--|
| | <ol style="list-style-type: none"> 11. External cloud returns the protected results about patient's arrhythmia classification to the MCI cloud, with a notification of absence of arrhythmias 12. MCI cloud forwards the (cleartext but pseudonymized) data about patient's profile and full patient's ECG to cardiologist, notifying him about the absence of detected arrhythmias 13. Cardiologist produces a report for the patient 14. MCI system forwards the produced report to the pharmacy 15. Pharmacist notifies the patient about the availability of the report |
| <p>Exceptions</p> | <p>1.0.E1 Patient's profile cannot be retrieved</p> <ol style="list-style-type: none"> 1. MCI app displays message: no data for selected patient 2. MCI app terminates the use case <p>1.0.E2 ECG signal is noisy and cannot be processed</p> <ol style="list-style-type: none"> 1. Pharmacist enters the patient's identifier (e.g., healthcare ID) 2. MCI app retrieves patient's profile and displays it 3. Pharmacist enters the CardioMonitor identifier to pair it with patient 4. Pharmacist sets up the CardioMonitor and gives it to patient 5. The CardioMonitor starts the ECG data recording, which will last for a fixed amount of time (e.g., 24 hours) 6. When the monitoring period is over, pharmacist retrieves the CardioMonitor from the patient and downloads the ECG data 7. MCI app uploads the ECG data on the MCI cloud 8. MCI cloud encrypts the ECG data 9. MCI cloud forwards encrypted ECG data to the external cloud 10. The PAPAYA platform on the external cloud analyses the ECG data so as to detect arrhythmias 11. The PAPAYA platform recognizes that ECG is noisy and thus notifies the pharmacist to repeat the ECG acquisition from patient |
| <p>Assumptions</p> | <p>Sensitive data is safely moved from the trusted to the untrusted party (and vice versa).</p> |



D2.1 – Use Case Specification Dissemination Level – PU

Project No. 786767

| | |
|--|---|
| | Processing in the untrusted environment is protected thanks to the use of privacy-enhancing technologies. |
|--|---|

2.1.3 Privacy requirements

MCI is considered as the Data Owner and controller. Pharmacist and Cardiologist are trusted parties that have the patient’s consent on processing his/her data. Raw data is transformed in protected data structures; once protected, they can be provided to third-party entities. The untrusted area is the external cloud environment (i.e., the PAPAYA platform in this project) where data are handled in a protected way. In here, data are processed so as to classify the heartbeats and hence detect arrhythmia.

In the following, we present the requirements related to privacy.

Table 3 Privacy Requirements for UC1

| Type | Description |
|------------|---|
| Data | Patients’ data shall be pseudonymized when sent to the cardiologist |
| Data | Patients’ data shall be protected via PETs before sending them to the external cloud for classification |
| Data | Re-identification of patients on data outsourced to the PAPAYA platform shall not be possible |
| Functional | Consent shall be handled by the system, as a lawful basis for processing |
| Functional | Performing any other analytics for other purposes rather than the one specified in the consent (i.e., analysis of ECG data to detect arrhythmias) shall be infeasible |
| Functional | Processing shall be denied when a valid consent from the patient is not provided |

2.2 Multiple source architecture: Privacy-preserving stress management

In this section we present the description of the second use case in the Healthcare umbrella, i.e., the *privacy-preserving stress management* use case (UC2).

2.2.1 Introduction and current solution limitations

The *privacy-preserving stress management* use case (UC2) targets workers that, for various reasons (related either to work or to daily life) suffer from stress, which affects their life negatively.

As it is difficult for people who suffer from stress to identify raising stress and anxiety levels at their onset, they often find themselves in the unfortunate situation of being left with a severe stress



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

situation and little or no help from the surrounding environment to solve this. It would be thus very helpful to have an automatic solution that would help anxious and stressed people to recognize symptoms at their onset, and prevent their raise by proposing mitigation actions to keep stress levels in check. Nevertheless, to the best of our knowledge, there currently exist few (if not none) solutions for the automatic management of stress (i.e. automatic detection of stress symptoms and correct elaboration of the data acquired).

It would be thus very useful to apply machine learning approaches to this situation: where a machine can identify a stress situation, there is no need from the external environment in identifying problematic symptoms.

MCI is currently testing a sensorized T-shirt that is able to collect some health-related parameters. By training a neural network on such collected parameters, an automatic classification would be able to notify users about raising stress levels, and propose on-the-fly (short-term) countermeasures. Nevertheless:

- collecting data from a single person would require a really long time to reach a significant dataset size, and would end up in building limited models that are not really representative of what a “stress situation” is;
- collecting data from a large set of people would require a really large amount of computational and storage resources, which could be provided by means of a cloud environment.

As discussed in the previous section, the more we move away from trusted premises (as in the case in which we move the computation to the cloud side), the more the data privacy is endangered. This use case also involves health-related parameters, which are classified as “special categories of data” (as of Article 9), and thus need special protection measures to preserve workers’ privacy.

2.2.2 Use case definition

In this section, we detail the definition of the use case. To tackle the aforementioned limitations, we will hypothesize to have the neural network training outsourced to PAPAYA, so that highly-demanding computations are tackled in a large cloud environment (the PAPAYA platform) while still ensuring workers’ privacy.

2.2.2.1 Story

Dataset collection

ITSoft is a large IT company that employs different programmers and software engineers in Northern Italy. ITSoft has two headquarters: one in Trento and one in Milan.

Matteo and Alex are two employees in ITSoft-Trento, while Paolo and Luca are two employees in ITSoft-Milan. All of them suffer from high stress levels, due to the hard period at work and the



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

constantly pressing deadlines, which force them to stay at work until late and consume extra energy in their daily routine.

ITSoft director hears about a study proposed by MCI that would help to solve workers' stress issues. During such a trial, MCI would collect workers' health-related data (under consent) and train some machine learning algorithm that would automatically identify stress situations. As the outcome of this trial would be valuable for ITSoft workers, the director decides to participate.

A workplace aggregator for data collection (called MC HealthCorner) is thus installed on both premises (both in Trento and in Milan), and a bunch of sensorized T-shirts is sent to ITSoft headquarters to be distributed among workers that decide to participate.

Matteo, Alex, Paolo and Luca all decide to participate to the trial, because they think it would be of help for them and for all their workers who, as them, feel stressed while working. They sign the consent, fill in all the privacy preferences (e.g., by specifying the hours in which they want to be monitored), wear the sensorized T-shirt they are given with, and go to work. Every time they recognize that the stress level is raising, they use the MCI app to tag the current moment as stressful and go back to work. Their tagging actions actually result in:

- storing both the tag and the current health-parameters as read from the T-shirt;
- sending the so-created tuple to the MC HealthCorner, which serves as workplace aggregator for the collection of workplace-related datasets.

Model training

When the datasets are of sufficient size (e.g., are considered sufficient by an expert in the ML field), each MC HealthCorner trains its own neural network on the collected data, and sends part of the neural network to the PAPAYA platform. The PAPAYA platform builds a collective model, which (due to the larger amount of data available) performs with higher accuracy with respect to the local ones produced by the MC HealthCorner. This collaboratively computed model, once trained, is sent back to the MC HealthCorners, so as to be used by workers.

Classification of stress situations

Now that each MC HealthCorner in ITSoft has its own copy of the collectively trained model, each worker in ITSoft can start using it to automatically identify stress situation.

Lorenzo, who works at ITSoft-Milan, wears one of MCI sensorized T-shirts, downloads the MCI app and goes to work. Data from the T-shirt are sent to the phone and bounced back to the MC HealthCorner for stress/non-stress classification, in a real-time fashion. Every time the MC HealthCorner classifies the current parameters as stress-related, Lorenzo's phone displays a notification reporting the stress alert. In that case, Lorenzo understands that it is time to have a break at work, and follows the exercises and instructions shown by the MCI app to lower his stress.



D2.1 – Use Case Specification Dissemination Level – PU

Project No. 786767

2.2.2.2 Players

In the following, we list the stakeholders playing a role in this use case.

Table 4 UC2 Players

| | |
|--------------------------------|--|
| Worker | <p>The worker (i.e., the ‘data subject’ in GDPR perspective) uses the stress trainer and classification tool to:</p> <ul style="list-style-type: none"> • either enrich the dataset of collected health-related data • or automatically detect stress situations while at work <p>His GDPR role is data subject.</p> |
| Employer | <p>This subject provides the stress management service to his employees, by installing a MC HealthCorner on premises and distributing sensorized T-shirts to workers.</p> |
| MC HealthCorner | <p>The MC HealthCorner functions as workplace aggregator, to:</p> <ul style="list-style-type: none"> • collect workers’ data from each premises <p>classify on-the-fly workers’ health-related parameters (as coming from the T-shirts) so as to identify stress situations</p> |
| MCI | <p>MCI is responsible for serving as a contact point for all players, especially for what concerns data subjects’ rights in GDPR term.</p> <p>Its GDPR role is data controller.</p> |
| External cloud provider | <p>The external cloud environment is used to perform NN models aggregation as coming from several workplaces. PAPAYA could play a big role here.</p> <p>Its GDPR role is data processor.</p> |

2.2.2.3 Involved data

In this section, we list the personal data whose processing is involved in the current use case. In this section, we indicate as identifier all the information that, alone, can be used to identify a person.

- **Biographical data:** first and last name (identifiers); healthcare ID (identifier); date of birth; gender; race; contacts (identifier)
- **Periodic data:** height; weight; blood pressure; glycemia; questionnaires (identifier)
- **Wearable data:** steps; breath rate; heart rate; oximetry (not yet implemented); skin bioimpedance (not yet implemented); body temperature (not yet implemented)
- **Worker’s annotation** (textual)



D2.1 – Use Case Specification Dissemination Level – PU

Project No. 786767

2.2.2.4 Processing and analytics

According to Galatzer-Levy et al.², various environmental and biological systems interact to influence individual differences in response to stress. Considering that the integration of so many heterogeneous and diffuse information to study and assess stress levels represents a significant and groundbreaking challenge, different approaches to stress management should be taken into account. With this in mind, a wide range - at least at an earlier stage - of supervised machine learning techniques must be considered to achieve the goals that we have set. One possible approach is to develop a classification model based on simple NN (i.e. Multi-Layer Perceptron (MLP)). The model will be trained on labelled (stress/no stress) data and will be able to classify previously unseen samples as stress/no stress. Processing would be subjected to what GDPR dictates. Specifically, implementation of appropriate technical and organisational measures to secure personal data shall be put in place when data are transferred between environments, being them either trusted or untrusted (i.e., *encryption upon transferring*), or stored in untrusted environments (i.e., *encryption at rest*).

2.2.2.5 Protocol

In this section, we present the structure of the use case, including preconditions, postconditions and the protocol flows.

In the following, we show the dataset collection protocol:

Table 5 Description of UC2 - Dataset collection

| | |
|-------------------------|---|
| ID and name | UC-STR-1 Collect stress-related dataset |
| Primary actor | Worker |
| Secondary actors | MC HealthCorner, MCI |
| Description | Workers in a workplace contribute to the creation of a stress-related dataset, by collecting their health-related parameters via a sensorized T-shirt and tagging them with a stress/non-stress label |
| Preconditions | <p>PRE-1 Worker's company is registered to the service</p> <p>PRE-2 Worker is registered to the service</p> <p>PRE-3 A consent to treat worker's data is signed by the worker</p> <p>PRE-4 Worker wears the MCI sensorized T-shirt</p> <p>PRE-5 Worker has filled in the privacy preferences</p> |

² Isaac R. Galatzer-Levy, Kelly V. Ruggles, and Zhe Chen, *Data Science in the Research Domain Criteria Era: Relevance of Machine Learning to the Study of Stress Pathology, Recovery, and Resilience*, Chronic Stress, Volume 2: 1–14, 2018



D2.1 – Use Case Specification Dissemination Level – PU

Project No. 786767

| | |
|-------------------------|--|
| Postconditions | POST-1 A dataset is collected on the MC HealthCorner |
| Normal flow | <p>1.0 Dataset collection when in stress situation</p> <ol style="list-style-type: none"> 1. Worker opens the MCI mobile app 2. Worker enters the tag “stress” for the current situation 3. MCI mobile app verifies that the data gathering step complies with the data subject’s privacy preferences 4. MCI mobile app reads health-related data from the T-shirt 5. MCI mobile app sends the tag and the read health-related data to the MC HealthCorner, as a new entry for the dataset 6. MC HealthCorner enriches the dataset with the new entry |
| Alternative flow | <p>1.1 Dataset collection with in non-stress situation</p> <ol style="list-style-type: none"> 1. Worker opens the MCI mobile app 2. Worker enters the tag “non-stress” for the current situation 3. MCI mobile app verifies that the data gathering step complies with the data subject’s privacy preferences 4. MCI mobile app reads health-related data from the T-shirt 5. MCI mobile app sends the tag and the read health-related data to the MC HealthCorner, as a new entry for the dataset 6. MC HealthCorner enriches the dataset with the new entry |
| Exceptions | <p>1.0.E1 Dataset collection when privacy preferences are non-fitting</p> <ol style="list-style-type: none"> 1. Worker opens the MCI mobile app 2. Worker enters the tag “non-stress” for the current situation 3. MCI mobile app verifies that the data gathering step is NOT compliant with the data subject’s privacy preferences 4. MCI mobile app end the use case |
| Assumptions | - |

In the following, we describe the model training protocol:

Table 6 Description of UC2 - Model training



D2.1 – Use Case Specification Dissemination Level – PU

Project No. 786767

| | |
|-------------------------|--|
| ID and name | UC-STR-2 Train collective model |
| Primary actor | External cloud provider |
| Secondary actors | MC HealthCorner |
| Description | A collective model is trained on cloud premises and redistributed to the local nodes, so that it can be used to identify on-the-fly stress situations |
| Preconditions | <p>PRE-1 Company is registered to the service</p> <p>PRE-2 MC HealthCorner is installed on the company premises</p> <p>PRE-3 MC HealthCorner contains workers' data</p> |
| Postconditions | POST-1 The collective model is available on the MC HealthCorner |
| Normal flow | <p>1.0 Model training</p> <ol style="list-style-type: none"> 1. Do iteratively (until accuracy doesn't improve or predefined number of iterations): 2. MC HealthCorner trains its own neural network 3. MC HealthCorner obfuscates and sends part of the neural network to the external cloud (i.e., PAPAYA platform) 4. PAPAYA platform integrates models from each healthcorner into collective model 5. MC HealthCorner downloads the collective model from the PAPAYA platform |
| Alternative flow | - |
| Exceptions | - |
| Assumptions | <p>Sensitive data is safely moved from the trusted to the PAPAYA platform (and vice versa).</p> <p>Processing in the untrusted environment is protected with privacy-enhancing technologies.</p> |

In the following, we describe the classification protocol:

Table 7 Description of UC2 - Stress detection and classification

| | |
|----------------------|---|
| ID and name | UC-STR-3 Classify worker's data to identify a stress condition |
| Primary actor | MC HealthCorner |



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

| | |
|-------------------------|--|
| Secondary actors | Worker |
| Description | The trained collective model is used to classify on-the-fly worker's data, to understand if he is in a stress condition and if he needs suggestions to lower the stress level |
| Preconditions | <p>PRE-1 Worker's company is registered to the service</p> <p>PRE-2 Worker is registered to the service</p> <p>PRE-3 A consent to treat worker's data is signed by the worker</p> <p>PRE-4 Worker wears the MCI sensorized T-shirt</p> <p>PRE-5 MC HealthCorner contains the trained collective model</p> <p>PRE-6 Worker has filled in the privacy preferences</p> |
| Postconditions | POST-1 A stress label is sent to worker, with the mitigation actions to be applied to lower the stress |
| Normal flow | <p>1.0 Stress classification in stress condition</p> <ol style="list-style-type: none"> 1. Worker's mobile phone collects current health-related parameters from worker's T-shirt 2. System verifies that the data gathering step complies with the data subject's privacy preferences 3. Worker's health-related parameters are sent to MC HealthCorner for classification 4. MC HealthCorner use collectively trained model to classify worker's health-related parameter as stress-related 5. MC HealthCorner sends "stress" notification to worker's phone 6. Mobile phone shows stress alert, plus the set of exercises to be followed to lower the stress |
| Alternative flow | <p>1.1 Stress classification in non-stress condition</p> <ol style="list-style-type: none"> 1. Worker's mobile phone collects current health-related parameters from worker's T-shirt 2. System verifies that the data gathering step complies with the data subject's privacy preferences 3. Worker's health-related parameters are sent to MC HealthCorner for classification |



D2.1 – Use Case Specification Dissemination Level – PU

Project No. 786767

| | |
|--------------------|---|
| | 4. MC HealthCorner use collectively trained model to classify worker's health-related parameter as non-stress-related |
| Exceptions | <p>1.0.E1 Stress classification when privacy preferences are non-fitting</p> <ol style="list-style-type: none"> 1. Worker's mobile phone collects current health-related parameters from worker's T-shirt 2. System verifies that the data gathering step is NOT compliant with the data subject's privacy preferences 3. System end the use case |
| Assumptions | <p>Sensitive data is safely moved from the trusted to the PAPAYA platform (and vice versa).</p> <p>Processing in the untrusted environment is protected with privacy-enhancing technologies.</p> |

2.2.3 Privacy requirements

MCI is considered as the Data controller. The data subject (i.e., the worker) is requested to give his/her consent, furthermore the data subject will configure his/her privacy preferences and only the data which complies with them could be processed. Raw data is transformed in protected data structures; once protected, they can be provided to third-party entities. The untrusted area is the external cloud environment where data are handled in an encrypted format. In here, data are processed so as to train the joint model.

In the following, we present the requirements related to privacy.

Table 8 Privacy requirements for UC2

| Type | Description |
|-------------------|--|
| Data | Neural network model input shall be protected via PETs before outsourcing them to the external cloud |
| Data | Re-identification of workers from models outsourced to the PAPAYA platform shall not be possible |
| Functional | Consent shall be handled by the system, as a lawful basis for processing |
| Functional | Performing any other analytics for other purposes rather than the one specified in the consent shall be infeasible |
| Functional | Processing shall be denied when a valid consent from the worker is not provided |



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

| | |
|-------------------|---|
| Functional | Privacy preferences (e.g., decision on when to monitor workers and what data to monitor) shall be handled by the system |
| Functional | Data collection shall be performed only when compliant to the privacy preferences specified by the workers |
| Functional | Retrieving data from single workers shall not be feasible |
| Functional | Computation of models on data coming from a single worker shall not be feasible (to avoid re-identification) |
| Functional | Analytics shall not be performed before multisource data are aggregated |
| Functional | Right to erasure shall be supported by the system |



Project No. 786767

3 Mobility and phone usage use cases

In this chapter, we report the use case description for the mobile and phone usage umbrella, namely, the Privacy-preserving mobility analytics use case (UC3) and the Privacy-preserving mobile phone usage analytics (US4) use case.

3.1 Single source architecture: Privacy-preserving mobility analytics

In this section, we present the description of the first use case in the mobile and phone usage umbrella, i.e., the *privacy-preserving mobility analytics* use case (UC3).

3.1.1 Introduction and current solution limitations

The *privacy-preserving mobility analytics* use case targets Orange's Third-Party Customers (TPC) interested in some analytics based on individuals' mobility for their own business. For example, mobility analytics can be profitable to several kinds of TPCs such as tourism development agencies, tourist offices, amusement parks, hotels, exhibition centers, stadiums capable of hosting all types of events (e.g.: festivals); etc. Indeed, the provision of insights on the visitors/tourists and their mobility patterns have strong implications for such organizations. Understanding these patterns could help those TPCs managing infrastructure planning, enhance visitors' experience or tailor tourism offerings, in order to increase their revenues and better satisfy visitors' needs.

On its side, Orange Mobile Operator (OMO), as a telecom operator, has access to a great amount of data, continuously generated by Orange Network Users (ONU) (whether they are Orange subscribers or not, e.g., in case of roaming), through their use of mobile phones and Orange's network infrastructure.

In fact, each time an individual uses his smartphone for a call, an SMS or an Internet access, the communication necessarily goes through mobile network antennas, so as to process it to the right place (the contacted person or server). From the antenna, the communication is sent to the mobile network operator for communications' conveyance, billing and legal obligations. From the communication, the mobile network operator can extract (i) the identity of the individual (through the mobile phone number), (ii) the timestamp of the communication and (iii) the individual's location since it knows the one of the antenna.

In other words, OMO is sitting on a goldmine of valuable information, which, if enriched with sociodemographic data, is of tremendous value for the TPCs.

Traditionally, these TPCs have gained visitors insights and mobility patterns through field surveys, questionnaires and interviews. Such tools, while being used for decades, are cumbersome to implement, time consuming and expensive, inaccurate and biased or representative of only a small sample of the population. By leveraging its operated data, OMO can offer a wide range of key insights about visitors and tourists to the TPCs, via dedicated analytics of the mobile network



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

data generated by the ONUs. Compared to conventional surveys, these insights include real-time, accurate and population-wide information. Besides, on Orange's side, no special and costly infrastructure are needed to be deployed for the analytics of the data, since the existing infrastructure owned by OMO serves the initial purpose of conveying communications and gathers the precious data.

Therefore, in this use case, both the TPC and OMO maintain a business relationship in which (1) the TPC expresses its needs for some visitor insights in the form of a data analytics request that is sent to OMO and (2) OMO sells to the TPC the analytics results without revealing the raw data that have been used to perform the analytics.

For our study, two kinds of analytics are considered.

- **Audience measurements:** it consists in counting the number of people in one or several areas of observation during a period of observation using Bloom filters [1]. This type of analytics also derives the origin of the visitors and the duration of their visits. It may also permit to know whether some individuals have been in two different areas during the period of observation.
- **Trajectories analysis:** it extracts information on mobility patterns, that is, information on how people travel from origin O to destination D and the amount of people flowing from O to D.

In both cases, our aim is to process probe data in a privacy-preserving way, i.e. the data processing will let the data resistant to inferences and re-identification attacks.

The limitations of the system are dictated by current and future data privacy regulations (namely, the GDPR³ and the ePrivacy⁴), that necessitate to define a legal basis for such a data process. After that, a risk analysis related to data protection, which will depend on both the legal basis and the technical solutions that will be used, may identify potential remaining risks. Mobility analytics may also depend on the French legislation, and in particular the notion of “on-the-fly anonymization” (“*anonymisation à bref délai*”) and the one of further processing for statistical purposes. The PAPAYA platform, that will be created and maintained by the project's partners, offers a unique opportunity for OMO to turn its data and insights into value, while complying with the abovementioned regulations and preserving users' privacy. By integrating this use case into the PAPAYA platform, OMO will benefit from easy-to-access and easy-to-use modules for privacy-preserving mobility analytics, which apply the appropriate protection mechanisms, designed with the expertise of the PAPAYA partners.

³ GDPR: <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=celex%3A32016R0679>

⁴ Proposal for the ePrivacy regulation for the protection of personal data in electronic communications <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX%3A52017PC0010>



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

3.1.2 Use case definition

In this section, we detail the use case description.

3.1.2.1 Story

The city of Deauville, France, organizes each year the Deauville American Film Festival, a film festival dedicated to American cinema.

For several years, the Deauville festival has faced, like all other festivals, the disengagement of the major film studios and a drop in its attendances. To try to stop this decline, the director decides to launch a study of attendance; he wants to better understand where visitors come from (either local people or tourists), how many unique visitors attend the festival and how long they stay there. The director thinks that such insights on festival’s visitors will help boost the vitality of the event.

The director of the Deauville festival has heard about OMO’s privacy-preserving Visitor Analytics offer. Privacy protection is a fundamental issue for the reputation of the festival. Therefore, he contacts his OMO correspondent to ask him for some analytics about the visitors. The request is formed by the director of the festival and specifies the geographical area (including points of origin and destination) and the time period of the observation.

The OMO then sets up the data collection and processing actions necessary for the production of mobility indicators, thanks to the PAPAYA platform. OMO has the role of Data Owner and Data Controller as stipulated in the GDPR. It is the guarantor that that must assess that at the end of the process, the data produced are perfectly resistant to inferences and re-identification attacks.

After the phase of observation and processing, OMO sends data indicator reports with a graphical presentation. These indicators concern statistical data on population flows. These data can also be enriched by socio-demographic criteria information such as the age group, the male / female distribution, and the socio-professional categories. Deauville festival’s director then exploits the data files according to the desired items, in a privacy-preserving way, and in total accordance with European and French laws on data protection.

3.1.2.2 Players

In the following we list the stakeholders playing a role in this use case.

Table 9 UC3 Players

| | |
|--|---|
| <p>Orange Network Users (ONU)</p> | <p>The ‘data subject’ in GDPR perspective. They connect to Orange network and use it. Either they have signed a service contract with the mobile operator, or they are foreign users in roaming that have no contract with the mobile operator, but can be contacted, e.g., by SMS.</p> |
|--|---|



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

| | |
|--|---|
| <p>Orange Mobile Operator (OMO)</p> | <p>It provides a network infrastructure and produces probe data of users' activities. OMO is the owner of the probe data.</p> <p>In this use case, OMO also processes probe data in a privacy-preserving way and produces mobility analytics to increase its own business. Resulting statistical indicators are in total accordance with European and French laws on data protection.</p> <p>To perform the mobility analytics, OMO runs an instance of the PAPAYA statistics service.</p> <p>Its GDPR role is Data Owner, Data Controller, and Data Processor.</p> |
| <p>Third Party Customer (TPC)</p> | <p>It is interested in mobility analytics results for its business.</p> <p>It expresses an analytics request to OMO which sells him the data analytics results without revealing the data that have been used to perform the analytics.</p> <p>It has no GDPR role.</p> |

3.1.2.3 Involved data

In this section, we list the personal data whose processing is involved in the current use case.

Probe data record signaling information transiting between ONUs and the Orange network (more precisely, network antennas that cover a specific area called a cell). This information is primarily collected for the purposes of conveying communication, monitoring the network infrastructure or billing. Raw data extracted from each probe data event, are:

- Orange Network User identifier (MSISDN);
- Timestamp (date and hour);
- Antenna number.

An on-the-fly pseudonymisation processing is done, so that pseudonymised data extracted from each probe event become:

- **Hash** of Orange Network User identifier (MSISDN);
- Timestamp interval;
- Antenna **localisation**.

3.1.2.4 Processing and analytics

As mentioned earlier, we consider two kinds of analytics:



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

1. **Audience measurements** (counting people using Bloom Filters).
2. **Trajectories analysis** (identifying and clustering trajectory patterns).

Audience measurements

Bloom filters (BF) [1] are simple space-efficient probabilistic data structures that are used to test whether an item x is a member of a given set S (of size n). They are typically defined by an array of m bits initially set to 0, and k hash functions h_1, \dots, h_k with range $\{1, \dots, m\}$. For an element $x \in S$, bits at positions $h_1(x), \dots, h_k(x)$ in the array are set to 1. To test whether an unknown item y is in S , we verify the bits at positions $h_1(y), \dots, h_k(y)$. If at least one of the bits is 0, then for sure $y \notin S$. However, if they are all set to 1, then there is a *certain probability* that $y \in S$. Indeed, whereas false negative never occur, false positives are possible since the bits set to 1 could have been computed by another item coincidentally. Nevertheless, the false positive rate can be kept at minimum by carefully choosing parameters m and k .

In our use case, a BF is attached to each antenna's (or group of antennas') location by OMO. Pseudonymised probe data are read in real time to achieve the counting. Unique people connected to an antenna (or group of antennas) are inserted in the BF for a given time interval (specified in the TPC's request).

After the counting period is over, BFs are encrypted. Encryption key is deleted so that re-identification attacks on BFs become unlikely.

Statistical operations like intersection of the BFs between two observed areas or cardinality computation of the BFs will be done on encrypted data. The results of these operations are then transmitted to the TPC that can visualize the insights about the visitors.

To generate the keying material and to execute the statistical operations on the encrypted BFs, OMO invokes the dedicated modules provided by the PAPAYA platform.

Trajectories analysis

To obtain fine-grained mobility patterns, the idea is to work directly on the probe data. Pseudonymised probe data are protected on the fly (using cryptographic techniques such as encryption or secure two-party computation). When the crypto-based protection phase is over, the keying material is deleted so that no one can have access to the data.

Once the probe data are protected, OMO performs a trajectory clustering using generic algorithm, such as k-means, or dedicated algorithms such as TRACLUS [2] and only the aggregated results will be disclosed.

Similarly, to the case of audience measurements, OMO runs an instance of the PAPAYA platform, which generates the keying material and performs the analytics operation. Specifically, OMO invokes the dedicated privacy-preserving clustering module of the PAPAYA platform which outputs the resulting clusters.



D2.1 – Use Case Specification Dissemination Level – PU

Project No. 786767

3.1.2.5 Protocol

In this section, we present the structure of the use case, including preconditions, post-conditions and flow.

In the following, we show the audience measurements protocol:

Table 10 Description of UC3 - Audience measurements

| | |
|-------------------------|---|
| ID and name | UC-Mobility-1 Audience measurements |
| Primary actor | Third-Party Customer (TPC) |
| Secondary actors | Orange Mobile Operator (OMO) running the PAPAYA platform |
| Description | A TPC needs audience analytics for some business reason. It uses OMO's service to analyse raw data probe and construct mobility indicators. |
| Preconditions | <p>PRE-1 The TPC forms an analytics request specifying the area and the period of observation.</p> <p>PRE-2 OMO registers to the PAPAYA platform.</p> <p>PRE-3 Instance of (dedicated to OMO) PAPAYA service performing statistics on BFs is running on PAPAYA platform.</p> <p>PRE-4 PAPAYA agent is running on OMO premises.</p> |
| Postconditions | POST-1 A report for the TPC, produced by OMO, is available for its download |
| Normal flow | <p>1.0 Counting phase</p> <ol style="list-style-type: none"> 1. Bloom filters are installed to capture concerned antennas' flow 2. Bloom filters are filled up during the observation period 3. When observation period is over, Bloom filters are encrypted by the dedicated module of the PAPAYA Agent. 4. Encryption key is deleted <p>2.0 Statistics phase</p> <ol style="list-style-type: none"> 1. OMO invokes the privacy-preserving PAPAYA module which computes cardinalities, unions and intersections of Bloom Filters according to Third Party Customer request 2. OMO produces and sends a report to TPC |
| Alternative flow | - |
| Exceptions | - |



D2.1 – Use Case Specification Dissemination Level – PU

Project No. 786767

| | |
|--------------------|---|
| Assumptions | - |
|--------------------|---|

In the following, we show the trajectories analysis protocol:

Table 11 Description of UC3 - Trajectories analysis

| | |
|-------------------------|--|
| ID and name | UC-Mobility-2 Trajectories analysis |
| Primary actor | Third Party Customer (TPC) |
| Secondary actors | Orange Mobile Operator (OMO) running the PAPAYA platform |
| Description | A TPC needs trajectories analytics for some business reason. It uses OMO service to analyse raw data probe and construct mobility indicators. |
| Preconditions | <p>PRE-1 The TPC forms an analytics request specifying the area and the period of observation.</p> <p>PRE-2 OMO registers to the PAPAYA platform.</p> <p>PRE-3 Instance of (dedicated to OMO) PAPAYA service performing trajectory clustering is running on PAPAYA platform.</p> <p>PRE-4 PAPAYA agent is running on OMO premises.</p> |
| Postconditions | POST-1 A report for the TPC, produced by OMO, is available for its download |
| Normal flow | <p>1.0 Data collection</p> <ol style="list-style-type: none"> 1. Pseudonymized probe data are protected on the fly (using cryptographic tools, such as encryption or secure two-party computation), via the dedicated module in the PAPAYA agent. 2. When protection phase is over, the keying material is deleted. <p>2.0 Statistics phase</p> <ol style="list-style-type: none"> 1. OMO executes the module of the PAPAYA platform dedicated to trajectories clustering, according to Third Party Customer request. |
| Alternative flow | - |
| Exceptions | - |
| Assumptions | - |



D2.1 – Use Case Specification Dissemination Level – PU

Project No. 786767

3.1.3 Privacy requirements

In order to make these protocols compliant with legal and ethical considerations, some privacy requirements must be satisfied and evaluated with a Privacy Impact Assessment (PIA). This PIA will list the possible attacks on protected data and estimate the seriousness and likelihood of each one.

In the following, we present the user requirements related to privacy.

Table 12 Privacy requirements for UC3

| Type | Description |
|------------|--|
| Data | Likelihood to re-identify a user counted in the Bloom Filters must be close to null |
| Data | Likelihood to infer information about a user counted in the Bloom Filters must be close to null |
| Data | Likelihood to single out a user in a cluster of trajectories must be close to null |
| Data | Likelihood to re-identify a user in a cluster of trajectories must be close to null |
| Data | Likelihood to infer information about a constituent user of a cluster of trajectories must be close to null (for example if, at a given moment, all users living near antenna A are going near antenna B and that antenna B is close to a place of worship, we can infer with a high probability the religion of people leaving living near A) |
| Functional | OMO has to check how to apply the right of information to respect transparency. For example by a SMS to inform data subjects. For example by a « welcome message » in the case of roaming OMO could be able to level: exercise of the right to object |

3.2 Multiple source architecture: Privacy-preserving mobile usage analytics

This section details the second use case in the mobile and phone usage umbrella, which exploits mobile phone application usage for statistical purposes (UC4).

3.2.1 Introduction and current solution limitations

The preponderance of smart devices, such as smartphones, has boosted the development and use of mobile applications (apps) in the recent years. The ubiquity of mobile usage empowers users to access information, whenever they want and wherever they are, and impacts the society by modifying everyday life habits. This prevalence also induces a large volume of mobile usage and app usage data. The analysis of these data could lead to a better understanding of users'



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

behaviors in using their smartphone, and the apps they have installed. Besides, if these data are coupled with a given context (location, time, date, sociological data, ...), the insights resulting from these analytics could be very meaningful and valuable for a wide range of stakeholders ranging from city services to retail stores, not to mention application developers or sociologists.

Orange owns a network that accommodates a large amount of data flowing through individuals and servers. Orange's ability in harnessing mobile usage data to extract meaningful insights could create value by sharing analytics results to some Third Party Customers (TPC).

However, mobile and app usage data are very sensitive, and are today considered as personal data. The collection and use of these data pose serious concerns associated with individuals' privacy. A few of recent examples of privacy breaches with respect to app usage shed light on how sensitive these data can be. For instance, some of the documents provided by Edward Snowden in 2014 revealed that the NSA has been using the mobile game Angry Birds to collect users' personal data such as age, gender and location⁵. In April 2018, the Norwegian research institute SINTEF reported that Grindr, a gay dating application, was sharing sensitive data such as HIV status (that the users disclose – or not – in their public profile) or locations to two third-party companies⁶. In early 2018, the Facebook-Cambridge Analytica data scandal hit the headline, when it was revealed that the consulting firm has harvested Facebook's data (given by users who had signed up for an application named "This is Your Digital Life") to profile US voters⁷.

To reconcile harnessing of data to get meaningful insights and privacy of users, we investigate in this use case the possibility to conduct privacy-preserving mobile data usage statistics that will prevent any inference or reidentification risks. In this use case, users will give their consent and express privacy preferences before data collection. They will then encrypt their (private and sensitive) data before sending them to the data processor.

The motivation for Orange to use the PAPAYA platform is that Orange as a mobile operator cannot share data or allow third-party consumers access to Orange's system without data protection and compliance with the regulatory obligations. The PAPAYA platform offers several cryptography-based modules that are useful and relevant for the use case of mobile usage data analytics without compromising privacy and only provides the results of this processing to third parties.

⁵ J. Ball, "Angry Birds and 'leaky' phone apps targeted by NSA and GCHQ for user data." The Guardian, January 28, 2014 <https://www.theguardian.com/world/2014/jan/27/nsa-gchq-smartphone-app-angry-birds-personal-data> [Last accessed: March 13, 2019]

⁶ "Grindr shared information about users' HIV status with third parties." The Guardian, April 3, 2018 <https://www.theguardian.com/technology/2018/apr/03/grindr-shared-information-about-users-hiv-status-with-third-parties> [Last accessed: March 13, 2019]

⁷ C. Cadwalladr and E. Graham-Harrison, "Revealed: 50 million Facebook profiles harvested for Cambridge Analytica in major data breach." The Guardian, March 17, 2018 <https://www.theguardian.com/news/2018/mar/17/cambridge-analytica-facebook-influence-us-election> [Last accessed: March 13, 2019]



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

3.2.2 Use case definition

3.2.2.1 Story

We consider a team of social scientists that are interested in having a better knowledge of the habits of users of mobile phones and in particular insights on how people consume applications such as Instagram. Traditionally, this kind of study is conducted through tools such as face-to-face interviews, diaries, questionnaires or surveys. Participants would be asked to self-report the frequency or the duration of their use of Instagram via questions such as “How frequently do you open your Instagram application? Less than a month; Once a month; Once a week; etc.”. This approach, while being widely spread in the social science research community, suffers from several drawbacks. First, participants’ answers to the questionnaire may be biased by several factors: subjectivity (people may underestimate or overestimate their use of Instagram), human memory limitations (it is hard to remember when or how often they open the application) or willingness (for fear of other’s eyes). Secondly, as a consequence of the first drawback, questions must be thoughtfully chosen and asked in order to avoid biases. Additionally, questionnaires and surveys only cover a small sample of the population. All these flaws of conventional tools for social science studies lead to ineffective implementation (right questions to ask), poor representativeness (small sample) and inaccurate results (biased answers).

Therefore, the aforementioned limitations influence our team of social scientists to study behaviors of Instagram users at source, based on their mobile usage data, that consist of timestamps of application opening and closing, actions performed in the application (such as “like” button), configuration parameters, etc. Collection and analytics of such data present several advantages compared to traditional tools: mobile usage data analytics deal with objective data that cannot be biased, cover a bigger and more scalable sample and result in more fine-grained and more accurate analytics. Besides, the usage data can be enriched with additional data such as phone data (OS, used network, i.e. 3G, 4G or Wi-Fi, ...) and demographic data (age, size of family, ...) to obtain even more fine-grained insights. Hence mobile usage data analytics either can be used in combination with traditional tools (surveys, questionnaires) or can even potentially supersede questionnaire-based methods.

Literature in social research [3, 4, 5] gives us possible examples of Instagram usage studies that can leverage mobile usage data analytics:

- study the divergence of Instagram usage between different generations (baby-boomers, generation X, millennials, etc.),
- study the psychological impact of Instagram usage (e.g., eating disorder, addiction, self-confidence, etc.) in function of age and geographical location of users.

The team of social scientists is aware that Orange provides a service for privacy-preserving analytics on mobile usage data collected from registered users. To perform such a collection, Orange provides to users an application to be installed on their mobile phones which will take



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

care of the collection (let us call this application the Orange Monitoring App, or OMA). Besides, to encourage users to participate to the study, Orange (or in another protocol the TPC) will offer incentives as rewards (discounts, vouchers, etc.).

In this context, the team of social scientists contacts Orange with a specific analytics request. For instance, they will request the average number of times Instagram is used per period of 1 day per age of users or the average duration of use of Instagram. The analytics requests specify the attributes on which they would like the analytics operations to be performed (for instance, age, day of the week, month, duration of use etc.) as well as the operation to be performed on the data. Orange then performs the requested analytics on the aggregated data collected during the specified observation period and sends back the analytics results to the research team. Having the average number of times Instagram is used or the average duration of usage of Instagram (or any other insights), the research team can derive social conclusions about Instagram usage.

As the usage data could be highly sensitive in terms of users' privacy, the previously described use case will include mechanisms for privacy protection, in accordance with the GDPR. Consent is a key condition to collect and process users' data (Articles 6, 7 and 8 of the GDPR). Therefore, before each new study, users are requested to give their (informed and explicit) consent on the collection and processing of their data by Orange. They should agree on (i) the attributes that will be collected from them; (ii) the duration of the collection; (iii) the purpose of the processing (the requested analytics) that will be performed on the collected data and (iv) the period of storage of their data. The users are also empowered with the ability to specify privacy preferences on collection and usage of their data. On the other hand, as mentioned earlier, to incentivize users to share their data, they could be offered rewards from Orange.

Furthermore, to fulfil the principle of privacy-by-design requested by the GDPR (Articles 25 and 32), the present use case will leverage technical solutions for data protection. Before the usage data leave the users' mobile phone, they will be encrypted using some specific encryption key defined for the study. Orange will perform the requested analytical operation on the encrypted collected data, coming from different users, hence without having a plain access to the content of the data. The results of this operation are then shared in unencrypted form with the research team. In this use case, we will focus on specific analytics operations, namely simple summary statistics such as computing the mean.

3.2.2.2 Players

Table 13 UC4 Players

| | |
|-----------------------------------|---|
| Orange Network Users (ONU) | <p>The 'data subject' in GDPR perspective.</p> <p>They use their mobile phones and the applications installed on them (such as Instagram), thus generating mobile usage & app usage data.</p> |
|-----------------------------------|---|



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

| | |
|-----------------------------------|--|
| | <p>To participate to the statistics study, they give their consent to Orange and install the OMA application that will monitor their usage of applications such as Instagram.</p> <p>They aggregate, enrich (with phone and demographic data) and then encrypt their usage data using a specific encryption key.</p> <p>The ONUs run an instance of the client side PAPAYA platform, which is responsible for key generation and data encryption.</p> <p>They receive rewards for sharing their usage data.</p> |
| Orange | <p>Orange is the mobile operator.</p> <p>It receives requests for analytics on usage data from the TPC and sends them back the results of the operations (B2B relationship). Orange collects encrypted data from users and obviously aggregates them. Then, Orange performs the requested analytics (statistics operations) on the aggregated encrypted data.</p> <p>To perform the requested analytics, Orange runs an instance of the PAPAYA statistics service.</p> <p>Its GDPR role is Data Owner, Data Controller and Data Processor.</p> |
| Third-party Customer (TPC) | <p>It is interested in mobile usage data analytics results for its business or studies.</p> <p>It expresses an analytics request to Orange which sends back the data analytics results without revealing the data that have been used to perform the analytics or compromising users' privacy.</p> <p>It has no GDPR role.</p> |

3.2.2.3 Involved data

The data involved in the processing are the attributes related to users' usage of their mobile phones and applications. These include, among other attributes:

- Usage data:
 - Identifiers: phone numbers, location
 - Other: timestamps (start/end time of app usage), mobile OS, configuration parameters (such as language, private/public account), network (Wi-Fi, 3G, 4G, ...), app actions (open, close, like, post, update, etc.), device sensors (accelerometers, microphone), etc.
- Sociodemographic data:



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

- Quasi-identifiers: age (age range or precise age), gender, relationship status, education level, income range, etc.

Users participating in the study are required to install an app, the Orange Monitoring App, which will perform the data collection. Through this application, users receive information about new studies and give their consent to the sharing of their usage data. Consent will be asked for each study, hence for each data collection. In the data collection process, the app will consider (i) the number of users that are invited to participate to the analytics, (ii) the number of users that positively respond to the invitation and (iii) the number of users that actually share their data. Only data of users from (iii) will be collected. The data collected through the Orange Monitoring App is first locally aggregated, in accordance with user preferences, and then encrypted before leaving the phone (via the PAPAYA client side platform).

3.2.2.4 Processing and analytics

The analytics are operated by Orange, as Data Processor, based on a request issued by the TPC. These operations will be performed on encrypted data, aggregated from multiple sources (multiple users). Neither Orange, nor the TPC will have access to the plain usage data. The scenario will focus only on statistical operations that can be expressed in terms of inner products or quadratic functions (such as mean, sum, variance, range, etc.). More elaborated operations such as linear regression could also be investigated.

3.2.2.5 Protocol

Table 14 Description of UC4

| | |
|-------------------------|---|
| ID and name | UC-MobileUsage-1 Mobile app usage statistics |
| Primary actor | Third-Party Customer (TPC) |
| Secondary actors | Orange, Orange Monitoring App (OMA) |
| Description | A TPC is interested in obtaining insights on mobile usage. It requests Orange to collect and analyse mobile usage data from users that consent to participate to a study and to share their data. Orange performs the requested analytics (basic statistics) and returns the result to the TPC. |
| Preconditions | <p>PRE-1 Orange registers to the PAPAYA platform.</p> <p>PRE-2 Instance of (dedicated to Orange) PAPAYA service performing statistics is running on PAPAYA platform.</p> <p>PRE-3 Orange and the TPC define a business relationship between each other where Orange provides an analytics service to the TPC.</p> <p>PRE-4 The TPC forms an analytics request specifying the period of observation, the attributes to collect and the type of statistics.</p> |



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

| | |
|------------------------------|--|
| | <p>PRE-5 Orange validates the technical feasibility of the requested analytics.</p> <p>PRE-6 The users install the OMA, which integrates an instance of the PAPAYA client side platform.</p> <p>PRE-7 The users give their consent to the collection and processing of their usage data.</p> |
| <p>Postconditions</p> | <p>POST-1 A report for the TPC produced by Orange, is available for its download.</p> |
| <p>Normal flow</p> | <p>1.0 Initialisation</p> <ol style="list-style-type: none"> 1. TPC sends to Orange an analytics request that specifies the period of observation, the attributes to collect and the statistics to compute. 2. Orange studies the technical feasibility of the analytics request, depending on the attributes, the analytics, and the available encryption mechanisms, in accordance with operative legal requirements. If the conclusion is that this is not possible, the process is stopped. 3. Orange sends an invitation and a consent request to a panel of users to participate to TPC’s study by forwarding the analytics request (via the OMA). 4. Users respond to the invitation and give their consent to the collection and purpose of processing. 5. Users respond to a questionnaire prepared by the Privacy Engine (embedded in the OMA). 6. The Privacy Engine extracts privacy preferences from the users’ answers to the questionnaire. 7. The OMA generates keying material for the user, by calling the dedicated module in the PAPAYA client side agent (embedded in the OMA). <p>2.0 Statistics phase</p> <ol style="list-style-type: none"> 1. During the observation period specified in the analytics request, the OMA collects the data (the attributes) listed in the request. 2. OMA performs a local aggregation of the collected data. 3. Aggregated data is enriched with other kind of data (phone and sociodemographic data). 4. OMA calls the PAPAYA client side agent to encrypt the enriched aggregated data and sends it to Orange. |



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

| | |
|-------------------------|--|
| | <ol style="list-style-type: none"> 5. Orange aggregates data received from the users actually participating to the study. 6. Orange invokes the dedicated module of the PAPAYA platform which performs the statistics operation specified in the query. 7. Orange sends the results to the TPC. |
| Alternative flow | - |
| Exceptions | <p>1 E1 No consent received from the user for the collection and processing of her usage data for the statistics purpose</p> <ol style="list-style-type: none"> 1. TPC sends to Orange an analytics request that specifies the period of observation, the attributes to collect and the statistics to compute. 2. Orange sends an invitation and a consent request to a panel of users to participate to TPC’s study by forwarding the analytics request (via the OMA). 3. Users respond to the invitation. One or more users DO NOT give their consent to the collection and purpose of processing. 4. OMA does not collect data from the unenrolled users. <p>2 E2 Orange deems that the requested analytics is not feasible regarding the available cryptographic mechanisms.</p> <ol style="list-style-type: none"> 1. TPC sends to Orange an analytics request that specifies the period of observation, the attributes to collect and the statistics to compute. 2. Orange studies the technical feasibility of the analytics request, depending on the attributes, the analytics, and the available encryption mechanisms, in accordance with operative legal requirements. The conclusion of this assessment is that the requested analytics is not possible. 3. Orange sends this conclusion to TPC and stops the process. |
| Assumptions | N/A |

3.2.3 Privacy requirements

This privacy-preserving mobile usage analytics use case raises several privacy concerns. In this use case, the users are requested to give their consent to the collection and processing of their data for a specific study, hence a specific statistical analysis. Therefore, no data will be collected or processed from users who disagree to give consent.



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

In case consent is given, as collected information is privacy-sensitive, they should be protected with regard to the current legislation. In this use case, we shall implement a functional encryption algorithm (see deliverable D3.1 on Preliminary Design of Privacy-Preserving Data Analytics) that, beyond protecting data confidentiality, allows Orange to derive some partial information (sum, average, etc.) of the encrypted aggregated data. Hence, one of the privacy requirements is that the possibility to decrypt each single user's data should be impossible. Similarly, it must be impossible for Orange to compute the analytics operations on a dataset originating from a single user; otherwise, this user will straightforwardly be re-identified. This implies that analytics must not be performed before multi-source data are aggregated. Finally, as consent is given for one study or one data processing, it must be infeasible to perform any other analytics on the data that have been collected for other purposes.



Project No. 786767

4 Threat detection use case

In addition to the four use cases described so far, we intend to investigate another use case (UC5) which, to our opinion, is relevant for PAPAYA. Unlike the previously described story lines, this use case does not necessarily apply to personal data that fall under the GDPR legislation, but focuses on business-sensitive data, whose confidentiality is of paramount importance for companies. Hence, the threat detection use case tackles the problem of detecting threats in systems or networks via dedicated analytics algorithms while respecting the confidentiality of the data used during detection.

4.1 Introduction and current solution limitations

Nowadays, organizations are increasingly interconnected to each other and thanks to paradigms such as Internet, mobile devices or cloud computing, they find new business opportunities in this ubiquitous model. However, this omnipresence in the cyberspace also attracts the attention of smarter and more threatening cybercriminals than before. In this threat landscape, machine learning can help detect advanced threats and contribute to fight cybercrime. Indeed, by collecting data from systems and networks and then by applying dedicated machine learning algorithms for threat detection, organizations can be more prepared to thwart attacks and automate their responses to incidents.

One of the traditional approaches for threat detection is malicious activity recognition, such as signature-based detection. The idea of this technique relies on the existence of a dataset of malicious attack signatures. Incidents are detected by matching the signature of the tested object (network traffic packets or system call sequences, for example) to one of the signatures in this dataset. Even though this approach is the backbone of many current security services, it suffers from several drawbacks: high rate of false alarms and impossibility to detect zero-day attacks (attacks whose signature is not yet in the dataset). Fortunately, another technique for threat detection is gaining interest in the cybersecurity community thanks to its promising results in zero-day attack detection: anomaly detection (also called outlier detection). In a nutshell, anomaly detection identifies observations (events, transactions, data etc.) that deviate from an expected pattern (the *normal* behavior). If a previously-unseen item tested by the detection algorithm falls outside the predefined normal pattern, then this item is classified as an anomaly. This wording might imply that anomaly detection is an easy task. However, practically, it is not as simple as it sounds. Indeed, the main challenge in this approach is how to model the *normal* behavior. There exist a couple of techniques that attempt to solve this challenge⁸:

- unsupervised anomaly detection algorithms such as isolation forests [6, 7] detect outliers in an unlabeled data set by isolating items that are rare and different from the normal items that appear in “dense” regions. Isolation forests consist of a collection of decision trees whose branches are randomly partitioned. Because of this random partitioning, outliers

⁸ See also <http://cucis.ece.northwestern.edu/projects/DMS/publications/AnomalyDetection.pdf>



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

stand closer to the root of the trees than normal items. Besides, compared to the two following techniques, unsupervised anomaly detection algorithms do not need to profile the normal behavior.

- supervised anomaly detection algorithms require a labelled training data set in which items are either labelled as “normal” or “anomalous” (which assumes the existence and availability of such a data set). The labelled data are used to train a classifier, such as a neural network. Hence, the model profiles both normal and abnormal items. Each new input tested with a supervised anomaly detection algorithm is classified as either normal or anomalous. Note that this kind of approach has a major drawback: normal data might be overrepresented in the training data set, since anomalies are, by definition, scarce.
- hybrid (semi-supervised) anomaly detection algorithms consider training data that are labelled only for the normal class. Hence, the model captures only the normal behavior and scores the probability that a new item belongs to this class. Anomalous items are those elements that the probability score is very low (or below some predefined threshold).

Therefore, defining the normal behavior is dependent on the available data. An approach [8, 9] to improve the model of normal behavior (in the case of (semi-) supervised learning) or to better detect outliers (in the case of unsupervised learning) makes different organizations combine their data sets, which are then submitted to and processed by a central service. Indeed, first, a normal behavior for a company A may be different from the one of a company B. Hence, by sharing their data, normal behavior will be better profiled thanks to their combined data sets. Moreover, having a better view of what normality means, outliers will be even more isolated in the case of combined data. To illustrate this approach, the authors in [9] called the central service the “detective”. He is in charge of detecting misbehavior. He gathers information from “witnesses”, e.g. the organizations that mutualize their data. These witnesses alone do not have the capability to detect whether a particular behavior is normal or suspect. Nevertheless, their data can be shared with the detective who is savvier than the witnesses and who can profile a global model for normality from the joint data sets and obtain a broader picture of the possible anomalies.

However, the organizations (the “witnesses”) may be reluctant to disclose their data sets, for privacy, confidentiality and competitive reasons. Hence, to incentivize companies to contribute to the mutualized data set, privacy and confidentiality of the sensitive data against the central service, the other organizations and external attackers must be ensured by means of cryptographic protection. In the same time, if the organizations eventually provide their protected data to the joint data set, they would like to be ensured that protection does not sacrifice the accuracy of the anomaly detection algorithm implemented by the central service. Besides, the latter does not want to share the details of its algorithm, since it represents its core business asset.

Another aspect is that, once the model has been defined, the detective may also want to protect his expertise in machine learning algorithms for threat detection so that he will not be prone to



D2.1 – Use Case Specification Dissemination Level – PU

Project No. 786767

send his model and algorithm to witnesses. Then, if one witness wants to later test his own observations with the output model, there will be a mutual freeze since neither the detective nor the witness will want to send its own knowledge to the other party.

This threat detection use case will allow us to investigate the possibility to use advanced cryptographic solutions such as homomorphic encryption or secure multiparty computation to incite organizations to share their data, since they would be guaranteed that their sensitive information is kept secret to other parties, while being still useful for the detection purposes.

4.2 Use case definition

4.2.1 Story

Orange plays the role of the central entity: it offers a service for threat detection to its business clients. As a matter of illustration, we consider that the service implements a supervised learning algorithm. Let us assume that companies A, B, C and D are these customer companies. The Orange service operates in two phases: research and commercial phases.

Research phase. In this phase, Orange wants to train its anomaly detection algorithm based on partners' data. Without loss of generality, companies A, B and C are contributing customers (CC) (i.e. the “witnesses” in the analogy sketched above). In other terms, these three companies submit their data to Orange's service in order to train its machine learning algorithm to profile the normal behavior. For secrecy reasons, Orange does not want to disclose its algorithm and its parameters to the CCs. In the same time, the CCs are reluctant to share their data with Orange. Hence, both Orange and the CCs will cryptographically protect their assets.

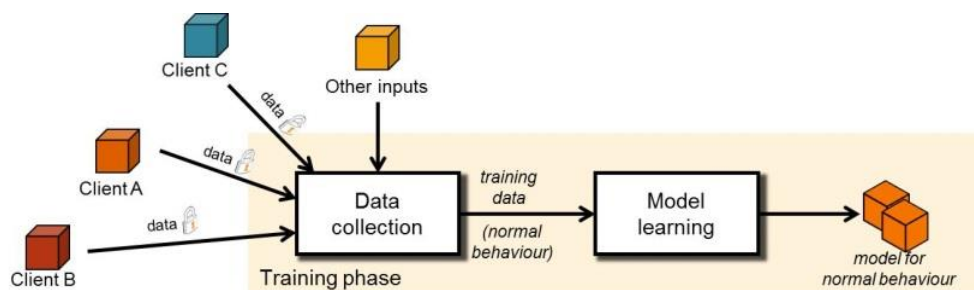


Figure 1 UC5 research phase

Commercial phase. Later on, company A or D (being part of so-called Third-Party Customers, TPC) will pay the service of Orange to detect some potential threats in its IT system or network. To do so, A (resp. D) provides some new (i.e., unseen before) data sets to be tested by Orange's algorithm. Orange applies the model learned in the previous phase and returns the anomaly score to the TPC. Based on this score, the TPC makes a decision about the potentially detected threats.



D2.1 – Use Case Specification Dissemination Level – PU

Project No. 786767

As in the research phase, the TPC will cryptographically protect its data set to preserve their confidentiality.

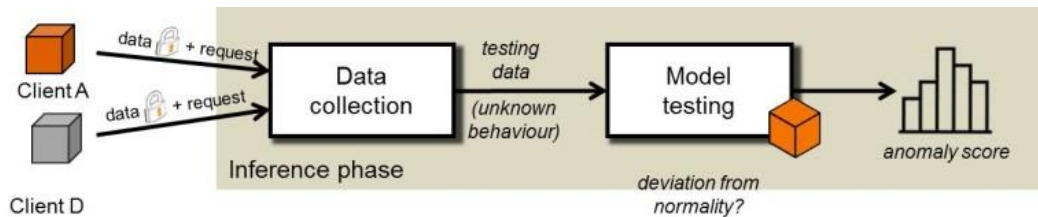


Figure 2 UC5 commercial phase

The PAPAYA will provide different kind of services that can be employed by Orange for the purpose of this threat detection use case: (1) one service for the research phase, that is used by Orange to train its machine learning model; and (2) another service for the commercial phase that is invoked by Orange to apply the learned model on CC data. In case that there is only a commercial phase (i.e. the underlying machine learning algorithm is unsupervised, see section 4.2.4), the PAPAYA platform will also provide a dedicated module, which again will be applied by Orange for the purpose of secure threat detection.

4.2.2 Players

Table 15 UC5 Players

| | |
|--|--|
| <p>Orange</p> | <p>Orange provides a paid service for threat detection to TPCs.</p> <p>This service implements an anomaly-based threat detection algorithm, which is trained with data provided by the Contributing Customers (or CCs, see next row). In the research phase, Orange obviously trains its algorithm, without having access to the plain data provided by the CCs.</p> <p>In the commercial phase, Orange applies the model learned in the research phase to the data set submitted by the TPCs.</p> <p>Orange runs an instance of the PAPAYA platform which integrates dedicated module for oblivious training and inference over encrypted data.</p> |
| <p>Contributing Customer (CC)</p> | <p>The CCs are business partners of Orange.</p> <p>They contractually provide data to Orange to improve its anomaly-based detection algorithm.</p> <p>They protect their data sets with advanced cryptographic solutions.</p> |



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

| | |
|-----------------------------------|---|
| | The CCs run an instance of the PAPAYA client side agent platform, which offers a dedicated module for data protection. |
| Third-Party Customer (TPC) | <p>The TPC consumes Orange’s threat detection services.</p> <p>They suspect anomalous behavior in their IT systems and requests Orange to confirm the anomalies.</p> <p>They provide Orange with data to be tested, and protect those using advanced cryptographic solutions.</p> <p>The TPCs get from Orange the output of the thread detection algorithms and make a decision about the potentially detected threats.</p> <p>A TPC can play the role of a CC and a CC can play the role of a TPC.</p> |

4.2.3 Involved data

The data that are shared by the parties include one of the following items or a combination of them:

- Network traffic data (network configuration, DNS logs, network packets, IP addresses, port numbers, etc.)
- Security-related events (access controls logs, access credentials, firewall logs, etc.)
- Web history (accessed URLs, user info, browser info, etc.)

These data are considered as sensitive since they contain information about the system and the procedures involved in the normal operational workflow of the organizations. Some of the data such as access credentials or user info are personal data and fall under the GDPR umbrella. This point will not be treated within the PAPAYA project.

4.2.4 Processing and analytics

At the time of writing this deliverable, we consider two possible anomaly-based threat detection algorithms.

- Unsupervised learning: isolation forests. This algorithm has recently gained popularity in the context of anomaly detection. As it is an unsupervised algorithm, it does not require a training phase, and hence no labelled training data is needed. The basic operation in isolation forests is comparisons of data value with a threshold. These comparisons are computed while building the trees that form the forest. At the end of the tree building, each level in the trees is mapped to an anomaly score: the closer to the root, the higher the anomaly score (i.e. the more probable the item found in that level is an anomaly).



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

- Supervised learning: neural networks. This algorithm requires a labelled training data set to train the model of the neural network (this is the training phase). This model will be able to classify any submitted item into normal or abnormal classes (this is the classification phase).

While designing the cryptographic solution to privacy-preserving anomaly detection, one has to consider the intrinsic operations of the algorithm. Indeed, the choice of the cryptographic building blocks depends on the nature of these operations and the feasibility of the building block to support the operations efficiently. In WP3, we will consider three candidate cryptographic tools for privacy-preserving anomaly detection, namely fully homomorphic encryption, functional encryption and secure multiparty computation. These three cryptographic tools are defined in details in deliverable D3.1 and present pros and cons. They can be combined to provide an efficient and accurate solution. These considerations are investigated in WP3. The PAPAYA platform will integrate a privacy-preserving anomaly detection module which will implement the selected cryptographic tool.

4.2.5 Protocol

Table 16 Description of UC5

| | |
|-------------------------|---|
| ID and name | UC-TD-1 Threat Detection |
| Primary actor | Orange |
| Secondary actors | The contributing customers (CCs) and the third-party customers (TPC). |
| Description | <p>A TPC would like to check whether some malicious factors might threaten its system and normal operations. The company uses the Orange service for threat detection which is based on anomaly detection.</p> <p>If a supervised algorithm is used, Orange trains its model based on CCs' data. This model is then used to detect whether the requesting TPC's data are subject to threats.</p> <p>If an unsupervised algorithm is used, Orange directly performs the anomaly detection on the TPC's data, without any training phase.</p> |
| Preconditions | <p>PRE-1 TPCs and Orange have a contractual business relationship in which Orange offers a service and the TPCs pay for this service to obtain threat detection insights.</p> <p>PRE-2 CCs and Orange have a contractual business relationship in which Orange improves its anomaly detection model based on CCs' data.</p> |



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

| | |
|-----------------------|---|
| | <p>PRE-3 Orange runs an instance of the PAPAYA service for threat detection on the PAPAYA platform.</p> <p>PRE-4 The CCs (and possibly the TPC run an instance of the PAPAYA client side agent).</p> |
| Postconditions | <p>POST-1 Orange prepares a report that gives the results of the threat detection procedure and sends it to the requesting TPC.</p> |
| Normal flow | <p>Case A: Unsupervised learning</p> <ol style="list-style-type: none"> 1. Detection phase 1. TPC encrypts its data using its own encryption key (via the PAPAYA client side agent). 2. TPC sends the encrypted data to Orange's service and a request for threat detection. 3. Orange applies the unsupervised anomaly detection algorithm (isolation forests) on the encrypted data, using the dedicated module from the PAPAYA platform. 4. Orange sends the anomaly detection results to TPC. <p>Case B: Supervised learning</p> <ol style="list-style-type: none"> 1. Research phase 1. CCs encrypt their data using their own respective encryption keys (via the PAPAYA client side agent). 2. CCs send the encrypted data to Orange. 3. Orange trains its anomaly detection model (neural networks), by invoking the dedicated module from the PAPAYA platform. 4. At the end of the research phase, Orange keeps its improved model private. 2. Commercial phase 1. TPC encrypts its data using its own encryption key (via the PAPAYA client side agent). 2. TPC sends the encrypted data to Orange service and a request for threat detection. |



D2.1 – Use Case Specification

Dissemination Level – PU

Project No. 786767

| | |
|-------------------------|--|
| | <ol style="list-style-type: none"> 3. Orange applies the learned model on the encrypted data, by invoking the dedicated module from the PAPAYA platform. 4. Orange sends the anomaly detection results to TPC. |
| Alternative flow | - |
| Exceptions | - |
| Assumptions | - |

4.3 Data sensitivity requirements

During the research phase, the PAPAYA system should not permit Orange to obtain any information about CCs data. Similarly, it should be impossible for the participating CCs to obtain any information about Orange’s algorithm and parameters.

During the commercial phase, the PAPAYA system should not permit Orange to obtain any information about TPC data. Similarly, it should be impossible for the requesting TCP to obtain any information about Orange’s algorithm and model. The anomaly detection results consist of anomaly scores that do not leak much information about the data. Hence, we allow Orange to see these scores in the clear.



Project No. 786767

5 Conclusions

In this document, we presented the use cases' definitions of the healthcare umbrella and the ones of the mobile and phone usage umbrella. This work fulfills the objectives of task T2.1 (providing valuable use cases for final users) and complements with the outcome of tasks T2.2 (that provides user requirements) and T2.3 (that provides platform requirements).

The consortium identifies:

- two use cases under the healthcare umbrella, where personal data are processed so as to detect arrhythmias in ECG signals (UC1) and stress situations in workers (UC2). Both use cases use neural networks as for the analytics operation. UC1 mainly focuses on the classification phase for arrhythmia detection, whereas UC2 targets the training phase whereby the neural network model is computed by multiple sources, jointly;
- two use cases under the mobile and phone usage umbrella, where personal data are processed so as to extract mobility analytics (UC3) and mobile usage analytics (UC4). The underlying analytics operations mainly consists of statistical operations such as counting and clustering algorithms. In this umbrella, an additional player, namely a third party consumer who mainly has access to the results of the analytics, is provided.
- These four use cases are complemented with an additional use case (UC5), under the mobile and phone usage umbrella, where **business-sensitive data** rather than personal data are processed, so as to detect threats in systems and networks based on anomaly detection algorithms.

The definition of the use cases presented in this document was performed in accordance with the directives dictated by the involved stakeholders and end users, via interviews and cyclic reviews of the work done. In this way, the definition of use cases is coherent with the expectations of the key players in these sectors, adding a tangible value to the work performed in PAPAYA.

The outcome of this document provides an entry point for the design and implementation of the PAPAYA platform (as of WP3 and WP4) and a basis for the validation of the work done in PAPAYA (as of WP5).



Project No. 786767

6 References

- [1] B. H. Bloom, «Space/time trade-offs in hash coding with allowable errors,» Communications of the ACM, vol. 13, n° 17, pp. 422--426, 1970.
- [2] J.-G. Lee, J. Han et K.-Y. Whang, «Trajectory clustering: a partition-and-group framework,» chez Proceedings of the 2007 ACM SIGMOD international conference on Management of data, Beijing, China, 2007.
- [3] K. J. Fietkiewicz, E. Lins, K. S. Baran et W. G. Stock, «Inter-generational comparison of social media use: Investigating the online behavior of different generational cohorts,» chez 49th Hawaii International Conference on System Sciences (HICSS), Koloa, Hawaii, USA, 2016.
- [4] C.-c. Yang, «Instagram use, loneliness, and social comparison orientation: interact and browse on social media, but don't compare,» Cyberpsychology, Behavior, and Social Networking, vol. 19, n° 112, pp. 703--708, 2016.
- [5] K. Lup, L. Trub et L. Rosenthal, «Instagram #instasad?: exploring associations among instagram use, depressive symptoms, negative social comparison, and strangers followed,» Cyberpsychology, Behavior, and Social Networking, vol. 18, n° 15, pp. 247--252, 2015.
- [6] F. T. Liu, K. M. Ting et Z.-H. Zhou, «Isolation forest,» chez 2008 Eighth IEEE International Conference on Data Mining, Pisa, Italy, 2008.
- [7] F. T. Liu, K. M. Ting et Z.-H. Zhou, «Isolation-based anomaly detection,» ACM Transactions on Knowledge Discovery from Data , vol. 6, n° 11, 2012.
- [8] H. A. Ringberg, «Privacy-Preserving Collaborative Anomaly Detection,» PhD. Thesis, 2009.
- [9] M. Allman, E. Blanton, V. Paxson et S. Shenker, «Fighting coordinated attackers with cross-organizational information sharing,» chez Fifth Workshop on Hot Topics in Networks, Irvine, USA, 2006.