**Autonomous Vehicles to Evolve to a New Urban Experience**

---

## DELIVERABLE

## D4.5 Second iteration In-vehicle services

# Disclaimer

This document reflects only the author's view and the European Commission is not responsible for any use that may be made of the information it contains.

# Document Information

| Grant Agreement Number | 769033 |
|---|---|
| Full Title | Autonomous Vehicles to Evolve to a New Urban Experience |
| Acronym | AVENUE |
| Deliverable | D4.5 Second iteration in-vehicle services |
| Due Date | 01.05.2020 |
| Work Package | WP4 |
| Lead Partner | Autonomous Mobility |
| Leading Author | Christian Zinckernagel |
| Dissemination Level | Public, restricted |

# Document History

| Version | Date | Author | Description of change |
|---|---|---|---|
| 0.1 | 14.04.2020 | Christian Zinckernagel, AM | First draft of D4.5 |
| 0.2 | 24.04.2020 | Antonios Lalas, CEERTH | First draft of D4.5 |
| 0.3 | 27.04.2020 | Christian Zinckernagel, AM | Second draft of D4.5 |
| 0.4 | 29.04.2020 | Antonios Lalas, CEERTH | Second draft of D4.5 |
| 0.5 | 04.05.2020 | Dimitrios Tsiktsiris, Anastasios Vafeiadis, Athanasios Papadakis, Nikolaos | Service descriptions, contributions of section 2.4 |

| | | Dimitriou, Antonios Lalas, CERTH | |
|---|---|---|---|
| 0.6 | 07.05.2020 | Christian Zinckernagel, AM | Service descriptions |
| 0.7 | 06.05.2020 | Dimitrios Tsiktsiris, Maria-Eleni Kadoglou, Leon Vitanos, Nikolaos Dimitriou, Antonios Lalas, CERTH | Contribution of section 2.5-2.6 and contributions to section 2.7 |
| 0.8 | 07.05.2020 | Antonios Lalas, Konstantinos Votis, CERTH | Third draft of D4.5 |
| 0.9 | 11.05.2020 | Christian Zinckernagel, AM | Final draft of D4.5 |
| 0.10 | 20.05.2020 | Kevin Salvi, MobileThinking | Internal WP4 review |
| 0.11 | 24.05.2020 | Christian Zinckernagel, AM | Internal WP4 corrections |
| 0.12 | 09.06.2020 | Clement val, Ceesar | Offical AVENUE review |

# Table of Contents

# List of figures

# Acronyms

| | |
|---|---|
| ADS | Automated Driving Systems |
| AI | Artificial Intelligence |
| API | Application Protocol Interface |
| AV | Autonomous Vehicle |
| BMM | Business Modelling Manager |
| CAV | Connected and Autonomous Vehicles |
| CB | Consortium Body |
| CERN | European Organization for Nuclear Research |
| D7.1 | Deliverable 7.1 |
| DC | Demonstration Coordinator |
| DI | The department of infrastructure |
| DMP | Data Management Plan |
| DSES | Department of Security and Economy Traffic Police |
| DTU test track | Technical University of Denmark test track |
| EAB | External Advisory Board |
| EC | European Commission |
| EC | European Commission |
| ECSEL | Electronic Components and Systems for European Leadership |
| EM | Exploitation Manager |
| EU | European Union |
| EUCAD | European Conference on Connected and Automated Driving |
| F2F | Face to face meeting |
| FEDRO | Federal Roads Office |
| FEDRO | (Swiss) Federal Roads Office |
| FOT | (Swiss) Federal Office of Transport |
| GDPR | General Data Protection Regulation |
| GIMS | Geneva International Motor Show |
| GNSS | Global Navigation Satellite System |

| | |
|---|---|
| HARA | Hazard Analysis and Risk Assessment |
| IPR | Intellectual Property Rights |
| IT | Information Technology |
| ITU | International Telecommunications Union |
| LA | Leading Author |
| MEM | Monitoring and Evaluation Manager |
| OCT | General Transport Directorate of the Canton of Geneva |
| ODD | Operational Domain Design |
| OEDR | Object And Event Detection And Response |
| OFCOM | Federal Office of Communications |
| PC | Project Coordinator |
| PEB | Project Executive Board |
| PGA | Project General Assembly |
| PRM | Persons with Reduced Mobility |
| PSA | Group PSA (PSA Peugeot Citroën) |
| PTO | Public Transportation Operator |
| PTO | Public Transport Operator |
| PTS | Public Transportation Services |
| QRM | Quality and Risk Manager |
| QRMB | Quality and Risk Management Board |
| RN | Risk Number |
| SA | Scientific Advisor |
| SAE Level | Society of Automotive Engineers Level (Vehicle Autonomy Level) |
| SAN | Cantonal Vehicle Service |
| SDK | Software Development Kit |
| SMB | Site Management Board |
| SoA | State of the Art |
| SOTIF | Safety Of The Intended Functionality |

| SWOT | Strengths, Weaknesses, Opportunities, and Threats. | WP | Work Package |
| TM | Technical Manager | WPL | Work Package Leader |
| UITP | Union Internationale des Transports Publics | | |

# Executive Summary

This deliverable will introduce the preparation work, done by Amoblity (AM) and CERTH, for prototyping and testing 4 in-vehicle services in the AVENUE project. For each chosen service the technological framework is described, including the maturity of the technology and the equipment needed to prototype the service elements. Furthermore a prototyping plan will be introduced for each service and the result framework will be defined.

The five chosen services that will be introduced and tested in the next 8-12 months are:

- Security trust services (example: Prevention of night aggressions)
- Automated passenger presence
- Follow your kid/grandparents
- Shuttle environment assessment
- Smart feedback system

As a part of the AVENUE goals are to be able to drive fully autonomous without the presence of a safety driver, the role of the safety driver is described. The chosen services all target some of the given services that the safety driver offers in person. The five services are all services that need to be automated in order to be able to remove the operator.

# 1 Introduction

AVENUE aims to design and carry out full-scale demonstrations of urban transport automation by deploying, for the first time worldwide, fleets of autonomous minibuses in low to medium demand areas of 4 European demonstrator cities (Geneva, Lyon, Copenhagen and Luxembourg) and 2 to 3 replicator cities. The AVENUE vision for future public transport in urban and suburban areas, is that autonomous vehicles will ensure safe, rapid, economic, sustainable and personalised transport of passengers. AVENUE introduces disruptive public transportation paradigms on the basis of on-demand, door-to-door services, aiming to set up a new model of public transportation, by revisiting the offered public transportation services, and aiming to suppress prescheduled fixed bus itineraries.

Vehicle services that substantially enhance the passenger experience as well as the overall quality and value of the service will be introduced, also targeting elderly people, people with disabilities and vulnerable users. Road behaviour, security of the autonomous vehicles and passengers' safety are central points of the AVENUE project.

At the end of the AVENUE project four year period the mission is to have demonstrated that autonomous vehicles will become the future solution for public transport. The AVENUE project will demonstrate the economic, environmental and social potential of autonomous vehicles for both companies and public commuters  while assessing the vehicle road behaviour safety.

## 1.1  On-demand Mobility

Public transportation is a key element of a region's economic development and the quality of life of its citizens.

Governments around the world are defining strategies for the development of efficient public transport based on different criteria of importance to their regions, such as topography, citizens' needs, social and economic barriers, environmental concerns and historical development. However, new technologies, modes of transport and services are appearing, which seem very promising to the support of regional strategies for the development of public transport.

On-demand transport is a public transport service that only works when a reservation has been recorded and will be a relevant solution where the demand for transport is diffuse and regular transport  is inefficient.

On-demand transport differs from other public transport services in that vehicles do not follow a fixed route and do not use a predefined timetable. Unlike taxis, on-demand public transport is usually also not individual. An operator or an automated system takes care of the booking, planning and organization.

It is recognized that the use and integration of on-demand autonomous vehicles has the potential to significantly improve services and provide solutions to many of the problems encountered today in the development of sustainable and efficient public transport.

## 1.2  Autonomous Vehicles

A self-driving car, referred in the AVENUE project as **an Autonomous Vehicle** (**AV**) is a vehicle that is capable of sensing its environment and moving safely with no human input.  The choice of Autonomous vs Automated was made in AVENUE since, in the current literature, most of the vehicle concepts have a

person in the driver's seat, utilize a communication connection to the Cloud or other vehicles, and do not independently select either destinations or routes for reaching them, thus being "automated". The automated vehicles are considered to provide assistance (at various levels) to the driver. In AVENUE there will be no driver (so no assistance will be needed), while the route and destinations will be defined autonomously (by the fleet management system). The target is to reach a system comprising of vehicles and services that independently select and optimize their destination and routes, based on the passenger demands.

In relation to the SAE levels, the AVENUE project will operate SAE Level 4 vehicles.



©2020 SAE International

## 1.2.1    Autonomous vehicle operation overview

We distinguish in AVENUE two levels of control of the AV: micro-navigation and macro-navigation. Micro navigation is fully integrated in the vehicle and implements the road behaviour of the vehicle, while macro-navigation is controlled by the operator running the vehicle and defines the destination and path of the vehicle, as defined the higher view of the overall fleet management.

For micro-navigation Autonomous Vehicles combine a variety of sensors to perceive their surroundings, such as 3D video, lidar, sonar, GNSS, odometry and other types sensors. Control software and systems, integrated in the vehicle, fusion and interpret the sensor information to identify the current position of the vehicle, detecting obstacles in the surround environment, and choosing the most appropriate reaction of the vehicle, ranging from stopping to bypassing the obstacle, reducing its speed, making a turn etc.

For the Macro-navigation, that is the destination to reach, the Autonomous Vehicle receives the information from either the in-vehicle operator (in the current configuration with a fixed path route), or from the remote control service via a dedicated 4/5G communication channel, for a fleet-managed operation. The fleet management system takes into account all available vehicles in the services area, the passenger request, the operator policies, the street conditions (closed streets) and send route and stop information to the vehicle (route to follow and destination to reach).

## 1.2.2    Autonomous vehicle capabilities in AVENUE

The autonomous vehicles employed in AVENUE fully and autonomously manage the above defined, micro-navigation and road behaviour, in an open street environment. The vehicles are autonomously capable to recognise obstacles (and identify some of them), identify moving and stationary objects, and autonomously decide to bypass them or wait behind them, based on the defined policies. For example with small changes in its route the AVENUE shuttle is able to bypass a parked car, while it will slow down and follow behind a slowly moving car. The AVENUE vehicles are able to handle different complex road situations, like entering and exiting round-about in the presence of other fast running cars, stop in zebra crossings, communicate with infrastructure via V2X interfaces (ex. red light control).

The shuttles used in the AVENUE project technically can achieve speeds of more than 60Km/h. However this speed cannot be used in the project demonstrators for several reasons, ranging from regulatory to safety. Under current regulations the maximum authorised speed is 25 or 30 Km/h (depending on the site). In the current demonstrators the speed does not exceed 23 Km/h, with an operational speed of 14 to 18 Km/h. Another, more important reason for limiting the vehicle speed is safety for passengers and pedestrians. Due to the fact that the current LIDAR has a range of 100m and the obstacle identification is done for objects no further than 40 meters, and considering that the vehicle must safely stop in case of an obstacle on the road (which will be "seen" at less than 40 meters distance) we cannot guarantee a safe braking if the speed is more than 25 Km/h. Note that technically the vehicle can make harsh break and stop with 40 meters in high speeds (40 -50 Km/h) but then the break would too harsh putting in risk the vehicle passengers. The project is working in finding an optimal point between passenger and pedestrian safety.

# 1.3 Preamble

WP4: Development, Adaptation and integration of Passenger Transport and in-, out-of-, vehicle services, aims to design, develop, adapt and integrate services to support users of autonomous vehicles before the trip, during the trip, and at the end of the trip. The main objective of WP4 is to provide services in order to demonstrate that the user experience can be seamless and secure, and that people embrace this new technology. Hence, we have to include the following services:

- Adapt and integrate existing transport services
- Develop autonomous vehicle specific services
- Provide services that foster the acceptance of driverless vehicles by both passengers and people interacting with the shuttles
- Introduce safety related services

The target of task T4.2 is to develop, teste and integrate innovative in-vehicles services, in collaboration with the four operators in Lyon, Luxembourg, Geneva and Copenhagen, respectively and with technical providers of the AVENUE consortium. The in-vehicle services should in combination with the out-of-vehicle services support a holistic service for travellers commuting with the AVs.

In-vehicle services are services developed to improve the user experience when traveling with autonomous vehicles. The services are user-centric and focus on supporting travellers (including PRM) with smart solutions while sitting inside the vehicle - in this case the Navya autonomous vehicle.

This deliverable, D4.5: Second iteration in-vehicle services, aims at introducing the preoperational work conducted with the purpose of prototyping and testing innovative and state-of-the-art technologies in autonomous shuttles.

This deliverable introduced the preparation of five selected in-vehicle services that will be prototyped and tested over the next 8-12 months. All five services are based on camera and sensor technologies, developed by CERTH. The services will be tested mainly on the two Amoblity routes (Copenhagen and Oslo). Once validated and tested the services will be further tested on the other AVENUE sites.

# 2 In-vehicle services

The end user in-vehicle services from the proposal can be seen in the following including type and timeframe:

1.  Intelligent ticket control (In-vehicle service)
2.  In-vehicle entertainment (In-vehicle service)
3.  Virtual personality interaction (In-vehicle service)
4.  Emergency automatic call system (In-vehicle service)
5.  Enhance the sense of security and trust (In-vehicle service)
6.  Prevention of night aggressions (In-vehicle service)

7.  Visualization in real time of the path / destination (In-and-out-of-vehicle service)
8.  Follow my kid / grandmother (In-and-Out-of-vehicle service)
9.  Mutual help facilitation (In-and-Out-of-vehicle service)

10. Passenger presence (In-vehicle service + AV functionality service)

## 2.1 Existing in-vehicle technical facilities

As of today, the Navya autonomous vehicle is equipped with some functions which will be integrated in the provision of services, targeted at the passengers. Therefore, the functions are highlighted with the purpose of including them in the further service development process.

The in-vehicle technical functions are as follows:

Open/close doors button
Manual buttons inside the vehicle, allowing the passengers to press them. The doors do also open automatically at each stop.

Wheelchair ramp button
Manual button inside the vehicle. Passenger or safety driver can use the button to activate the automatic wheelchair ramp. This function is not installed in all Navya autonomous vehicles as a standard. It must be installed on the shuttle by Navya.

Real time visualisation of route and stops (service screen)
There is a screen inside the vehicle showing a map of the route, including the stops. The screen is also used by the safety driver to restart the vehicle or detect issues during operation. The screen is currently not designed for passengers, but for the safety driver. In the future, once the safety driver is removed from the vehicles, the screen will function as the main interaction point between the travellers and the service.

Speaker system

Speakers located inside the vehicle, enabling contact from Navya supervision to passengers during emergencies. Potentially in the future, the PTO's could use the system for service announcements, safety protocols or regular communication - especially once the safety driver is removed from the autonomous vehicles.

Emergency call button

A manual button located under the service screen allows the passengers to get in contact with Navya supervision. The button could potentially in the future be used to link the travellers with the operators' own supervision - leaving the contact to Navya to the operators instead.

Emergency stop button

Manual button located on each side of the large window in the middle. Manually brakes the vehicle.

Emergency door opening handle

Manual handle located on each side of the doors. Manually forces the doors to open.

Emergency Glass breaking hammer

A hammer clearly located inside the vehicle is available to be used in case passengers are blocked and need to exit from the window of the e-minibus

Emergency First Aid Kit

A First Aid kit is located inside the e-minibus and is easily accessible in case of need.

Fire extinguisher

Inside the e-minibus a fire extinguisher is installed for use in case of any emergency.

Interior camera

Inside the e-minibus there is a camera that records the inside of the e-minibus. The camera is a fisheye camera and can be used by the operators to assess the situation inside the shuttle.

For the prototyping and testing of the chosen five services, the speaker system, the emergency call button and the interior camera will be used for different tests.

## 2.2 AVENUE goal

As a part of the AVENUE goals and vision, the safety driver has to be taken out of the vehicles when possible. In order to do so it is essential to understand the role of the safety driver in shuttles. What kind of services are he/she performing, how is the existence of a safety driver perceived by the travellers and so forth. As a part of understanding this, the role of the safety driver is shortly described (mostly from D4.4 including some new insights) - as T4.2 focusses on in-vehicle services, the following description also focuses on the role of the safety driver inside the shuttle.

## 2.2.1    Safety driver

This presence of the safety driver enables the operators to gather valuable insight and observation about the users' behaviour and interaction with in-vehicle services, and thereby gather a solid foundation for further development and refinement of existing and new in-vehicle services. The safety driver in-vehicle services are identified with the purpose of understanding what is necessary to automate in order to take the safety driver out of the vehicle. Furthermore, by identifying the role of the safety driver it is possible to pair up the proposal services.

The analysis of the in-vehicle services provided by the safety driver is mainly based on the initial insights from the State of the Art analysis and the experience of Autonomous Mobility. Once all operators are in operation, the safety driver inside the vehicles will be further analysed and compared across the operation sites - contributing to a more refined and exhaustive definition of the in-vehicle services that needs to be designed in the future.

The following table 1 shows the current in-vehicle services provided by the safety driver. Each service is shortly described, with the purpose of initiating a thorough investigation of what is necessary once the safety driver is removed from the vehicle. For each part one or more of the in-vehicle proposal services are paired, where they make sense.

| **Safety driver** Current in-vehicle services provided by the safety driver | | | |
|---|---|---|---|
| # | In-vehicle service | Description | Proposal services that potentially can automate the service |
| 1 | Operational information | Providing information about the project, route, vehicle, laws etc. | - Visualization in real time of the path / destination (Mobile application) |
| 2 | Safety + risk mitigation | Safety driver presence and authority establishes a safe environment for the passengers. + Active risk control (joystick, eyes etc.) | - Security and trust services<br>- Prevention of night aggressions<br>- Emergency automatic call system<br>- Mutual help facilitation |
| 3 | Entertainment | The safety driver can interact with the passengers and start conversations about anything + tell stories about the technology etc. | - In-vehicle entertainment<br>- Virtual personality interaction |
| 4 | Travel assistance | The safety driver can assist passengers in boarding and exiting the vehicle + ensure that the bus stops where the passengers want to get on/off. | - Intelligent ticket control<br>- Follow my kid / grandmother<br>- Passenger presence |
| 5 | Area guidance | The safety driver can assist passengers in finding restaurants, shops etc. in the local area of the route. | - Visualization in real time of the path / destination (Mobile application)<br>- Virtual personality interaction |
| 6 | Branding | Branding of operating company, logo, name, etc. Can be used to assure users that the operation is safe and offered by a trusted company. Things are under control. | - Visualization in real time of the path / destination (Mobile application) |

| 7 | Provides data on various areas | User behaviour/insights + counting passenger + road user behaviour | - Visualization in real time of the path / destination (Mobile application)<br>- Virtual personality interaction |
|---|---|---|---|

*Table 1 - Safety driver (in-vehicle services)*

# 2.3 In-vehicle service focus

The services that will be focused on in the next 8-12 months are as follows:

- Security trust services (example: Prevention of night aggressions)
- Automated passenger presence
- Follow your kid/grandparents
- Shuttle environment assessment
- Smart feedback system

The reason for choosing these services are two-fold:

- They are all essential services necessary to remove the safety driver from the shuttle, as safety and efficient operation are critical factors. For example ensuring safety for the travellers have been identified in WP2 as the single most important factor for choosing to ride with the autonomous shuttle.
- They are all based on the same foundation of technology with cameras, sensors and algorithms, meaning that we can harness from testing multiple types of services, once we have installed the necessary equipment. The analysis is also similar and therefore the results can be compiled together. The technology, equipment and algorithms are provided by CERTH and tested and prototyped together with Amobility.

During the next 8-12 months the above described services will be tested, further prototyped and analysed on the Oslo and Copenhagen site of Amobility. The purpose of the development of services is to prototype and test the services as a concept. The services will probably not be fully developed and implemented, but recommendations for further improvements and development will be the results of the service development and testing. For the development and testing a close collaboration will be established with MobileThinking as they have developed the mobile application that can be connected to the in-vehicle services over time.

# 2.4 Service: Enhance the sense of security and trust

## 2.4.1  Concept of service

The service "Enhance the sense of security and trust" aims to address the new reality that is formed in autonomous shuttles mobility infrastructures as a result of the absence of the bus driver and the increased threat from terrorism in European cities. Typically, drivers are trained to handle incidents of passengers' abnormal behaviour, incidents of petty crimes, etc. according to standard procedures adopted by the transport operator. Surveillance using sensors such as cameras (cameras of different technologies can be used so that passengers' privacy is protected) and microphones, as well as smart software in the bus will maximize the feeling of security and the actual level of security.

Several concerns of the end users regarding the Safety and Robustness of the autonomous vehicles that are directly linked to the final User Acceptance of the new technology, can be identified. The prospective passengers fear several possible instances that could arise in case there is no driver in the bus. Indicatively:

- No one will be in the bus to perform first aid if required
- Feeling of discomfort being all alone in the bus at night, especially in certain neighbourhoods
- No authority figure present to keep passengers calm (eg. school kids)
- Vandalism, bag snatching, indoor fighting, unaccompanied luggage

To address the aforementioned concerns on social and personal safety and security into the vehicle, certain measures need to be implemented. For example the detection of unaccompanied luggage and of other personal belongings may raise a notification or an alert to the supervisor and/or the suitable authorities. This may be followed by appropriate notifications and/or instructions to the passengers, while the vehicle may also implement respective actions.

Moreover, implementing a solution for enhancing the safety and security inside the autonomous buses will support safekeeping not only the users of the autonomous public bus but also the vehicle itself. In this section, the implementation of a video, depth and audio analytics software module for an embedded security subsystem or for cloud-based services of the system are described along with appropriate planning for the deployment and test of the service into the pilot sites of the Avenue project.

### 2.4.1.1  Use case

The service addresses the timely, accurate, robust and automatic detection of various petty crime types or misdemeanors as well as the assistance of authorized end-users towards the re-identification of any offenders. A misdemeanor is any "lesser" criminal act in some common law legal systems. Misdemeanors are generally punished less severely than felonies, but theoretically more so than administrative infractions (also known as minor, petty, or summary offences) and regulatory offences. Many misdemeanors are punished with monetary fines.

The petty crimes that are targeted for identification by the sensors include: petty theft like bag snatching and pickpocketing, vandalism, aggression, illegal consumption of cigarettes, public intoxication, simple assault and disorderly conduct. These are explained in more detail:

Petty theft: In common usage, theft is the taking of another person's property or services without that person's permission or consent with the intent to deprive the rightful owner of it.

Vandalism: Vandalism is the action involving deliberate destruction of or damage to public or private property.

Aggression: Aggression is overt or covert, often harmful, social interaction with the intention of inflicting damage or other unpleasantness upon another individual.

Public intoxication: Public intoxication, also known as "drunk and disorderly" and drunk in public, is a summary offense in some countries rated to public cases or displays of drunkenness.

Simple assault: An assault is the act of inflicting physical harm or unwanted physical contact upon a person or, in some specific legal definitions, a threat or attempt to commit such an action.

Disorderly conduct: Disorderly conduct makes it a crime to be drunk in public, to "disturb the peace", or to loiter in certain areas.

In the context of Avenue project, the following use cases have been identified to be further examined and addressed:

**Use Case 1: Unaccompanied Luggage Monitoring**
- Inside the autonomous shuttle there is unaccompanied luggage which remains unmoved for a long time.
- The video cameras installed in the autonomous shuttle acquire the color depth images and the data are fed into the system's video analytics algorithms for further analysis.
- In case the algorithms identify that the total time of the luggage that remains unmoved in the vehicle passes the predefined time frame, a notification is sent to the security operator.
- The security operator monitors the clips that are captured and evaluates the criticality of the situation and whether to intervene or not.

**Use Case 2: Bag Snatch Detection**
- While a commuter is in the autonomous shuttle a petty crime takes place in the form of assault.
- The commuter is attacked by another person who is attempting to snatch the bag from the commuter.
- The aggressive incident is captured by the sensors in the bus, both microphone and cameras, and fed to the video analytics component for further analysis.
- The system sends a security alert to the operator or security supervisor.
- The course of action of the operator is a human decision, that means, whether he/she will decide to stop the autonomous shuttle or will notify the passengers via the radio system.

**Use Case 3: Vandalism Detection**
- A young person onboards the autonomous shuttle during an itinerary performed by the vehicle during the night hours.
- The youngster attempts to perform a vandalism action on the shuttle, through painting graffiti on the windows.

- The night mode of the cameras installed in the vehicle acquire the data that will be fed to the video analytics algorithms for further analysis.
- The person is warned by the radio system of the autonomous shuttle or security personnel intervenes by stopping the bus.

**Use Case 4: Sound Events Detection**
- A person onboards the autonomous shuttle during an itinerary performed by the vehicle.
- The person attempts to terror the passengers on the shuttle, through breaking the windows. The passengers may also scream during this behavior or call for help.
- The microphones installed in the vehicle acquire the data that will be fed to the audio analytics algorithms for further analysis.
- The person is warned by the radio system of the autonomous shuttle or security personnel intervenes by stopping the bus.

## 2.4.2     Stakeholders (development/prototyping team)

In this section, the relevant stakeholders involved in the development and prototyping of the Enhance the Sense of Security and Trust Service are presented. Specifically:

- **CERTH** develops software and algorithms for the detection of abnormal events using video and sound analysis techniques. The service is implemented by adopting a platform for petty crimes and incident detection systems developed for the active detection of abnormal behaviour, as well as suspicious objects in conjunction with a variety of sensors. Sensors such as cameras and microphones, along with machine learning algorithms are employed for the timely, accurate and robust detection of petty crimes and incidents. CERTH is also involved in the installation of the related sensors and software into the autonomous shuttles on the operators facilities.
- **Bestmile** provides multiple integration interfaces towards the different stakeholders in the project (e.g. vehicle manufacturers, public transport operators) into its cloud platform. In this service, Bestmile provides connectivity between the vehicle and the related operators regarding the notifications that are generated by the detection software.
- **MobileThinking** develops the AVENUE mobile app for the end users that will be involved into the project. In this service, MobileThinking develops the interconnection with the AVENUE mobile app by providing the notifications generated from the system to the end users.
- **Amobility** is the operator of the autonomous vehicles in the Copenhagen site and is handling all daily operation of the vehicles and everyday contact with the end users. Amobility is leading the T4.2 and will be coordinating and developing the testing framework for the service validation and prototyping of the services in the vehicles. As an operator Amobility provides the autonomous shuttles that will be used for the AVENUE pilot activities and its facilities for performing data capture activities, deploying the service and testing its performance through short or longer evaluation periods. The shuttles that will be used to test and prototype the service are located in Copenhagen,Nordhavn and Oslo, Omrøya. Once the prototyping and testing framework has been developed and tested, the prototyping session will be coordinated on the other AVENUE sites, if possible. Hence via TPG, SLA and KEOLIS.

## 2.4.3     Technical requirements

The "Enhance the sense of security and trust" service exhibits a variety of advantages, which include the onsite intelligence of the platform using low cost devices that can decide with minimal latency whether the cloud side of the system must be notified, saving thus network/cloud resources. Moreover, the onboard side analytics can confirm or reject a video clip as true or false positive respectively, provides the relevant information to the end user, allowing him to focus his/her attention on crime incidents. In addition, the identification components will ease the task of authorities contributing to the timely apprehension of any offender.

In the next sections, the indicative technical specifications of the system are presented to host the algorithms that are developed based on RGB, audio and depth sensors for the embedded and cloud side of the surveillance platform. In this context, various analysis approaches for early petty crime detection are adopted, such as

1. Video analytics
2. Depth analytics
3. Audio analytics

The embedded system is

- Easy to install and use (normally less than an hour).  Adaptable to the network conditions.
- Supports a multitude of video/audio formats and codecs (98 codecs; including all the most widely used ones).
- Supports standalone operation when connection with the cloud is not an option.
- Allows real time re-configuration of the system's components.

Certain sensors need to be installed in the autonomous vehicles which include:

- Embedded camera support
- USB cameras support
- IP cameras support
- USB, IP microphones support

The sensors of the embedded system can include various types, have enhanced interoperability and support integration with legacy cameras. When an event is detected a clip is generated. The streams will be cut to 1 minute packages and will be analysed.

The system utilizes standard of-the-shelf **cameras** and **microphones**, so no specialized protocols and/or ports are required. In most cases the sensors will be directly connected with the Embedded System through a USB port, whereas IP cameras could also be supported, with the latter being connected through Ethernet or Wi-Fi with the Embedded System. The following tables present the detailed characteristics of indicative sensors that can be used. At this point it has to be noted that since the system can support a number of different cameras and microphones, the minimum requirements are presented and not a specific model.

**Sensor/Device: Video Sensing**

| Name | The system supports most of-the-shelf cameras |
|---|---|
| Short Description | The sensor can be utilized to detect abnormal events in the indoor environment of the autonomous shuttles. |
| Measurement | The sensor must be capable to acquire colour images and should have a night mode. The Sensor will be able to acquire data that will be fed to the Video analytics algorithms for further analysis. The sensor is attached to a PC using the USB-interface. IP cameras are also supported. |
| Functionality | The sensor will be part of the system in order to detect abnormal events (e.g. theft, fighting, anti-social behaviour) on the autonomous shuttles. Moreover, in case an abnormal event is detected, data will be sent to the AVENUE platform in order to identify the suspect(s). |

Indicative specifications are presented next. Detailed specifications will be derived based on discussions with Navya, the manufacturer of the autonomous vehicles. The dimensions and the weight of the sensor vary and the mounting can be made manually. In terms of resolution, the sensor should provide minimum HD resolution images (1280x720 pixels) at 24 frames per second and the operating temperature range applies (e.g. 5 to 35 degrees Celsius). The sensor will be powered via the USB 2.0 or Power over Ethernet (PoE) connection. An IP Camera should support widely used streaming protocols (e.g RTSP).

The RGB video stream should be 24-bit HD (ideally FullHD 1080p) resolution. The data rate is at minimum 30 Hz but it highly depends on the actual application use for real-time performance. The data stream is continuous and can be acquired also as required (or on demand). The API and the SDK of the camera is flexible enough to acquire the images as soon as they are required on the aforementioned maximum data epoch rate.

The transmission frequency should be 30 Hz. No specific software is required. The sensor acquires data that could be used to identify a person. However, within the system special countermeasures will be taken into account in order to address all ethical, legal and privacy issue.

## Sensor/Device: Depth Imaging

| Name | Kinect Sensor v2 (Microsoft) |
|---|---|
| Short Description | The sensor can be utilized on the interior of public transportation (i.e. autonomous shuttles) |
| Measurement | The sensor is capable to acquire colour, infrared and depth images. Within the system, only the depth images will be used in each case. The depth sensor will be able to acquire depth data that will be fed to the system's video analytics algorithms for further analysis (Event detection, person identification). The sensor is attached to a PC using the USB-interface (USB 3.0) |

**Functionality**   The sensor will be part of the system in order to detect abnormal events (e.g. theft, fighting) on indoor environments of autonomous shuttles. Moreover, in case an abnormal event is detected, data will be sent to the AVENUE platform in order to notify the related stakeholders.

Indicative specifications are presented next. Detailed specifications will be derived based on discussions with Navya, the manufacturer of the autonomous vehicles. The Kinect Sensor head size is 292.1 x 304.8 x 184.15 (in mm). It weighs less than 1kg and the mounting should be made manually. The operating sensor range for depth images is from 0.4m to 8m (extended range with lower accuracy). The measurement resolution is given for both color and depth sensing lenses: Colour VGA motion camera 1920 x 1080 pixel resolution @30 FPS Depth Camera 512 x 424 pixel resolution @30 FPS Infrared Camera 512 x 424 pixel resolution @30 FPS The corresponding field of view is the following: Horizontal field of view- 70 degrees Vertical field of view- 60 degrees.

At the maximum range the random error of depth measurements increases quadratically. It can reach even 4cm at the maximum range (8 meters). The Kinect Sensor is pre-calibrated with its lens. Operating temperature range applies (e.g. 5 to 35 degrees Celsius). The power is supplied from the mains by way of an AC adapter 12 Watts power demand through a standard USB 3.0 port interface. The output of the sensor is raw data for color and depth images (RGB and depth images with pixel values representing depth in mms). The RGB video stream uses 32-bit FullHD resolution (1920 × 1080 pixels) with a Bayer color filter, while the monochrome depth sensing video stream is in resolution of 512 × 242 pixels) with 16-bit depth. The data rate is about 30 Hz but it highly depends on the actual application use for real-time performance. The data stream is continuous and can be acquired also as required (or on demand). The API and the SDK of the camera is flexible enough to acquire the images as soon as they are required on the aforementioned maximum data epoch rate. The transmission frequency is 30 Hz.

The sensor acquires data that could be used to identify a person. However, within the system only the depth images will be used for analysis and special countermeasures will be taken into account in order to address all ethical, legal and privacy issues. The Kinect Sensor will be part of the system, as a more advanced proposal for volume crime detection in indoor environments.

## Sensor/Device: Audio Sensing

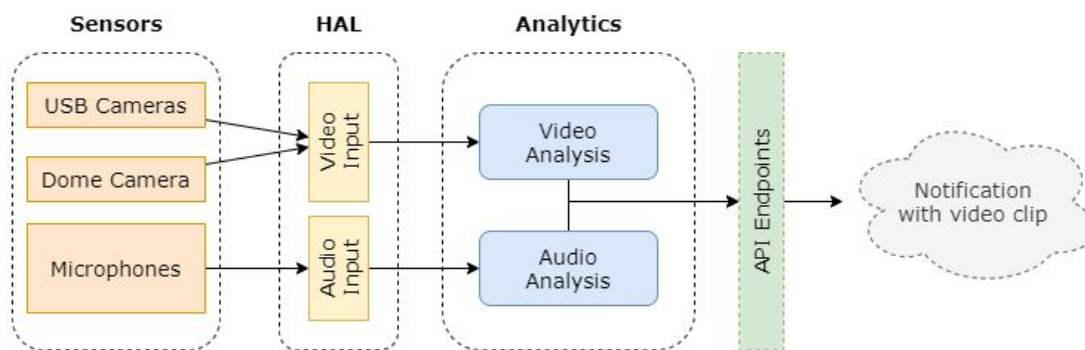| | |
|---|---|
| **Name** | The system supports most of-the-shelf microphones |
| **Short Description** | The sensor can be utilized on the interior of public transportation (i.e. autonomous shuttles) |
| **Measurement** | The sensor must be capable to acquire audio at 16-bit audio at a sampling rate of 16 kHz. The sensor is attached to a PC using the USB-interface. |
| **Functionality** | The sensor will be part of the system in order to detect abnormal events (e.g. screaming, breaking glass) on the indoor environments of autonomous shuttles. Moreover, in case an abnormal event is detected, data will be sent to the AVENUE platform in order to notify the related stakeholders. |

Indicative specifications are presented next. Detailed specifications will be derived based on discussions with Navya, the manufacturer of the autonomous vehicles. The dimensions and weight of the sensor vary and the mounting should be made manually. The operating temperature range applies (e.g. 5 to 35 degrees Celsius). The data connection is USB 2.0 through which the powering is performed. The data format is WAVE format. The sensor acquires data that could reveal a person's identity. However, within the system special countermeasures will be taken into account in order to address all ethical, legal and privacy issues.

## Other requirements

A PC to host the system and all developed algorithms is required. Power supply and data connectivity are also required in this case.

### 2.4.3.1 Technology

In the following sections, the research, involving algorithms and experiments conducted, is presented for the "Enhance the sense of security and trust" service. As depicted in Figure 1, the first layer of sensors connects to the Hardware Abstraction Layer (HAL). The HAL implements the IP and the USB protocol supporting IP and USB camera, respectively. The input data is converted and transformed in a compatible format and passed into the analytics algorithms. The prediction and a short video clip (if required) are then transferred via the API endpoints into the cloud. The operator has access to the data and acts accordingly.



**Figure 1. High level overview of the "Enhance the sense of security and trust" service**

In this context, various analysis approaches for early petty crime detection are implemented, such as:
1. Video analytics (described at Section 2.4.3.1.1)
2. Audio analytics (described at Section 2.4.3.1.2)
Each technique is presented and explained, along with its results and its major shortcomings.

### 2.4.3.1.1 Video analysis

This section focuses on the services that are going to be deployed in the autonomous shuttle, using images from video feeds. Specifically, we describe the datasets that have been used and the algorithms for the image-based abnormal event detection starting from the pre-processing steps to the classification deep neural network modules.

For video analysis, three different approaches were implemented: (a) a Pose Long Short-Term Memory (LSTM) Classifier, (b) a Spatiotemporal Autoencoder and (c) a Spatiotemporal LSTM Classifier.

## Background

Automatic awareness of human actions and its interaction with the environment has become a prominent study area in recent years. To perform such a demanding mission, many scientific areas rely on modeling human activity in its various dimensions (emotions, relational attitudes, actions, etc.). In this context, the identification of a person's activity tends to be necessary to the comprehension of specific acts. Thus, a great interest has been granted to human action recognition, especially in real-world environments. There are various attempts on human action recognition based on RGB video and 2D/3D skeleton data. The RGB video-based action recognition methods[1][2][3][4] mainly focus on modeling spatial and temporal representations from RGB frames and temporal optical flow.

Despite RGB video-based methods have achieved promising results, there still exist some limitations, e.g., background clutter, illumination changes, appearance variation, and so on. 3D skeleton data represents the body structure with a set of 3D coordinate positions of key joints. Since skeleton sequence does not contain color information, it is not affected by the limitations of RGB video. Such robust representation allows to model more discriminative temporal characteristics about human actions.

Moreover, Johansson et al.[5] have given an empirical and theoretical basis that key joints can provide highly effective information about human motion. Besides, the Microsoft Kinect[6] and advanced human pose estimation algorithms[7] make it easier to gain skeleton data. For skeleton based action recognition, the existing methods explore different models to learn spatial and temporal features. Song et al.[8] employ a spatial-temporal attention model based on LSTM to select discriminative spatial and temporal features. The Convolutional Neural Networks (CNNs) are used to learn spatial-temporal features from

---

[1] Karen Simonyan and Andrew Zisserman. Two-stream convolutional networks for action recognition in videos. In NIPS, 2014

[2] Limin Wang, Yuanjun Xiong, Zhe Wang, Yu Qiao, Dahua Lin, Xiaoou Tang, and Luc Van Gool. Temporal segment networks: Towards good practices for deep action recognition. In ECCV, 2016

[3] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. In ICCV, 2015

[4] Pichao Wang, Wanqing Li, Philip Ogunbona, Jun Wan, and Sergio Escalera. Rgb-d-based human motion recognition with deep learning: A survey. Computer Vision and Image Understanding, 2018

[5] Gunnar Johansson. Visual perception of biological motion and a model for its analysis. Perception & Psychophysics, 1973

[6] Zhengyou Zhang. Microsoft kinect sensor and its effect. IEEE Multimedia, 2012.

[7] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In CVPR, 2017

[8] Sijie Song, Cuiling Lan, Junliang Xing, Wenjun Zeng, and Jiaying Liu. An end-to-end spatio-temporal attention model for human action recognition from skeleton data. In AAAI, 2017

**Figure 1. Screenshots of the utility: (a) Initialization screen (b) Pose labelling using mouse clicks on the skeleton's circle and keyboard navigation**

## LSTM Classification via Pose Estimation

The pipeline of this approach consists of 4 stages as shown in Figure 2. In stage one, the pose of each person in the frame is extracted (15 keypoints).  In the second stage, a skeleton tracking algorithm is performed, associating persons across multiple frames. Regarding the third stage, the detected and tracked human body key-points are represented as trajectories and during the fourth stage they are "fed" into a Multi-Layer Classifier which classifies each action into normal or abnormal.

**Figure 2. Pipeline of the Pose Estimation Classification**

**Pose Estimation:** For the first stage, we are using a custom implementation of the Regional multi-person pose estimation by Fang etal.[15] for training (better accuracy) and OpenPose[16] for testing (higher performance). The integrated implementation uses VGG19, a convolutional neural network model proposed by K. Simonyan and A. Zisserman[17]. We are using this model to improve the accuracy of the skeleton extraction for the training process and further perform data augmentation without sacrificing the data integrity. We generate noisy data with variable intensities, based on the extracted data from the backend, and we combine these data with the original ones as an augmentation technique. Extensive tests indicated that the model generalizes better and the accuracy improves. Although we initially implemented AlphaPose using VGG for training our model, multiple tests shown that we could use it for evaluation too, when specific parameters (resolution, heatmaps etc.) are tuned. We apply the pose estimation framework in the presence of inaccurate human bounding boxes. The generated pose proposals are refined by parametric pose non-maximum suppression to obtain the estimated human poses. In this stage, 17 different human body keypoints are detected and the number of people in each frame is obtained. The number of N frames for the feature generation along with the evaluation accuracy are depicted in the graph below (Figure 3). As we can see, a buffer size (window size) of 5 frames achieved the best accuracy on the evaluation test. Higher values may result in lower accuracy as the tracker may fail to consistently detect people when the shuttle is overcrowded. We are currently investigating some optimizations of the tracker, in order to further increase the size of the buffer.

---

[15] Fang, Hao-Shu, et al. "Rmpe: Regional multi-person pose estimation." Proceedings of the IEEE International Conference on Computer Vision. 2017.

[16] Cao, Zhe, et al. "OpenPose: realtime multi-person 2D pose estimation using Part Affinity Fields." arXiv preprint arXiv:1812.08008 (2018).

[17] Simonyan, Karen, and Andrew Zisserman. "Very deep convolutional networks for large-scale image recognition." arXiv preprint arXiv:1409.1556 (2014).

**Figure 3. Performance evaluation across different buffer sizes**

**Tracking:** An important step is to match cross-frame poses and form pose flows (tracking). Also, a novel pose flow non-maximum suppression is applied to reduce redundant pose flows and re-link temporal disjoint ones. This is an important step that associates poses indicating the same person, across multiple frames.

We implemented an online skeleton tracking algorithm based on distance and other heuristics, in order to meet the performance requirements of the real-time service. The algorithm is sorting the skeletons based on the distance between neck and image center, from small to large.

A skeleton near center will be processed first and be given a smaller human id. Later on, each skeleton's features will be matched between the current and the previous frame.



**Figure 4. Skeleton matching across two different frames (blended)**

The distance matrix (or cost) between the skeleton joints is main criterion for the matching function. Skeletons with the less distance are matched between the frames and are given the same id (Figure 4). In some cases, the skeleton detection framework might fail to detect a complete human skeleton from the image due to the restricted field of view of the camera in the autonomous vehicle.

This may cause some blanks in the joint positions, which should be filled with some values in order to maintain a fixed-size feature vector for the following feature classification procedure.

We evaluated some solutions for this issue:
- **Solution 1:** Discard this frame. However, the algorithm would never be able to detect the action when the person is standing sideways and not facing the camera.
- **Solution 2:** Fill in the positions with some value outside a reasonable range. Theoretically, when the classifier is strong enough, this method could work.
- **Solution 3:** Fill in a joint's position based on its relative position in the previous frame with respect to the neck.

In order to solve this issue, solution 3 was implemented (Figure 2), but the classifier's performance was degraded in some test cases. After extensive tests, we noticed that a previous joint position might be missing too, being replaced by the estimation of our algorithm. This led to "stuck" joints across various frames and confused both the tracker and the classifier.

To overcome this issue, we are using a default "idle" pose as an example for our algorithm. When a previous joint is missing, the value being replaced is relative to the default example pose (Figure 5). We chose a person sitting as the default, because it is the most common for the passengers in the AV.



**Figure 5. The two leg joints are being reconstructed (right) by their relative location in the previous frame (left)**

**Feature Extraction:** Regarding to the third stage, the detected and tracked human body key-points are converted into features and forwarded to an LSTM Neural Network. For extracting features, we store every person's skeleton data into a circular (ring) buffer deque (double-ended queue) of N frames (window size) into the Feature Generator class.

21

**Figure 6. Representation of the extracted features**

We consider the buffer as invalid if the newest appended skeleton does not contain at least the neck (Point 0) or one of the thigh bones (Point 7 or 10) shown in Figure 6, as the height of the skeleton (used for normalizing features) cannot be calculated. The feature extraction process occurs when the buffer is full.

| Pose Classification – Feature Extraction | |
|---|---|
| **Xs** | A direct concatenation of joints positions of the N frames. |
| **H** | Average height of the skeleton of the previous N frames. This height equals the length from Neck to Thigh. Used for normalizing features. |
| **X** | Normalized joint positions $[Xs - mean(Xs)]/H$ |
| **Vj** | Velocities of the joints $\{X[t]-X[t-1]\}$ |
| **Vc** | Velocity of the center $\{sum(Xc[t]-Xc[t-1])\}$ (**10x** weight) |

**Classification:** The LSTM model is capable of binary or multi-class softmax classification. It contains three hidden layers of size (32x64) with the rectified linear unit (ReLU) activation function. An overview is shown in Figure 7.



**Figure 7. Model Overview**

The model is trained end-to-end and regularized so that it distills the most compact profile of the normal patterns of training data and effectively detects abnormal events (Figure 8).



**Figure 8. Abnormal event detection (passengers are fighting)**

We applied a Principal Component Analysis (PCA) procedure to reduce the 314 initial features to 50 principal components. PCA is a statistical procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables (entities each of which takes on various numerical values) into a set of values of linearly uncorrelated variables called principal components. This transformation is defined in such a way that the first principal component has the largest possible variance (that is, accounts for as much of the variability in the data as possible), and each succeeding component in turn has the highest variance possible under the constraint that it is orthogonal to the preceding components. The resulting vectors (each being a linear combination of the variables and containing n observations) are an uncorrelated orthogonal basis set. PCA is sensitive to the relative scaling of the original variables. After the PCA procedure, we achieved a sum of 0.9981\% eigenvalues. The model converges faster and the accuracy improves. Due to the highly imbalanced classes, we compiled our model with custom **weighted** categorical cross-entropy function for loss calculation, and a weighted categorical accuracy method in order to acquire accurate metrics.

The evaluation performance of our model is depicted below (Figure 9). After 50 epochs, our model achieved a 96.22% accuracy. The time cost for feature extraction and classification is less than 0.05s per frame for the classifier, since the model is relative shallow.

**Figure 9. Model performance evaluation – (a) Training accuracy and (b) Training loss metrics**

## Spatiotemporal Autoencoder

The second approach is based on the principle that the most recent frames of video will be significantly different than the older frames, in case of an abnormal event. Our goal, inspired by Lu etal.[18], is to train an end-to-end model consisting of both a spatial feature extractor and a temporal encoder-decoder that combined learn the temporal patterns of the input volume of frames. So as to minimize the reconstruction error between the input and the output video volume reconstructed by the learned model, we trained our model using video volumes of only normal scenes. After our model's proper training, we expect to have low reconstruction error in a normal video volume as opposed to a video volume containing abnormal scenes. Finally, our system will be able to detect the occurrence of an abnormal event, by thresholding on the error produced by each testing input volume.

**Preprocessing:** At this stage, our task is to convert raw data into aligned and acceptable input for the model. To do so, each frame extracted from the raw videos is then resized to 64x64. To ensure that the input images are all on the same scale, the pixel values are scaled between 0 and 1 and each frame subtracted by its global mean image for normalization. The mean image is calculated by averaging the pixel values at each location of each frame in the training dataset. After that, the images are converted to grayscale to reduce dimensionality. Finally, the processed images are then normalized to have zero mean and unit variance. As mentioned before, we use video volumes as input to our model, where each volume consists of 10 consecutive frames with various skipping strides (Figure 10). As the number of parameters in this model is large, we also need a large amount of training data. Following Lu's practice, to increase the size of the training dataset, we perform data augmentation in the temporal dimension. To generate these volumes, we concatenate frames using sequences, namely being stride-1, stride-2, and stride-3. For example, the stride-1 sequence consists of frame numbers {1, 2, 3, 4, 5, 6, 7, 8, 9, 10}, whereas the first stride-2 sequence contains frame numbers {1, 3, 5, 7, 9, 11, 13, 15, 17, 19}. Now the input is ready for model training.



**Figure 10. Preprocessing and training pipeline**

**Feature Learning:** In order to learn the regular patterns in training videos, we propose a convolutional spatiotemporal autoencoder. Our architecture consists of two parts — a spatial autoencoder for learning

---

[18] Lu, Yiwei, et al. "Future Frame Prediction Using Convolutional VRNN for Anomaly Detection." 2019 16th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS). IEEE, 2019.

spatial structures of each video frame, and a temporal encoder-decoder for learning temporal patterns of the encoded spatial structures. As illustrated in Figure 11, the spatial encoder and decoder have two convolutional and deconvolutional layers respectively, while the temporal encoder is a three-layer convolutional long short-term memory (LSTM) model. Convolutional layers are well-known for their superb performance in object recognition, while the LSTM model is widely used for sequence learning and time-series modeling and has proved its performance in applications such as speech translation and handwriting recognition.

**Autoencoder:** There are two stages that form an autoencoder: encoding and decoding. Autoencoders set the number of encoder input units less than the input; thus, they were first used to reduce dimensionality. Usually, unsupervised back-propagation is used for training, minimizing the reconstruction error of the decoding results from the original inputs. Generally, an autoencoder can extract more useful features when the activation function is non-linear rather than some common linear transformation methods, such as PCA.



**Figure 11. Autoencoder Model Architecture**

**Spatial Convolution:** The primary purpose of convolution in a convolutional network is to extract features from the input image. Convolution can preserve the spatial relationships between pixels by using small squares of input data to learn image features. Mathematically, convolution performs dot products between the filters and local regions of the input. Assuming that we have some n x n square input layer, followed by the convolutional layer, then if we use a m x m filter W, the convolutional layer output will be of size (n-m+1) x (n-m+1).

During the training process, a convolutional network learns the values of these filters on its own, although parameters such as the number of filters, filter size, the number of layers before training still need to be specified. The larger the number of filters used, the more image features get extracted and the better the network becomes at recognizing patterns in unseen images. However, balance is key when it comes to the number of filters used, as more filters would add to computational time and exhaust memory faster.

**Recurrent Neural Network (RNN):** In a traditional feedforward neural network, we assume that all inputs (and outputs) are independent of each other. However, in tasks involving sequences, learning temporal dependencies between inputs are important, as e.g. a model of word predictor should be able to derive information from the past inputs.

An RNN works just like a feedforward network, except that the values of its output vector are influenced not only by the input vector but also on the entire history of inputs. In theory, RNNs can make use of information in arbitrarily long sequences, but in practice, due to vanishing gradients, they are limited to looking back only a few steps.

**Long Short-Term Memory (LSTM):** To overcome this problem, a variant of RNN is introduced:  a Long Short-term Memory (LSTM) model that incorporates a recurrent gate called forget gate. With the new structure, LSTMs prevent backpropagated errors from vanishing or exploding.  Therefore, LSTMs can work on long sequences and can be stacked together to capture higher level information.



**Figure 12. The structure of a typical LSTM unit. The blue line represents an optional peephole structure, which allows the internal state to look back (peep) at the previous cell state Ct−1 for a better decision**

**Convolutional Long Short-term Memory (ConvLSTM):** The Convolutional Long Short-term Memory (ConvLSTM) model, considered a variant of the LSTM architecture, was introduced by Shi et al. in[19] and has been recently utilized by Patraucean et al.[20] for video frame prediction. Compared to the usual fully connected LSTM (FC-LSTM), ConvLSTM has its matrix operations replaced with convolutions. ConvLSTM requires fewer weights and yields better spatial feature maps, by using convolution for both input-to-hidden and hidden-to-hidden connections.

---

[19] Xingjian, S. H. I., et al. "Convolutional LSTM network: A machine learning approach for precipitation nowcasting." Advances in neural information processing systems. 2015.

[20] Patraucean, Viorica, Ankur Handa, and Roberto Cipolla. "Spatio-temporal video autoencoder with differentiable memory." arXiv preprint arXiv:1511.06309 (2015).

**Figure 13. The zoomed-in architecture at time t, where t is the input vector at this time step. The temporal encoder-decoder model has 3 convolutional LSTM (ConvLSTM) layers.**

**Regularity Score:** Once the model is trained, its performance can be evaluated by feeding in testing data and checking whether it can detect abnormal events while maintaining a low false alarm rate. For a better comparison, we used the same formula as Hasan etal.[21] to calculate the regularity score for all frames. Our only difference is that the learned model is of a different kind. The reconstruction error of all pixel values in a frame of the video sequence is taken as the Euclidean distance between the input frame and the reconstructed frame. The reconstruction error or cost of a frame sequence is calculated as the difference between the ground truth (original frames) and the reconstructed frames (prediction – model output). The frame sequence is flagged as "abnormal" if the reconstruction cost exceeds a certain threshold.

---

[21] Hasan, Mahmudul, et al. "Learning temporal regularity in video sequences." Proceedings of the IEEE conference on computer vision and pattern recognition. 2016.

**Figure 14. Architecture overview of the autoencoder approach**

## Hybrid LSTM Classification

Even if our dataset consists of thousands of data, we will get some anomalies in certain occasions. So now humanly it is possible to manually go through the anomaly outputs and flag some of them as false positives. Therefore, we can let the previous autoencoder neural network model act as a **High Recaller**.

**Semi-supervised Learning:** The threshold is decreased so that almost all the actual anomalies are getting detected (high recall) along with other false positive anomalies (low precision). To achieve the semi-supervised approach, we designed a new model which includes the previous Encoder and an LSTM which acts as a classifier.



**Figure 15. Example of a prediction with a lower regularity threshold**

In real-time inference the anomalies predicted by the high recaller model (autoencoder neural network) are sent through the false positive reduction model (hybrid model). This combination of neural networks together should provide a deep neural network model with high recall and high precision.



**Figure 16. Model architecture of the hybrid model. The red container contains components of the previous autoencoder approach. The green components indicate the new hybrid model which acts as a classifier**

**Training:** The training process of the new experiment consists of 3 stages: At first, the autoencoder model (encoder + decoder) is being trained with unsupervised data to learn regularity. In the second stage the encoder weights are transferred to the hybrid model. The encoder's layers are marked as non-trainable. Finally, we perform supervised training only the LSTM classifier.



**Figure 17. Pipeline of the hybrid training procedure**

30

## 2.4.3.1.2    Sound analysis

This section focuses on the services that are going to be deployed in the autonomous bus, using information from acoustic sensors. Specifically, we describe the datasets that have been used and the algorithms for the audio-based abnormal event detection starting from the pre-processing steps to the classification deep neural network module.

In sound analysis, it is important to analyse the signal either at its raw form (waveform in time-domain) or convert the input signal to the frequency domain and extract audio features (e.g., spectrogram representation, mel-frequency cepstral coefficients, spectral roll-off, etc.). An advantage of the latter, is that the dimensionality and complexity of the data set is significantly reduced, but there is some information lost in the process. In this Section, we will first introduce the MIVIA Audio Events Dataset, and then look at the sounds from a time perspective, frequency perspective and see how to extract information. Finally, we look at some of the samplings and see if we can classify them from a human visual perspective.

**Background**

In the everyday life, we are continuously surrounded by sounds, and usually the human brain is able to recognize what kind of sound it is based on experience. Artificial neural networks are again based on studies of the human brain, and if the artificial neurons work as the biological ones, there should be possible to train a neural network recognizing sounds as well. Lately, immense efforts have been put into this subject with the purpose of translating voice into text, leaded by technology companies like Google, Microsoft, Amazon and so on; who develop voice controlled virtual assistants. The technology is promising; the time saving using voice commands compared to the keyboard commands is potentially large and has a great market[22].

Artificial neural networks that use an audio signal as an input have received a great research interest in the past ten years and have shown potential in speech recognition and environmental sound classification. Piczak[23] tested a very simple CNN architecture with environmental audio data and achieved accuracies comparable to state-of-the-art classifiers. Cakir et al.[24] used one dimensional (time domain) deep neural networks (DNNs) in polyphonic sound event detection for 61 classes to achieve an accuracy of 63.8%, which was a 19% improvement over a hybrid Hidden Markov Model/Non-negative Matrix Factorization method. Lane et al.[25] created a mobile application capable of performing very accurate speaker diarization and emotion recognition using deep learning. Recently, Wilkinson et al.[26] performed unsupervised separation of environmental noise sources adding artificial Gaussian noise to

---

[22] Jiang, Jiepu, Ahmed Hassan Awadallah, Rosie Jones, Umut Ozertem, Imed Zitouni, Ranjitha Gurunath Kulkarni, and Omar Zia Khan. "Automatic online evaluation of intelligent assistants." In *Proceedings of the 24th International Conference on World Wide Web*, pp. 506-516. International World Wide Web Conferences Steering Committee, 2015.

[23] Piczak, Karol J. "Environmental sound classification with convolutional neural networks." In *2015 IEEE 25th International Workshop on Machine Learning for Signal Processing (MLSP)*, pp. 1-6. IEEE, 2015.

[24] Cakir, Emre, Toni Heittola, Heikki Huttunen, and Tuomas Virtanen. "Polyphonic sound event detection using multi label deep neural networks." In 2015 international joint conference on neural networks (IJCNN), pp. 1-7. IEEE, 2015.

[25] Lane, Nicholas D., Petko Georgiev, and Lorena Qendro. "DeepEar: robust smartphone audio sensing in unconstrained acoustic environments using deep learning." In Proceedings of the 2015 ACM International Joint Conference on Pervasive and Ubiquitous Computing, pp. 283-294. ACM, 2015.

[26] Wilkinson, Bryan, Charlotte Ellison, Edward T. Nykaza, Arnold P. Boedihardjo, Anton Netchaev, Zhiguang Wang, Steven L. Bunkley, Tim Oates, and Matthew G. Blevins. "Deep learning for unsupervised separation of environmental noise sources." The Journal of the Acoustical Society of America 141, no. 5 (2017): 3964-3964.

pre-labeled signals and used auto-encoders to cluster. However, background noise in an environmental signal is usually non-Gaussian, making this method to work on specific datasets only.

Regarding the field of surveillance using audio sensors, there has been extensive research in order to robustly identify the audio scene that an event takes place. The main reason is that audio sensors (microphones) offer the following advantages:

- microphones are cheap, compared to the cameras and can be easily deployed in any kind of environment
- omnidirectional coverage
- specular reflections of the audio signal can be another form of audio input[27]

However, the main challenge that this research area is facing is that environmental sounds are unstructured, especially compared with the speech signals where a phoneme-based approach can be used for classification. The signal-to-noise ratio (SNR) is typically small for an environmental signal, especially when the device is placed far away from the acoustic source. Additional challenges include the recognition of overlapping events from one environment[28], audio event recognition using weakly labeled data[29] and lack of public datasets that contain information from multiple sensors[30].

In this report, we focus on the MIVIA audio events dataset[31], where sounds such as background noise, glass break, gun shot and scream are differentiated. After providing information about the dataset, we see how we can analyze the audio samples in the time and frequency domain and provide details about audio feature extraction. In the Classification Methods Section, we mainly focus on two-dimensional Convolutional Neural Networks (CNNs) that have shown great performance in recognition accuracy[32]. The focus of the work was to study the ability of 2D CNNs to generalize in different SNR conditions. For this work, the Adam[33] optimizer was used and the ReLU[34] activation function between the convolutional layers. Finally, all the accuracy scores are presented in the Results Section.

### MIVIA Audio Events Dataset

The MIVIA audio events data set is composed of 6000 events for surveillance applications, namely glass breaking, gun shots and screams. The 6000 events are divided into a training set (composed of 4200 events) and a test set (composed of 1800 events).

---

[27] Baum, Elizabeth, Mario Harper, Ryan Alicea, and Camilo Ordonez. "Sound identification for fire-fighting mobile robots." In 2018 Second IEEE International Conference on Robotic Computing (IRC), pp. 79-86. IEEE, 2018.

[28] Dikmen, Onur, and Annamaria Mesaros. "Sound event detection using non-negative dictionaries learned from annotated overlapping events." In 2013 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics, pp. 1-4. IEEE, 2013.

[29] Kumar, Anurag, and Bhiksha Raj. "Audio event detection using weakly labeled data." In Proceedings of the 24th ACM international conference on Multimedia, pp. 1038-1047. ACM, 2016.

[30] Chachada, Sachin, and C-C. Jay Kuo. "Environmental sound recognition: A survey." APSIPA Transactions on Signal and Information Processing 3 (2014).

[31] Foggia, Pasquale, Nicolai Petkov, Alessia Saggese, Nicola Strisciuglio, and Mario Vento. "Reliable detection of audio events in highly noisy environments." Pattern Recognition Letters 65 (2015): 22-28.

[32] Zhao, Jianfeng, Xia Mao, and Lijiang Chen. "Speech emotion recognition using deep 1D & 2D CNN LSTM networks." Biomedical Signal Processing and Control 47 (2019): 312-323.

[33] Kingma, Diederik P., and Jimmy Ba. "Adam: A method for stochastic optimization." arXiv preprint arXiv:1412.6980 (2014).

[34] Nair, Vinod, and Geoffrey E. Hinton. "Rectified linear units improve restricted boltzmann machines." In Proceedings of the 27th international conference on machine learning (ICML-10), pp. 807-814. 2010.

In audio surveillance applications, the events of interest (for instance a scream) can occur at different distances from the microphone that correspond to different levels of the signal-to-noise ratio. Moreover, in these applications the events are generally mixed with a complex background, usually composed of several types of different sounds depending on the specific environments both indoor and outdoor (household appliances, cheering of crowds, talking people, traffic jam, passing cars or motorbikes etc.).

The data set is designed to provide each audio event at eight different values of signal-to-noise ratio (namely -5dB, 0dB, 5dB, 10dB, 15dB, 20dB, 25dB and 30dB) and over-imposed to different combinations of environmental sounds in order to simulate their occurrence in different ambiences.

The sounds have been registered with an Axis P8221Audio Module and an Axis T83 omnidirectional microphone for audio surveillance applications are, sampled at 32000 Hz and quantized at 16 bits per PCM sample. The audio clips are distributed as WAV files. The training set has a duration of about 20 hours while the test set of about 9 hours.

The events of interest are organized in three classes (glass breaking, gun shots and screams) and their duration in the training and test sets is reported in the following table

**Table 1. The MIVIA audio events dataset**

|  | TRAINING SET | | TEST SET | |
| --- | --- | --- | --- | --- |
|  | #Events | Duration (s) | #Events | Duration (s) |
| **Background** | 10768 | 58371,6 | 4716 | 25036,8 |
| **Glass breaking** | 4200 | 6024,8 | 1800 | 2561,7 |
| **Gun shots** | 4200 | 1883,6 | 1800 | 743,5 |
| **Screams** | 4200 | 5488,8 | 1800 | 2445,4 |

In our experiments, we have focused on the 0 dB SNR, which is the second most challenging SNR value. At 0 dB the background noise can significantly mask the target sound event, such as the scream. Finally, we notice that the dataset is quite imbalance towards the background noise. Therefore, our system needed to be robust against the imbalances and the error was evaluated with the F1 macro average score defined as

$$F1_{macro} = \frac{2 \times Precision_{macro} \times Recall_{macro}}{Precision_{macro} + Recall_{macro}} \qquad \textbf{(1)}$$

By macro averaging, we refer computing the metric independently for each class and then take the average, hence treating all classes equally.

**Time Domain:** Sounds are longitudinal waves, which actually are different pressures in the air. They are usually represented as a function in the time domain, which means how the pressure vary per time. This function is obviously continuous, but since computers represent functions as arrays, we cannot save all the information.

When discretizing the audio signal, we lose some information. The loss of the information depends on the sampling rate, which is the number of sampling points per second (Figure 1). According to the Nyquist theorem, one should have twice as high sampling rate as the highest sound frequency to keep most of the information and avoid clipping of the signal. For instance, a human ear can perceive frequencies in the range 20-20000Hz, so a sampling rate around 40 kHz should be sufficient to keep all the information. Ordinary CD's use a sampling rate of 44.1 kHz, but for speech recognition, 16 kHz is enough to cover the frequency range of human speech.

**Figure 18. Quantized audio signal (continuous to discrete)**

**Frequency Domain:** Most of the times, a frequency representation of the signal provides a better picture than the time domain. On one hand, one losses information in the time domain, but on the other one, gets information about which frequencies are significant in the original audio signal. To transform from the time domain to the frequency domain, the Fourier Transformation is performed, defined by

$$\hat{f}(x) = \int_{-\infty}^{\infty} f(t)\, e^{-2\pi i x t} dx \quad \textbf{(2)}$$

In the above equation, $f(t)$ is the time function and $\hat{f}(x)$ is the frequency function. The Fast Fourier Transform (FFT) has been proposed in order to achieve higher speeds in calculating the Fourier Transformations by splitting the input signal into many sub-signals. In Figure 2, the FFT is calculated on the time domain in order to obtain the corresponding frequency domain.

**Figure 19. Waveforms of three harmonic functions (top) and their corresponding frequency responses (bottom)**

**Feature Extraction:** To include the range of frequencies that are relevant to identifying the given environmental sounds and to efficiently extract the audio features, we split the input signal into smaller frames for processing. First, we down-sampled the original 32 kHz audio signal to 16 kHz that allowed faster processing. Each frame had an FFT window size of 512 with a 256 hop length (50 % overlap). Therefore, we resulted in a 188 x 257 grayscale spectrogram image, for a 3 s recording, that was used as an input to the 2D CNN classifier.

(a)



(b)

**Figure 20.  Time (a) and frequency (b) representations of the background noise at 0 dB SNR**

(a)



(b)

**Figure 21. Time (a) and frequency (b) representations of a glass breaking at 0 dB SNR**

(a)



(b)

**Figure 22. Time (a) and frequency (b) representations of a gunshot at 0 dB SNR**

(a)



(b)

**Figure 23. Time (a) and frequency (b) representations of a person screaming at 0 dB SNR**

Figures 20-23 show the time and frequency domain representations of the selected target classes at 0 dB SNR (noisy signal). One can notice that the waveform of the background noise and the glass breaking look very similar in a very noisy environment. However, when observing the FFT spectrogram, the differences are clearer, since there are some high energies concentrated in frequencies between 4 and 10 kHz. Additionally, in this work, we have studied the audio signals with SNR values varying between -5 dB and 30 dB with 5 dB increments in between. Figures 24-27 show the time and frequency domain representations of the selected target classes at 30 dB SNR (clean signal). We notice that the signal at 30 dB SNR is much cleaner, compared to the one at 0 dB and therefore, the network would be able to classify the target classes with higher accuracy.



(a)



(b)

**Figure 24. Time (a) and frequency (b) representations of the background noise at 30 dB SNR**

(a)

(b)

**Figure 25. Time (a) and frequency (b) representations of a glass breaking at 30 dB SNR**

(a)



(b)

**Figure 26. Time (a) and frequency (b) representations of a gun shot at 30 dB SNR**

(a)



(b)

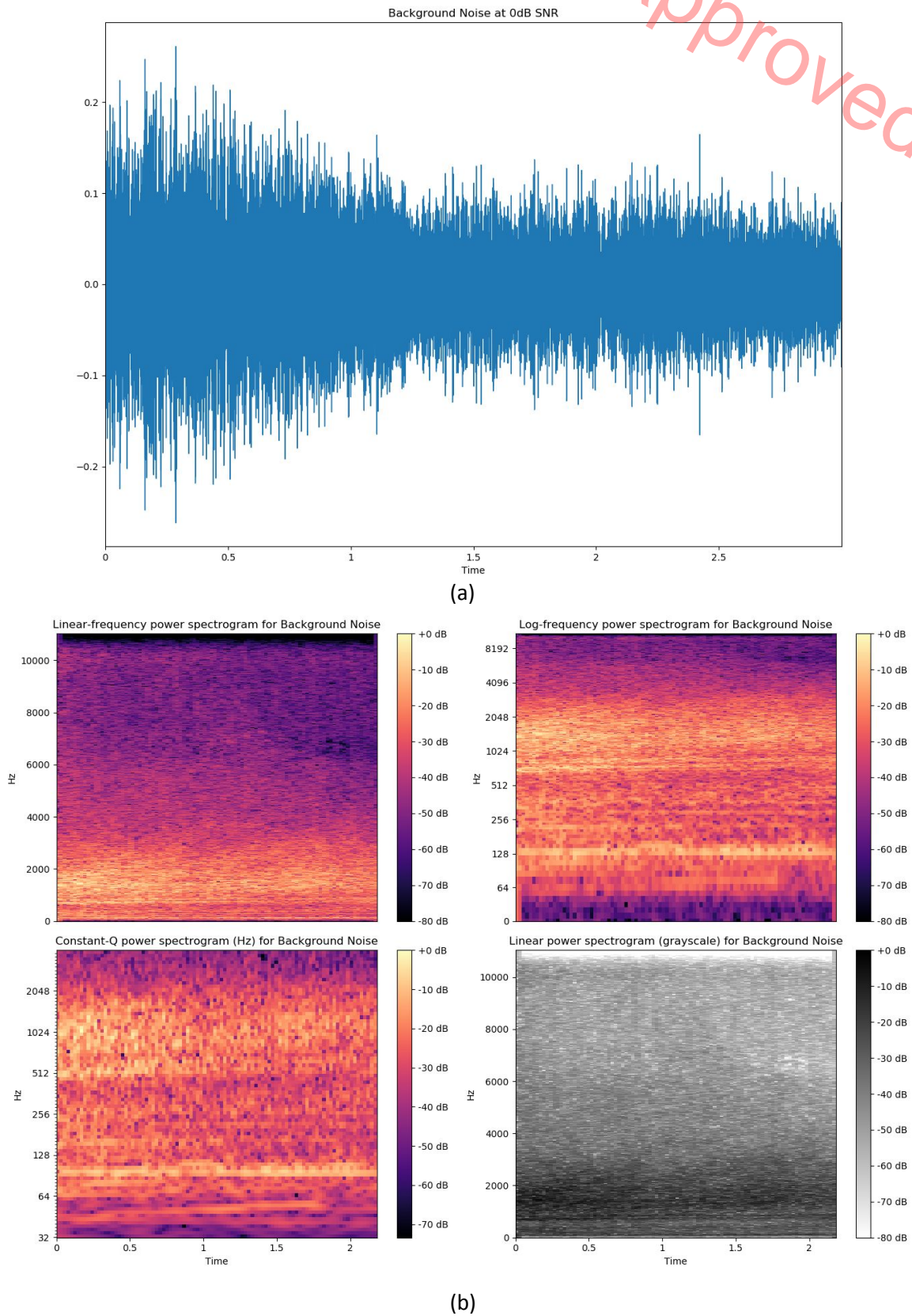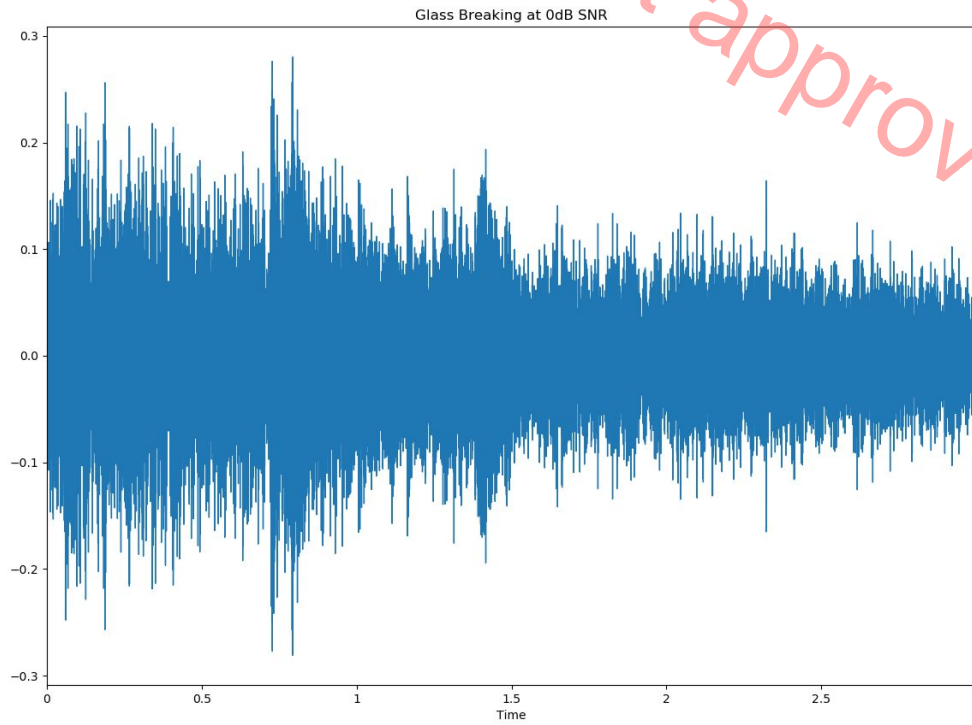**Figure 27. Time (a) and frequency (b) representations of a person screaming at 30 dB SNR**

## Classification Methods

In this Section, we will describe the Convolutional Neural Network architectures used, the optimizer that was selected and the activation function between the convolutions and max-pooling operations. CNNs are known to perform well in image classification and have achieved similar results in the field of audio-based event detection. In general, the input image is passed through various convolutional layers, before it is flattened and fed to a fully connected neural network that outputs the probabilities of the target classes.

**Convolutional Layer:** Initially, the image is passed through a convolutional layer, which is meant to reveal the structures and shapes in the image. They way that this is performed is to introduce a filter that slides over the entire image and multiplies with all the pixels (with overlap). Every time the filter is multiplied with a set of pixels, we sum all the multiplications and add the value to an activation map (convolution operation). The activation map is completed after the filter is multiplied with the entire image. A typical filter has dimensions of 16 × 16 or 32 × 32, depending on the shape of the input image. It is important in this case, to choose a filter that is large enough to cover all the structures of the image.

**Pooling Layer:** A pooling layer is a way to reduce dimensionality of the representation (down-sampling) such that there are not many parameters that need optimization, but it also helps to control overfitting. The activation map is divided into region of equal size and represent each region with one single number. Max-pooling is one of the most popular pooling techniques, which just represents each region

with the largest number in that region. However, it is possible to use average pooling, global pooling, etc.

**Batch Normalization:** To increase the stability of a neural network, batch normalization normalizes the output of a previous activation layer by subtracting the batch mean and dividing by the batch standard deviation. Consequently, batch normalization adds two trainable parameters to each layer, so the normalized output is multiplied by a "standard deviation" parameter ($\gamma$) and add a "mean" parameter ($\beta$).

**Dropout:** Dropout is widely used in neural network layers and is another way to prevent overfitting. This layer is simple in the sense that it randomly drops out units (activation maps) in the current layer by setting them to zero.

**Fully Connected Layer:** The last part of a convolutional network is often referred to as a fully connected layer, which is a regular feed-forward neural network (FNN). The output from the other layers needs to be flattened out before it is passed into the FNN. The activation function that is commonly used for classification is Softmax. This activation function assigns probabilities between the target classes and the probability that is closest to one is selected.

**DenseNet-121:** As a first part of our experiments, we have used a DenseNet-121 CNN architecture as seen in Figure 28 and Figure 29.



188 × 257 grayscale
spectrogram image

**Figure 29. Part of the DenseNet-121 architecture as depicted using the Netron software**

DenseNets[35] were introduced in order to solve the problem of the vanishing gradient in neural networks. The vanishing gradient occurs when the CNNs become big that the path for information from the input layer to the output layer significantly increases that the gradient vanishes during back-propagation. In DenseNets, we connect every layer directly with each other. This process ensures maximum information kept throughout the architecture. Additionally, by using this connection the DenseNets require fewer parameters to train than vanilla CNNs, since there is no need to learn redundant feature maps.

The DenseNet-121 was used for classification, as the most basic DenseNet yet powerful architecture. Each dense layer consists of two convolutional operation as follows
- 1 × 1 Convolutions (for feature extraction)
- 3 × 3 Convolutions (bringing down the feature depth/channel output)

The DenseNet-121 consists of six such dense layers in a dense block. This resulted in approximately seven million parameters, compared to the 44 million parameters of a vanilla CNN architecture.

**Optimization**

After the non-linear function describing the dataset is found, one typical wants to find the minimum or maximum to optimize various parameters. In order to perform that there are numerical optimization methods. In our experiments, we have used the Adam optimization method with an initial learning rate of 0.001.

Adam is often the preferred optimization method in machine learning, due to its computational efficiency, little memory occupation and implementation ease. Adam is a momentum-based method that only relies on the first derivative of the cost function. Required inputs are exponential decay rates $\beta_1$ and $\beta_2$, cost function c($\theta$) and initial weights $\theta$, in addition to the learning rate $\eta$ and eventually a

---

[35] Huang, Gao, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q. Weinberger. "Densely connected convolutional networks." In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 4700-4708. 2017.

regularization $\lambda$. Iteratively, the gradient of the cost function is calculated and this is used to calculate the first and second moment estimates.

### Activation

The activation function is used to activate the outputs, which often need to take some certain properties. For example, when doing classification, the outputs are probabilities and therefore take values between zero and one. A nonlinear activation function is often preferred to reinforce the non-linearity of the neural network.

Traditionally, the logistic function and the tanh function have been used as activation functions, since they were believed to work in the same way as the human brain. However, in 2012 Alex Krizhevsky et al. introduced AlexNet, taking image recognition to a new dimension. They used a function named Rectified Linear Units (ReLU), which is zero for negative values. The ReLU function is a modification of the pure linear function for positive values. This make the function able to recognize the non-linearity in the model, providing it a "clean" derivative given by the step function.

## 2.4.3.1.3    Research results

In the next section, we will present results from the vision and audio modalities. Future research should consider the potential effects of fusing video and audio modalities. Models can be fused both on decision-level and by concatenating their respective fully connected layers. Recent studies by Kampman et al.[36] and Ortega et al.[37] has proven that using multiple modalities combined allows interaction between them in a non-trivial way and greatly outperforms the individual performance. By combining the last network layers and fine-tuning the parameters, we can take advantage of the complementary information of visual and auditory modalities outperforming the current individual results.

### Video analysis results

We tested our model at the NTU RGB+D dataset and on the data captured by CERTH inside the AV's shuttle (Figure 30). In the images there are various debug layers enabled, such as skeleton points, lines, tracker id and bounding boxes of each detection. The predicted result is marked as green, when the classifier indicates it as "normal" and red when "abnormal", correspondingly. So far, we did not include NTU dataset samples in our training set, so it is safe to assume that our model can generalize across different people, view angles etc. Also in Figure 31, Figure 32, Figure 33 and Figure 34, we present some results on each use case along with the link to the full video. Classification results and metrics for each class are presented in Table 2.

### LSTM Classification results

P:    Person Identifier (e.g P1 for the Person with ID1)

Normal behaviour          Abnormal behaviour

---

[36] Kampman, Onno, et al. "Investigating audio, video, and text fusion methods for end-to-end automatic personality prediction." Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers). 2018.

[37] Ortega, Juan DS, et al. "Multimodal fusion with deep neural networks for audio-video emotion recognition." arXiv preprint arXiv:1907.03196 (2019).

**Figure 30. Evaluation on test data:**
**(a), (b), (c) Abnormal event detection (violence/passengers are fighting) using different camera angles from the NTU-RGB dataset.**
**(d), (e) Detection of Fighting / Bag Snatch real-world scenarios inside the shuttle.**

**Evaluation on Use Cases**

P: Person Identifier (e.g P1 for the Person with ID1)

Normal behaviour          Abnormal beh:

**Use Case 2 – Bag Snatch**

**Figure 31. The commuter is attacked by another person who is attempting to snatch his bag**
Full video: https://drive.google.com/open?id=15tLAGcjp5pzI5erHRcg0u2Mh53PSeIkt

**Use Case 2 – Fighting 1**



**Figure 32. A petty crime takes place in the form of assault**
Full video: https://drive.google.com/open?id=1fGktDygbtxWT20ILfqBXPs3fdy8mXozi

**Use Case 3 – Vandalism**

**Figure 33. A person attempts to perform a vandalism action on the bus, through painting a graffiti on the windows**
Full video: https://drive.google.com/open?id=125B7YX1w_yKqvr6zyiai22urGXK0nV0j

**Use Case 1 – Unaccompanied Luggage Monitoring**



**Figure 34. Unaccompanied luggage which remains as is unmoved for a long time**
Full video: https://drive.google.com/open?id=1rCV0nYH4LQwOv4hlw1vyeXpXqcdrByMw

Detection of certain events may raise a notification or an alert to the supervisor and/or the suitable authorities. This may be followed by appropriate notifications and/or instructions to the passengers, while the vehicle may also implement respective actions.

While a commuter is in the shuttle a petty crime takes place in the form of assault (Figure 31 – video: https://drive.google.com/open?id=15tLAGcjp5pzI5erHRcg0u2Mh53PSeIkt). Also, the commuter is attacked by another person who is attempting to snatch his bag (Figure 12 - video: https://drive.google.com/open?id=1fGktDygbtxWT20ILfqBXPs3fdy8mXozi). The system identifies the event and sends a security alert to the operator or security supervisor. The course of action of the operator is a human decision, that meaning, whether he/she will decide to stop the minibus or will notify the passengers via the radio system.

A young person onboards the shuttle during an itinerary performed by the vehicle during the night hours (Figure 33: https://drive.google.com/open?id=125B7YX1w_yKqvr6zyiai22urGXK0nV0j). The youngster attempts to perform a vandalism action on the shuttle, through painting a graffiti on the windows or smashing the window. The event is detected and the person is warned by the radio system of the shuttle or security personnel intervenes by stopping the shuttle.

Inside the shuttle there is unaccompanied luggage which remains as is unmoved for a long time (Figure 34: https://drive.google.com/open?id=1rCV0nYH4LQwOv4hlw1vyeXpXqcdrByMw). In case that the total time of the detected luggage that remains unmoved in the vehicle passes the predefined time frame, a notification is sent to the security operator. The security operator monitors the clips that are captured and evaluates the criticality of the situation and whether to intervene or not.

We compiled a 3 min single video which includes detections from all the above scenarios (video link: https://drive.google.com/drive/folders/1rBm3Ss2XuwbZhiVm5MS8QRYFMuNYrrVu)

**Table 2. Precision, Recall and F1-Score metrics for the two classes**

| Classification Report | | | |
|---|---|---|---|
| *Class* | *Precision* | *Recall* | *F1-Score* | *Support* |
| Normal | 0.99 | 0.99 | 0.99 | 1208 |
| Abnormal | 0.93 | 0.95 | 0.94 | 147 |
| | | | | |
| accuracy | | 0.99 | | 1355 |
| macro avg | 0.96 | 0.97 | 0.96 | 1355 |
| weighted avg | 0.99 | 0.99 | 0.99 | 1355 |

**Spatiotemporal Autoencoder results**

In order to visualize the predictions, we provide the preprocessed input frame for the current moment at the bottom-left. The frame is resized from 64x64 and grayscale and a mask overlay is obscuring the out-of-interest areas (road). At the right next to the input frame, the resized output of our model is shown. The third mini-frame demonstrates only the significant differences between the input and the output frames with white pixels which is then merged above the original frame for demonstration purposes.

**Figure 35. Evaluation on Use Case 2 – Fighting Scenario (1)**



**Figure 36. Evaluation on Use Case 2 – Fighting Scenario (2)**

## Sound analysis results

Regarding our experiments, we used a batch size of 16 images and set the initial epochs to 200. However, we noticed that for the DenseNet-121 only eight epochs were sufficient to achieve the optimal performance. We applied an Early Stopping function, where we checked the validation F1 macro averaged score for an improvement in five consecutive epochs. If no improvement was detected the network stopped the training in order to avoid overfitting.

The DenseNet-121 achieved a training F1-Score of 95.92% and a validation F1-Score of 88.74% (Figure 13-left) for the case of 0 dB SNR. Regarding the case of 30 dB, the DenseNet-121 achieved a training F1-Score of 96.84% and a validation F1-Score of 91.82% (Figure 13-right). As expected, the network is able to classify the target classes with higher accuracy in the case of 30 dB SNR and we notice that the training loss starts at smaller values, compared to the case of 0 dB SNR.

There are plenty of options to tweak in a CNN. The number of convolutional layers, max-pooling layers, filter sizes, activation functions. In our experiments, we focus on comparing the default DenseNet-121 architecture in various SNR settings.



**Figure 37. F1-Score (a) and categorical cross-entropy loss (c) for the 0 dB case of the DenseNet-121 architecture. F1-Score (b) and categorical cross-entropy loss (d) for the 30 dB case of the DenseNet-121 architecture**

Since the train and validation losses and F1-Scores are not sufficient for the complete evaluation of the proposed framework, we evaluated the Precision, Recall, and F1-Score for each class (Table 2) and calculated the receiver operating characteristic (ROC) curves for each class (Figure 14).

**Table 3. Precision, Recall and F1-Score metrics for the four classes at different SNR levels**

| Classification Report | | Classes | | | | | | | | | | |
| | | Background Noise | | | Glass Breaking | | | Gun Shot | | | Scream | | |
| | | Precision | Recall | F1-Score | Precision | Recall | F1-Score | Precision | Recall | F1-Score | Precision | Recall | F1-Score |
| SNR values | -5 dB | 79 % | 86 % | 83 % | 74 % | 82 % | 78 % | 89 % | 87 % | 88 % | 85 % | 60 % | 70 % |
| | 0 dB | 87 % | 90 % | 89 % | 86 % | 90 % | 88 % | 97 % | 96 % | 96 % | 88 % | 77 % | 82 % |
| | 5 dB | 91 % | 85 % | 88 % | 80 % | 94 % | 87 % | 92 % | 94 % | 93 % | 88 % | 87 % | 87 % |
| | 10 dB | 89 % | 91 % | 90 % | 89 % | 90 % | 90 % | 99 % | 97 % | 98 % | 88 % | 83 % | 86 % |
| | 15 dB | 93 % | 89 % | 91 % | 88 % | 93 % | 90 % | 97 % | 98 % | 98 % | 88 % | 90 % | 89 % |
| | 20 dB | 92 % | 90 % | 91 % | 88 % | 93 % | 90 % | 99 % | 99 % | 99 % | 89 % | 89 % | 89 % |
| | 25 dB | 91 % | 91 % | 91 % | 87 % | 89 % | 88 % | 99 % | 99 % | 99 % | 91 % | 89 % | 90 % |
| | 30 dB | 92 % | 90 % | 91 % | 88 % | 90 % | 89 % | 99 % | 99 % | 99 % | 88 % | 91 % | 90 % |

From Table 2 we can see that while the SNR value increases, the network is able to distinguish the four classes more accurately. Specifically, for all four classes, the network achieves the highest F1-Score for SNR values higher than 15 dB.

The class "gun shot" is one of the most easily distinguishable, while the "scream" class is the hardest to classify. The results from the classification report can also be summarized in the following ROC curve figure (Figure 14). Classes 0, 1, 2 and 3 correspond to background noise, glass breaking, gun shot and scream respectively.

(a) − 5 dB

(b) 0 dB

(b) 5 dB

(d) 10 dB

(e) 15 dB

(f) 20 dB

(g) 25 dB

(h) 30 dB

**Figure 38. ROC curves for the four classes at the selected SNR values**

Finally, in order to test the generalizability of the selected DenseNet architecture, we tested the network under three settings. The first one with the network trained on 30 dB SNR and tested on 0 dB SNR. The second one with the network trained on 0 dB SNR and tested on 30 dB SNR and finally with the network trained on 15 dB SNR and tested on 30 dB SNR. The classification report are summarized in Figures 15 – 17.

```
Classification Report
                 precision    recall  f1-score   support

Background Noise      0.48      0.99      0.64      2358
  Glass Breaking      0.80      0.11      0.20       900
        Gun Shot      1.00      0.00      0.00       900
          Scream      0.82      0.04      0.07       900

        accuracy                          0.49      5058
       macro avg      0.78      0.29      0.23      5058
    weighted avg      0.69      0.49      0.35      5058
```

**Figure 39. Classification report for DenseNet-121 trained on 30 dB and tested on 0 dB SNR**

```
Classification Report
                 precision    recall  f1-score   support

Background Noise      0.71      0.84      0.77      2358
  Glass Breaking      0.80      0.96      0.88       900
        Gun Shot      0.97      0.28      0.44       900
          Scream      0.83      0.86      0.84       900

        accuracy                          0.77      5058
       macro avg      0.83      0.74      0.73      5058
    weighted avg      0.79      0.77      0.74      5058
```

**Figure 40. Classification report for DenseNet-121 trained on 0 dB and tested on 30 dB SNR**

```
Classification Report
                 precision    recall  f1-score   support

Background Noise      0.95      0.84      0.89      2358
  Glass Breaking      0.82      0.98      0.89       900
        Gun Shot      0.98      0.99      0.99       900
          Scream      0.84      0.92      0.88       900

        accuracy                          0.91      5058
       macro avg      0.90      0.93      0.91      5058
    weighted avg      0.91      0.91      0.91      5058
```

**Figure 41. Classification report for DenseNet-121 trained on 15 dB and tested on 30 dB**

From Figure 15, we notice that the network trained on an environment with the least noise cannot distinguish the classes in a noisy environment. On the other hand, the network trained on a noisier environment (Figure 17) can distinguish the classes in the quietest environment settings almost as well as the network trained on the same environmental settings. Therefore, we notice that our network can generalize well in clean environments when trained in noisy ones.

Despite the classification reports and the ROC curves, we used the T-distributed Stochastic Neighbor Embedding (t-SNE) plots to further visualize the automatic features that were learnt by the proposed 2D

CNN architecture. The advantage of the t-SNE visualization, compared to a Principal Component Analysis, is that it uses the local relationships between points to create a low-dimensional mapping. This allows the t-SNE to capture non-linear structure of the given dataset (raw data and learnt features), since the neural network is learning non-linear representations of the dataset. Figure 18 shows the t-SNE visualization at 0 dB (noisy environment) and Figure 19 shows the t-SNE visualization at 30 dB (quiet environment).

1: Background Noise, 2: Glass Break, 3: Gun Shot, 4: Scream

1: Background Noise, 2: Glass Break, 3: Gun Shot, 4: Scream

(a)

(b)

1: Background Noise, 2: Glass Break, 3: Gun Shot, 4: Scream

1: Background Noise, 2: Glass Break, 3: Gun Shot, 4: Scream

(a)

(b)

From the above Figures we notice the randomness in the 2D space of the raw features, in both environments and on the right part of the Figures the ability of the proposed 2D CNN architecture to distinguish between the target classes and create clear clusters for each class.

Before deploying the system for the real-life application, more experiments will be conducted with other network architectures (e.g., other 2D CNNs, 1D CNNs, RNNs), ensemble methods will be explored and

statistical test will be performed in order to determine whether the results are independent or there is a degree of randomness.

## 2.4.3.2 Equipment

At previous sections, we presented three different solutions for the vision modality. In this section, we are going to compare those solutions in terms of requirements and equipment and enlist the advantages and drawbacks of each method.

Stacked LSTM classification is a supervised method that classifies the extracted skeleton key points based on a (supervised) dataset. It can reliably detect various types of events but depends on a previous skeleton extraction and tracking process, which may not be accurately feasible due to space and occlusion constraints. It does not depend on the camera setup (but it's less effective on top-down cameras with a wide field of view). The Convolutional LSTM autoencoder (unsupervised) extracts spatiotemporal features from a video sequence, able to learn regularity. It can detect abnormal events and activities and depends on the camera setup but due to its unsupervised nature it can be trained over-time and self-improve. It is ideal for crowded areas and static camera setups. In Hybrid classification (semi-supervised), the previous encoder acts as a high recaller, and the anomalies are sent through a false positive reduction model (hybrid model). This combination provides a deep neural network with high recall and high precision.

We developed different solutions due to the perspective of the final camera setup and the technical specifications of the lens/sensor (Table 4). Existing in-shuttle camera have a wide-angle field of view and a top-down perspective which is not suitable for pose estimation. In some occasions the whole body of the commuter can be occluded by his head, rendering impossible for the pose estimation algorithms to detect his pose in 2D. Although the LSTM Classification via Pose Estimation performs generally better, it cannot be applied with satisfactory results in the existing in-shuttle camera. We developed the autoencoder based solutions to prevent excluding this possibility. Autoencoders do not depend on pose estimation and extract their own features by learning regularity.

Taking the above into account, the relevant equipment for each case can be derived. The choice of the final solution is based on these characteristics, as well as on the availability of connection with the system and real-time data frames.

**Table 4. Equipment and proposed solutions**

| Requirements | Custom installation by CERTH | Integration to NAVYA's equipment |
|---|---|---|
| **Solution** | LSTM Classification via Pose Estimation | Spatiotemporal Autoencoder/Hybrid LSTM Classification |

| | | |
|---|---|---|
| **Camera** | 2x FullHD camera sensors (require installation) | Existing in-shuttle camera |
| **Camera Power supply** | Using USB protocol (requires wiring) | Already exists |
| **Data transmission** | Using USB protocol (requires wiring) | Requires access to the video stream via a standard protocol (e.g RTSP) or a documented API |
| **Host PC** | In-shuttle PC capable of real-time inference | In-shuttle PC capable of real-time inference |
| **PC Power supply** | 1000W (max) | 1000W (max) |

## 2.4.4 Prototyping plan

The prototyping plan of this service consists of the development phase and the deployment phase:

**Development phase:** In this phase, CERTH initially plans and analyses the requirements with inputs from all the stakeholders. Once the requirement analysis is done, the final software and hardware representation and documentation of the requirements are accepted from the project stakeholders. The next stage consists of designing the critical service components, the research and the development of the algorithms and the integration of the components across several platforms. Also, in this stage several data captures are performed for training the machine learning algorithms. CERTH will perform additional data capture sessions with different lighting conditions and operators. The captured data should demonstrate multiple scenarios, both with abnormal and normal events, and various capturing conditions. Also, this phase includes the re-evaluation of the results, using the new data acquired for retraining the machine learning algorithms.

**Deployment phase:** In the deployment phase, the individual service components are unified into a complete system. The system has discrete input and output and is ready to be integrated with other platforms. In the following months, CERTH will deploy the service in a demo AV, perform internal tests and verifications with all the stakeholders. In this stage, minor modifications and fine-tuning may be applied on the final setup, mainly on design or algorithms of the service, depending on the real operation conditions.

**Prototyping phase:** Amobility will in coordination with CERTH install cameras and sensors in the Amobility shuttles for testing and validation of the technologies, use cases etc.

## 2.4.5 Result analysis

In the evaluation phase, the service will be tested under real conditions, that can be performed on controlled routes of Amobility along with a safety driver inside the shuttle, depending on the GDPR permissions. Also, short-sessions for assessment with real passengers will be available. After the successful evaluation of the service in Amobility sites, the results will be used to fine tune the service if required and the service will be deployed and evaluated also in other sites of the involved operators in the AVENUE project towards a successful integration to the relevant pilot sites.

# 2.5 Service: Automated passenger presence

## 2.5.1 Concept of service

The service "Automated passenger presence" aims to address a basic problem of operators' services which is related to the occupation of their vehicles as well as the awareness of the number of people on-board in order to schedule the routes. Furthermore, the passengers would like to know in advance if there is an available seat or enough space on a shuttle to plan their boarding. Traditionally, but also nowadays, passenger counting is conducted manually via passenger surveys or human ride checkers. Typically, the driver or inspectors are responsible for performing enumeration of the onboard passengers, something not feasible in an autonomous shuttle. Automatic passenger counting has been rapidly emerging in recent years to address similar needs. An automated system is introduced capable to detect passenger presence in real-time with high accuracy, count onboard passengers and calculate vehicle occupancy. Surveillance using sensors such as cameras (cameras of different technologies can be used so that passengers' privacy is protected) and smart software in the bus will automate the detection of passenger presence.

Several concerns of the end users regarding the Safety and Robustness of the autonomous vehicles that are directly linked to the final User Acceptance of the new technology, can be identified. The prospective passengers may deal with several possible instances that could arise in case there is no staff inside the shuttle. Indicatively:

● No one will be in the bus to count the number of passengers with regard to the shuttle's capacity
● Continuous stops throughout the entire route, even in cases where the shuttle is fully occupied
● No authority figure present to alert passengers of their designated bus stop

To address the aforementioned concerns on social and personal safety and quality-of-service into the vehicle, certain measures need to be implemented. For example, counting the number of passengers being inside the autonomous vehicle could help avoid overcrowding in the shuttle, as well as meaningless stops in cases where the bus is in full capacity. This may be followed by appropriate notifications and/or instructions to the passengers, while the vehicle may also implement respective actions.

The service provides a video analysis of the vehicle internals, using the on-board camera, in order to identify the vehicle occupation, vehicle free space, as well as counting people on-board. Automatic assessment of space occupation using the on-board cameras is enabled. Capacity is set as an absolute number of space units. For example, each space unit is associated with one standing passenger. Occupancy is set as an absolute number of space units currently in the shuttle. For the operation manager, occupancy is visible on the dashboard of the AVENUE platform, whereas occupancy is displayed as real time information via the AVENUE mobile app, wherever the traveler is. Each passenger (normal, big size, wheelchair user, seated) can determine whether he/ she can fit in or not. Assessment for different cases can be provided to assist the passengers on determining whether to request onboarding or not. Automatic counting of people using the on-board cameras is also provided. Additional information can be derived from automated people-counting, while fusion of data related to

space occupation and counted number of people can provide more accurate information about the capacity and occupancy of the vehicle. Moreover, occupancy marked with information for the different user cases is displayed as real time information, wherever the traveler is, however it does not guarantee them a free spot by the time the shuttle reaches the station of their choice.

In this section, the implementation of a video analytics software module for an automated passenger presence counting subsystem or for cloud-based services of the system are described along with appropriate planning for the deployment and test of the service into the pilot sites of the AVENUE project.

### 2.5.1.1    Use case

The service addresses the timely, accurate, robust and automatic counting of the passenger number within the autonomous vehicle, as well as appropriate notifications and/or instructions to its passengers. The main goal is to monitor the number of passengers onboard at all times, so that the shuttle's capacity it is not exceeded at any time. The passenger would like to know if there is an available seat or space before getting on-board on the autonomous shuttle. Occupancy is displayed as real time information, wherever the traveler is, but it does not guarantee them a free spot when the shuttle arrives at the station of their choice. Within the suggested software module, analytics algorithms were developed for the timely, accurate, robust and automatic detection of onboard passengers. The service is able to estimate occupied seats and report the passenger capacity and vehicle occupancy of the autonomous shuttle continuously.

In the context of AVENUE project, the following use cases have been identified to be further examined and addressed:

**Use Case 1: Passenger Counting**
- The autonomous shuttle has a fixed capacity regarding the number of passengers it can carry.
- The video cameras installed in the autonomous shuttle acquire the color depth images and the data are fed into the system's video analytics algorithms for further analysis.
- In case the algorithms identify that the total number of passengers is reached, the shuttle stops receiving any others and appropriate notifications are sent to the AVENUE mobile app for the passengers that would like to board.

**Use Case 2: Route Optimization**
- Even though the shuttle is in full capacity, there may still be people waiting on a bus stop to go aboard.
- The bus only makes a stop when a passenger needs to get-off, while the route is modified to save time and cost.
- The number of onboard passengers is always being monitored, so that new passengers could get on, in case of availability.

**Use Case 3: Passenger Awareness**
- Even though the autonomous vehicle has reached its terminal, there could still be passengers onboard.
- The shuttle counts the number of passengers to make sure there is no one left.
- If there are passengers, the bus alerts them to get-off.

## 2.5.2 Stakeholders (development/prototyping team)

For the "Automated passenger presence" service the stakeholders are the same as in Stakeholders (development/prototyping team) (Section 2.4.2).

## 2.5.3 Technical requirements

For the "Automated passenger presence" service the technical requirements are the same as in Technical requirements (Service: Security trust services (C)), without the need for any audio sensing equipment (such as microphones).

### 2.5.3.1 Technology

Automatic passenger counting systems have evolved considerably within the past 40 years. Passenger flow data can be acquired with high accuracy outperforming manual ride checkers[38]. Devices that operate on 3D image streams are the industrial state-of-the-art technology. Latest generation devices offer an accuracy of around 99%[39] and technical progress is ongoing.

**Background**

So far, a wide range of competing automatic passenger counting technologies has been developed. Detection methods include infrared light beam cells, passive infrared detectors, infrared cameras, stereoscopic video cameras, laser scanners, ultrasonic detectors, microwave radars, piezoelectric mats, switching mats, and also electronic weighing equipment (EWE)[40].

Operators usually mount one or multiple sensors to collect data in each door area of public transport vehicles like buses, trams, and trains. The number of boarding and alighting passengers are counted separately by converting 3D video streams (infrared beam break) or light barrier methods, which are the most commonly used technologies.

In recent years also weight-based EWE approaches utilizing pressure measurements in the vehicle braking/air bag suspension system have emerged to estimate passenger numbers[41]. These relatively new approaches have proven to provide easy-to-acquire additional information since modern buses and powertrains are equipped with (intelligent) pressure sensors by default. However, those methods are considered not feasible in smaller scale vehicles, such as the autonomous buses in our case, due to design constraints and other restrictions.

**Passenger Detection**

As depicted in Figure 43, the first layer of sensors connects to the Hardware Abstraction Layer (HAL). The HAL implements the IP and the USB protocol supporting IP and USB cameras correspondingly. The input data are converted and transformed in a compatible format and passed into the analytics algorithms. The prediction is then transferred via the API endpoints into the cloud. Other systems and services have access to data.

---

[38] Hwang, M., Kemp, J., Lerner-Lam, E., Neuerburg, N., Okunieff, P.E.: Advanced public transportation systems: the state of the art update 2006. FTA report (2006)

[39] Hella Aglaia: Public Transport: HELLA Aglaia People Sensing (2018). http://people-sensing.com/public-transport.

[40] Kotz, A.J., Kittelson, D.B., Northrop, W.F.: Novel vehicle mass-based automated passenger counter for transit applications. Transp. Res. Record 2536, 37–43 (2015)

[41] Nielsen, B.F., Frølich, L., Nielsen, O.A., Filges, D.: Estimating passenger numbers in trains using existing weighing capabilities. Transportmetrica A Transp. Sci. 10(6), 502–517 (2014)

**Figure 43. High level overview of the "Automated passenger presence" service**

In order to detect timely and accurately human items, we implemented YOLO[42], a state-of-the-art and real-time object detection system. YOLO applies a single neural network to the full image, which divides the image into regions and predicts bounding boxes and probabilities for each region. These bounding boxes are weighted by the predicted probabilities.

YOLO model has several advantages over classifier-based systems. It looks at the whole image at test time so its predictions are informed by global context in the image. It also makes predictions with a single network evaluation unlike systems like R-CNN which require thousands for a single image. This makes it extremely fast, more than 1000x faster than R-CNN and 100x faster than Fast R-CNN.

### 2.5.3.1.1 Research results

The algorithm implemented monitors the shuttle in real time and records the position of each person in frames (Figure 45).



| Method | mAP-50 | time |
|---|---|---|
| [B] SSD321 | 45.4 | 61 |
| [C] DSSD321 | 46.1 | 85 |
| [D] R-FCN | 51.9 | 85 |
| [E] SSD513 | 50.4 | 125 |
| [F] DSSD513 | 53.3 | 156 |
| [G] FPN FRCN | **59.1** | 172 |
| RetinaNet-50-500 | 50.9 | 73 |
| RetinaNet-101-500 | 53.1 | 90 |
| RetinaNet-101-800 | 57.5 | 198 |
| **YOLOv3-320** | 51.5 | **22** |
| **YOLOv3-416** | 55.3 | 29 |
| **YOLOv3-608** | 57.9 | 51 |

**Figure 44. Performance comparison between various state-of-the art methods[43]**

---

[42] Redmon, Joseph, et al. "You only look once: Unified, real-time object detection." Proceedings of the IEEE conference on computer vision and pattern recognition. 2016.

[43] https://pjreddie.com/darknet/yolo/

**Figure 45. People counting using the on-board camera by NAVYA**

More results are also presented in Video analysis results sub-section of the "Enhance the sense of security and trust" service.

### 2.5.3.2    Equipment

The same equipment is utilised that is enlisted in Section 2.4.3.2 and 2.4.3 of the Security trust service without the need for any audio sensing equipment (such as microphones).

## 2.5.4    Prototyping plan

The prototyping plan of this service consists of the development phase and the deployment phase:

**Development phase:** In this phase, CERTH initially plans and analyses the requirements with inputs from all the stakeholders. Once the requirement analysis is done, the final software and hardware representation and documentation of the requirements are accepted from the project stakeholders. The next stage consists of designing the critical service components, the research and the development of the algorithms and the integration of the components across several platforms. Also, in this stage several data captures are performed for training the employed machine learning algorithms. CERTH will perform additional data capture sessions with different lighting conditions, angles and operators. The same data from the "Enhance the sense of security and trust" service's Prototyping plan can be reused. An initial setup of the equipment and the camera sensors is required, along with the interconnection of the components and software preparation. The captured data should demonstrate multiple camera angles and various capturing conditions. The next steps include the re-evaluation of the results, using the new data acquired for retraining the machine learning algorithms.

**Deployment phase**: In the deployment phase, the individual service components are unified into a complete system. The system has discrete input and output and is ready to be integrated with other

platforms. In the following months, CERTH will deploy the service in a demo AV, perform internal tests and verifications with all the stakeholders. In this stage, minor modifications and fine-tuning may be applied on the final setup, mainly on design or algorithms of the service depending on the real operation conditions.

**Prototyping phase:** Amobility will in coordination with CERTH install cameras and sensors in the Amobility shuttles for testing and validation of the technologies, use cases etc.

## 2.5.5    Result analysis

At the final step, an evaluation of the service under real conditions follows, that can be performed on daily or controlled routes of Amobility operator along with a safety driver inside the shuttle, depending on the GDPR permissions. Also, short-sessions for assessment with real passengers will be available. The results will also be compared with the manual counting done by Amoblity for validation purposes of the service's accuracy. After the successful evaluation of the service in Amoblity sites, the results will be used to fine tune the service if required and the service will be deployed and evaluated also in other sites of the involved operators in the AVENUE project towards a successful integration to the relevant pilot sites.

# 2.6 Service: Follow my kid/grandparents

## 2.6.1 Concept of service

The service "Follow my kid/grandparents" is designed to increase autonomy of non-fully autonomous people (Kids, Grandparent(s), disabled people etc.). It will allow carers or family members to be sure that their beloved family members are safe while moving around the city using public transports. On the other hand, it will increase confidence to the non-fully autonomous people to use public transports knowing that their family can "be with them". Surveillance using sensors such as cameras (cameras of different technologies can be used so that passengers' privacy is protected) and microphones, as well as smart software in the shuttle will maximize the feeling of security and the actual level of security.

Several concerns of the end users regarding the Safety and Robustness of the autonomous vehicles that are directly linked to the final User Acceptance of the new technology, can be identified. The prospective passengers fear several possible instances that could arise in case there is no driver in the bus. Indicatively:

- Passengers feeling discomfort travelling alone during nighttime
- Parents not being able to know if their kids have reached their destination safely
- Caregivers not being able to track passengers with dementia or other health issues

To address the aforementioned concerns on social and personal safety and security into the vehicle, certain measures need to be implemented. For example, third parties monitoring the route of minors or passengers with health issues could make their route much easier and less frightening. This may be followed by appropriate notifications and/or instructions to the third party, while the vehicle may also implement respective actions.

Moreover, implementing a solution for monitoring the routes of kids and patients will support safekeeping not only the users of the autonomous public shuttle but also the vehicle itself. In this section, the implementation of a video, depth and audio analytics software module for an embedded security subsystem or for cloud-based services of the system are described along with appropriate planning for the deployment and test of the service into the pilot sites of the AVENUE project.

### 2.6.1.1 Use case

The service and scenario proposes a full-fledged solution that allows designated "guardians" to follow the APT journeys of more vulnerable people, since the guardians can check the trip via a dashboard or mobile app, receive notifications via mobile app, add people to their "guarded" list, and share trips/position and E.T.A. with others.

In the context of AVENUE project, the following use cases have been identified to be further examined and addressed:

**Use Case 1:**
- Travelling without a guardian during nighttime can be unsettling for a vulnerable person.

- The video cameras installed in the autonomous shuttle acquire the color depth images and the data are fed into the system's video analytics algorithms for further analysis.
- When the vehicle's system identifies the passenger, the tracking begins.

**Use Case 2: Kids Monitoring**
- Parents need to be able to track their kids.
- The video cameras installed in the autonomous shuttle acquire the color depth images and the data are fed into the system's video analytics algorithms for further analysis.
- When the shuttle's system identifies the kid, the tracking begins, and the parents can monitor their route.

**Use Case 3: Patients Monitoring**
- Caregivers need to be able to track their patients, especially when they are not able to commute on their own.
- The video cameras installed in the autonomous shuttle acquire the color depth images and the data are fed into the system's video analytics algorithms for further analysis.
- When the autonomous bus's system identifies the patient, the tracking begins, and the caregivers can monitor their route.

## 2.6.2 Stakeholders (development/prototyping team)

For the "Follow my kid/grandparents" service the stakeholders are the same as in Stakeholders (development/prototyping team) (Section 2.4.2).

## 2.6.3 Technical requirements

In order to develop this service, we implement a facial recognition system capable of identifying or verifying a person from a video frame. There are multiple methods in which facial recognition systems work, but in general, they work by comparing selected facial features from given image with faces within a database. It is also described as a Biometric Artificial Intelligence based application that can uniquely identify a person by analyzing patterns based on the person's facial textures and shape. Many face recognition techniques require multiple data of the subject in the training dataset, in order to correctly identify the face of a person. From our perspective this is not possible as we rely in one single input image of the subject. To be able to overcome this problem, we are using one shot (or single shot) facial recognition algorithms.

### 2.6.3.1 Technology

In the following sections, the research, involving algorithms and experiments conducted, is presented for the "Follow my kid/grandparents" service. As depicted in Figure 52, the first layer of sensors connects to the Hardware Abstraction Layer (HAL). The HAL implements the IP and the USB protocol supporting IP and USB cameras respectively but also can request raw data by the API endpoints in order to perform face recognition. The input data is converted and transformed in a compatible format and passed into the analytics algorithms. The prediction is then transferred via the API endpoints into the cloud. The user has access to the data and acts accordingly.

**Figure 52. High level overview of the "Follow my kid/grandparents" service**

One-shot learning is a classification task where one, or a few, examples are used to classify many new examples in the future. This characterizes tasks seen in the field of face recognition, such as face identification and face verification, where people must be classified correctly with different facial expressions, lighting conditions, accessories, and hairstyles given one or a few template photos.

Modern face recognition systems approach the problem of one-shot learning via face recognition by learning a rich low-dimensional feature representation, called a face embedding, that can be calculated for faces easily and compared for verification and identification tasks. Historically, embeddings were learned for one-shot learning problems using a Siamese network. The training of Siamese networks with comparative loss functions resulted in better performance, later leading to the triplet loss function used in the FaceNet[44] system by Google that achieved state-of-the-art results on benchmark face recognition tasks.



**Figure 46. Architecture of a Siamese Neural Network**

Instead of directly classifying an input (test) image to one of the 10 people in the shuttle, this network instead takes an extra reference image of the person as input and will produce a similarity score denoting the chances that the two input images belong to the same person. Typically, the similarity score is squished between 0 and 1 using a sigmoid function; wherein 0 denotes no similarity and 1 denotes full similarity. Any number between 0 and 1 is interpreted accordingly. Notice that this network

---

[44] Schroff, Florian, Dmitry Kalenichenko, and James Philbin. "Facenet: A unified embedding for face recognition and clustering." Proceedings of the IEEE conference on computer vision and pattern recognition. 2015.

is not learning to classify an image directly to any of the output classes. Rather, it is learning a similarity function, which takes two images as input and expresses how similar they are.

A new passenger can be enrolled to the Follow My Kid service using a single image of his face which will be stored in a database. Using this as the reference image, the network will calculate the similarity for any new instance presented to it. Thus, we conclude that the network can predict the score in one shot.

### 2.6.3.2 Equipment

The same equipment is utilised that is enlisted in Section 2.4.3.2 and 2.4.3 of the Security trust service without the need for any audio sensing equipment (such as microphones).

## 2.6.4 Prototyping plan

The prototyping plan of this service consists of the development phase and the deployment phase:

**Development phase:** In this phase, CERTH initially plans and analyses the requirements with inputs from all the stakeholders. Once the requirement analysis is done, the final software and hardware representation and documentation of the requirements are accepted from the project stakeholders. The next stage consists of designing the critical service components, the research and the development of the algorithms and the integration of the components across several platforms. Also, in this stage several data captures are performed for training the employed machine learning algorithms. CERTH will perform additional data capture sessions with different lighting conditions, angles and operators. The same data from the "Enhance the sense of security and trust" service's Prototyping plan can be reused. An initial setup of the equipment and the camera sensors is required, along with the interconnection of the components and software preparation. The captured data should demonstrate multiple camera angles and various capturing conditions. The next steps include the re-evaluation of the results, using the new data acquired for retraining the machine learning algorithms.

**Deployment phase**: In the deployment phase, the individual service components are unified into a complete system. The system has discrete input and output and is ready to be integrated with the other platforms. In the next months, CERTH will deploy the service in a demo AV, perform internal tests and verifications with all the stakeholders. In this stage, minor modifications and fine-tuning may be applied on the final setup, mainly on design or algorithms of the service depending on the real operation conditions.

**Prototyping phase:** Amobility will in coordination with CERTH install cameras and sensors in the Amobility shuttles for testing and validation of the technologies, use cases etc.

## 2.6.5 Result analysis

At the final step, an evaluation of the service under real conditions follows, that can be performed on daily or controlled routes of Amobility operator along with a safety driver inside the shuttle, depending on the GDPR permissions. Also, short-sessions for assessment with real passengers will be available. The results will also be cross validated with manual person identification by Amoblity for validation purposes of the service's accuracy and compared to the mobile application (by MT) if possible. After the successful evaluation of the service in Amobility sites, the results will be used to fine tune the service if required and the service will be deployed and evaluated also in other sites of the involved operators in the AVENUE project towards a successful integration to the relevant pilot sites.

## 2.7 Service: Shuttle environment assessment

### 2.7.1     Concept of service

The service "Shuttle environment assessment" aims to maintain at acceptable levels the environmental conditions in the autonomous vehicle that may not be adequately controlled due to the absence of the shuttle driver. Minimum acceptable conditions and comfort, such as good air quality, acceptable odours and absence of smoke are necessary for the safe transport of the passengers, as well as the viability of the whole autonomous service, since lack of these conditions within the vehicle could significantly discourage potential passengers. After all, monitoring the environment conditions could enable for passengers' alert and warning services via notifications, thus enhancing the user experience and safety during their trips.

Under these circumstances, there are several instances that have to be considered in order for the prospective passengers to feel content and safe.  In particular, while there would be no driver inside the vehicle, variable problems might come up, such as the following:

● There will be no presence of the stuff inside the bus to prevent someone from lighting a cigarette
● In quite high or low temperatures, there will be no driver in order to regulate the air conditioning system
● In emergency situations there will be no one in charge of informing the operators and the competent authorities
● If the air concentration of CO2 inside the bus is high, and someone might get dizzy or exhibit breathing difficulties, there will be no driver so as to either open the windows or stop the shuttle

As far as it is considered, the buses should pose a comfortable environment for all the passengers. This feeling could undoubtedly be strengthened by controlling the temperature inside the vehicle. Besides, heating, ventilation and air conditioning control, now belong to the standard equipment on city buses. As a result, it is crucial for temperature sensors to be positioned in specific locations inside the vehicle. Simultaneously, these sensors will be connected with the air conditioning system and in cases that the temperature will exceed a suitable limit, the air condition will be put into operation. In that way, the existing temperature will be autonomously adjusted, in order to provide the appropriate indoor climate, namely not to be neither too hot nor too cold for the passengers.

In addition, detection of certain pollutants, such as CO2, NO2, or dust particles in the indoor environment, along with critical temperature variations, is important for the condition of certain passengers, especially ill people, such as asthma patients. Smoke in the vehicle, i.e. from a person that lights a cigarette, will deteriorate the passenger experience but may also put in danger the whole vehicle (danger of fire). Detection of certain events (air quality deterioration, smoke) may raise a notification or an alert to the passengers along with instructions on how to handle this situation, while the vehicle may also implement respective actions. Likewise, smoke in the vehicle will deteriorate the passenger experience, put in danger the whole vehicle (danger of fire) and may also result in cancelling the autonomous transport service. Detection of certain events (air quality deterioration, smoke) may raise a

notification or an alert to the supervisor and/or the suitable authorities (i.e. police, fire department). This may be followed by appropriate notifications and/or instructions to the passengers, while the vehicle may also implement respective actions.

### 2.7.1.1 Use case

The service is responsible for the timely, accurate, robust and automatic detection of any change in the air quality and the presence of smoke or fire, inside the vehicle. In cases there will be an alert on the system, notifications and instructions will be sent to the passengers to the operators and/or to the suitable authorities.

In light of this, several possible situations could take place in, counting from high level $CO_2$ and $NO_2$ concentrations at the indoors air composition to presence of humidity, smoke or even fire. Especially, exposure to carbon dioxide can produce a variety of health effects. These may include headaches, dizziness, restlessness, a tingling or pins or needles feeling, difficulty breathing, sweating, tiredness, and increased heart rate. Furthermore, detection of certain air quality indexes and pollutants in the indoor environment, along with critical temperature variations are necessary for providing a secure service to the passengers. In particular, a gas composition sensor could be used for checking the inside air quality. This sensor monitors the air intake of the heating, ventilation, and air conditioning system of the vehicle, while it detects undesirable gases and adjusts the system accordingly, by shutting off the intake and recirculating the indoor air back to the outside. Furthermore, humidity and temperature sensors could be used for measuring/regulating indoor air quality and adjusting the heating, ventilation, and air conditioning system settings accordingly.

Despite that, another possible scenario, mostly observed during the rainy days of the year, is the fogging of the windows due to the increased humidity. This might have a negative impact to the passengers' attitude, while reinforcing feeling of confinement. Consequently, including a fogging prevention sensor inside the vehicle might be an efficient solution. More specifically, fogging prevention sensors are used to prevent fogging of the windshield glass. These sensors consist of three sensing elements for sensing indoor temperature, windshield glass temperature, and cabin humidity. The fogging sensor feedback is used for adjusting the heating, ventilation, and air conditioning system to maintain the interior temperature higher than the windshield glass temperature. Hence, it prevents the windshield fogging up.

Summarizing, it is passengers' wish to travel in a clean and comfortable environment and be notified if the conditions have deteriorated. Moreover, operators would like to be notified when environmental conditions are considered harmful for the passengers. For a deeply explanation, some of the most representative use cases are indicatively displayed as following:

**Use Case 1: Lighting a Cigarette**
- Inside the autonomous vehicle a passenger lights a cigarette.
- The smoke detection sensors embedded inside the vehicle detects the smoke coming out of the cigarette
- The real-time sensor data is sent to a central PC, installed in the vehicle and in which the data processing take place
- With the real-time data processing the PC decides that it is an emergency case and sents the message of smoke detection to the operators

● The suitable operator evaluates the criticality of the situation and decides how to intervene (with an announcement from the loudspeakers or by taking more drastically measurements, for example by stopping the bus).

**Use Case 2: Exposure to Carbon Dioxide**
● While the bus is on its route, high levels of Carbone Dioxide are detected from the relevant sensor.
● The sensor data are sent in real time to the central PC installed in the vehicle and are processed.
● Carbone Dioxide in unusual levels of air concentrations might have an adversely effect on passengers' health. For instance, high level of CO2 is considered to be related with dizziness, restlessness or breathing difficulties and increased heart rate.
● In order to prevent an event concerning these health issues to take place in, such as a passenger's fainting, the PC is sending to the vehicle's central system the command to open the windows, so as the air come back to its normal composition. Simultaneously, the passengers are informed for the air composition through the corresponding mobile application in real-time.

**Use Case 3: High Temperature on the Autonomous Vehicle**
● It is a hot day in summer and indoors the temperature is starting raising.
● The suitable sensor measures the temperature and sends the data to the central PC installed in the vehicle.
● When the temperature exceeds a predefined level, the PC sends to the vehicle's central system the command to put the cooling system into operation.
● At the same time, passengers are able to be informed for the temperature inside the vehicle through the mobile application.

## 2.7.2    Stakeholders (development/prototyping team)

In this section, the relevant stakeholders involved in the development and prototyping of the Shuttle Environment Assessment Service are introduced. In particular:

● **CERTH** is conducting innovative research regarding the development of the environmental assessment in terms of setting up the sensors inside the vehicle. Simultaneously, CERTH will determine the sensors that will be needed, the specific role that every single sensor has to play and how this equipment will be efficiently organized and combined so as the operators will have access to the crucial information collected from sensors.
● **Bestmile** provides multiple integration interfaces towards the different stakeholders in the project (e.g. vehicle manufacturers, public transport operators) into its cloud platform. In this service, Bestmile provides connectivity between the sensors' system and the related operators regarding the notifications that are generated by the detection software.
● **MobileThinking** is responsible for the AVENUE mobile application development. This application will illustrate the environment-related measurements or some other relevant with the passenger's convenience information, in a user-friendly environment, while it will provide any necessary notifications to the passengers.
● **Amobility** is the operator of the autonomous vehicles in the Copenhagen nad Oslo site and is handling all daily operation of the vehicles and everyday contact with the end users. As
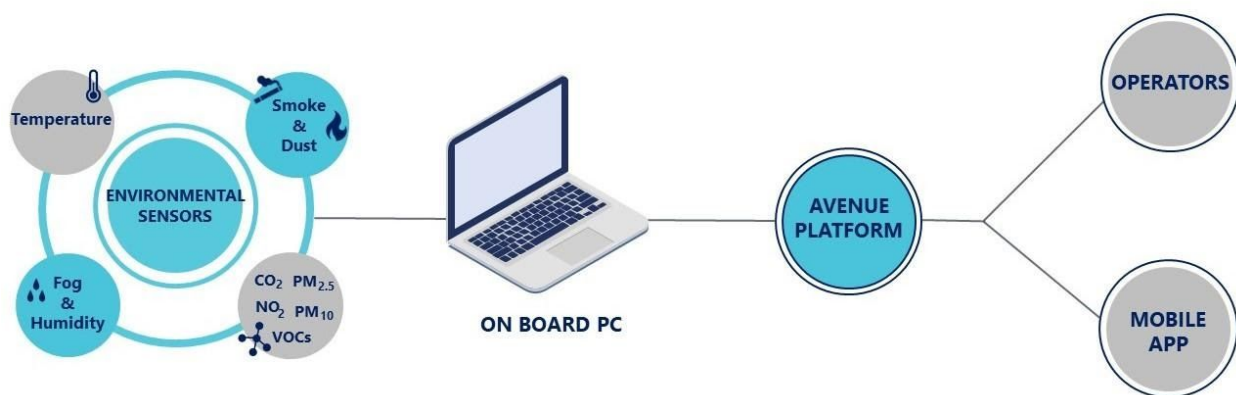
an operator Amobility provides the autonomous shuttles that will be used for the AVENUE pilot activities and its facilities for performing data capture activities, deploying the service and testing its performance through short or longer evaluation periods.

## 2.7.3 Technical requirements

The service is based on the collection of several vehicle telemetry information from sensors and related services. Accordingly, several sensors are required, such as:

- Environmental sensors (for assessing pollution, air quality in the vehicle, such as CO2, NO2, dust particles concentrations, temperature in the vehicle, humidity sensors, fogging prevention sensors etc.)
- Smoke sensors and smoking detection (for fire and also people that smoke in the vehicle)



**Figure 42:Environmental Assessment Process**

Obviously, a wide variety of sensors should be embedded into the vehicle for this purpose, along with proper connectivity with a central assessment system. Subsequently, the conditions inside the bus will be monitored in order to ensure stability. In more detail, collection of the sensor data, data processing and subsequently, data transmission to a PC and furthermore to operators seems to be essential so as for the condition stability inside the vehicle to be guaranteed. Likewise, it is necessary, an interface with the central AVENUE platform to be developed, as well as, a mobile application so that passengers will be able to have access to information regarding the status of the vehicle or other relevant information. As a result of such implementation, when it is required, notifications to the passengers and supervisors will be provided and in case of alarms, the supervisors or other suitable authorities, will be notified in order to deal with the situation or prevent an unfortunate event.

### 2.7.3.1 Technology

The technology utilised is based on the different sensors that will be used to develop and deploy this service. A detailed explanation of the technology of each sensor is out of the scope of this deliverable. Indicative sensors are presented in the next section.

## 2.7.3.2     Equipment

Overall, multiple environmental sensors will be needed, so as for the quality of the air to be determined, as well as sensors about smoke detection or fogging prevention. Moreover, a third PC is tented to be installed inside the vehicle, in order to collect and process the data of the sensors. Finally, there must be appropriate action to assure the power supply to the sensors.

In particular, regarding the sensors that will be used, some indicative models are proposed as it follows:

**Aeroqual**

- Series 500 – Portable Indoor Air Quality Monitor: It is compatible with 30 different sensors. It has the ability of real-time data observation through a screen, while all the metrics are stored in memory and can be transferred to PC via a USB. [45]

- Indoor Air Quality Test Kit (Starter): It measures PM2.5, PM10, CO2, VOC, temperature and humidity, as it uses the suitable sensor. Also, it can connect to a PC through a USB and the measurements are captured in a screen. [46]

**TSI**

- AIRASSURE PM2.5-AD INDOOR AIR QUALITY MONITOR IPM2.5-AD:  It can measure PM2.5. It is placed onto the wall and it includes a screen, in which the measurement is illustrated [47]
- Q-TRAK INDOOR AIR QUALITY MONITOR 7575: It has the ability of measuring CO2, CO, temperature and humidity. On its screen it could be able to capture up to 5 measurements, while VOC sensors could be connected. Furthermore, it can store up to 39-days data, with samples of 1-minuite frequent. Communication through Bluetooth is also available. [48]

**IOTSENS**

- IOTSENS city: It can measure PM1, PM2.5, PM10, temperature and humidity. It is a complete solution, as it has the ability to connect through LoRaWAN, Sigfox, NBIoT, GPRS, Wifi, while there are specific platforms and applications from the company. [49]

**Senseair**

---

[45] Aeroqual, Series 500 – Portable Indoor Air Quality Monitor, Available online at:
https://www.aeroqual.com/product/series-500-portable-indoor-monitor
[46] Aeroqual, Indoor Air Quality Test Kit (Starter), Available online at:
https://www.aeroqual.com/product/indoor-portable-monitor-pro-kit
[47]Aeroqual, AIRASSURE PM2.5-AD INDOOR AIR QUALITY MONITOR IPM2.5-AD, Available online at:
            https://tsi.com/products/indoor-air-quality-meters-instruments/installed-iaq/airassure-pm2-5-ad-indoor-air-quality-monitor-ipm2-5-ad/
[48]Aeroqual, Q-TRAK INDOOR AIR QUALITY MONITOR 7575, Available online at:

https://tsi.com/products/indoor-air-quality-meters-instruments/indoor-air-quality-meters/q-trak-indoor-air-quality-monitor-7575/
[49] IOTSENS, IOTSENS city, Available online at:   http://www.iotsens.com/solution/smart-city/

- EXPLORAPM2.5: It can measure PM2.5, temperature and humidity. It connects through LoRaWAN and there is a specific cloud-based platform and application.[50]

- EXPLORACO2: It is able to measure CO2, temperature and humidity. It connects through LoRaWAN and there is a specific cloud-based platform and application. In addition, it has the ability to connect with an application, developed from externals, using the open API.[51]

- SENSEAIR AERCAST: It is able to measure CO2, temperature and humidity. It connects through BLE as well as other choices such as LoRa, ZigBee, $\kappa\,\alpha\,\iota$ N B I o T.[52]

**Awair**

- Awair Element: It is able to measure CO2, VOCs, PM2.5, temperature and humidity. It connects through Bluetooth and includes an application for the control of the measurements. Additionally, a LED is available to indicate dimensions.[53]

- Awair Glow C:  It is able to measure VOCs, temperature and humidity.  It is connected to the socket and devices such as a fan, dehumidifier, heater can be connected to it, which are activated when the limit we have set for a measurement is exceeded. It has an application and has the ability to connect to Google Home.[54]

- Awair 2nd Edition: It is able to measure VOCs, CO2, PM2.5, temperature and humidity. It has an application for controlling the measurements and can be connected to Google Home, Alexa and Ecobee. It has a Led to indicate the measurements.[55]

- Awair Omni: It has the ability to measure VOCs, CO2, PM2.5, ambient light, ambient noise, temperature and humidity. It has a Led to indicate the measurements. It can connect to WiFi, Bluetooth, Cellular, Ethernet. It is mounted on the wall.[56]

**Airthings**

---

[50] Senseair, EXPLORAPM2.5, Available online at:
 https://senseair.com/products/aercast-explora-family/explorapm2-5/

[51] Senseair, EXPLORACO2, Available online at:
 https://senseair.com/products/aercast-explora-family/exploraco2/

[52] Senseair, SENSEAIR AERCAST, Available online at:
 https://senseair.com/products/aercast-explora-family/aercast/

[53]  Awair, Awair Element, Available online at:
 https://getawair.com/pages/awair-element

[54] Awair, Awair Glow C, Available online at:
 https://getawair.com/pages/awair-glow

[55] Awair, Awair 2nd Edition, Available online at:
 https://getawair.com/pages/awair-2nd-edition

[56] Awair, Awair Element, Available online at:
 https://getawair.com/pages/awair-for-business

- Wave Plus: It has the ability to measure radon, CO2, VOCs, pressure, temperature and humidity. It connects via Bluetooth and has its own application and platform for displaying and processing data.[57]

- Wave Mini: It has the ability to measure VOCs, temperature and humidity. It has Led for visual display and is connected via Bluetooth. It can be used in conjunction with Google Assistant and has its own application.[58]

**Netatmo**

- Netatmo: It has the ability to measure air quality, noise, temperature and humidity. Moreover, it is able to connect via Bluetooth and has its own application.[59]

**Kaiterra**

- Laser Egg: It has the ability to measure PM2.5, temperature and humidity. It has WiFi connectivity, its own application and can connect to other smart devices via Apple Home Kit. It has a screen on which the measurements are illustrated, as well as the weather forecast.[60]

- Laser Egg + Chemical: It has the ability to measure PM2.5, VOCs, temperature and humidity. It has WiFi connectivity, has its own application and can connect to other smart devices via Apple HomeKit. It has a screen on which the measurements are illustrated, as well as the weather forecast.[61]

- Laser Egg + CO2: It has the ability to measure PM2.5, CO2, temperature and humidity. It has WiFi connectivity, has its own application and can connect to other smart devices via Apple HomeKit. It has a screen on which the measurements are illustrated, as well as the weather forecast.[62]

## 2.7.4    Prototyping plan

The prototyping plan of this service consists of the development phase and the deployment phase:

**Development phase:** In this phase, CERTH initially plans and analyses the requirements with inputs from all the stakeholders. Once the requirement analysis is done, the final software and hardware representation and documentation of the requirements are accepted from the project stakeholders. The

---

[57] Airthings, Wave Plus, Available online at:
https://www.airthings.com/wave-plus

[58] Airthings, Wave Mini, Available online at:
https://www.airthings.com/wave-mini

[59] Netatmo, Available online at:
https://www.netatmo.com/en-us/aircare/homecoach

[60] Kaiterra, Laser Egg, Available online at:
https://www.kaiterra.com/en/laser-egg/

[61]Kaiterra, Laser Egg + Chemical, Available online at:
europe.kaiterra.com/collections/frontpage/products/laser-egg-plus-chemical

[62] Kaiterra, Laser Egg + CO2, Available online at:
https://www.kaiterra.com/en/laser-egg-co2

next stage consists of designing the critical service components, the research and the development of the algorithms and the integration of the components across several platforms. At this point, all the individual sensors that will be used, have to be specifically determined and either to be built or bought. Following, they will be connected all together at the AVENUE Platform.

**Deployment phase**: In the deployment phase, the individual service components are unified into a complete system. The system has discrete input and output and is ready to be integrated with the other platforms. In the next months, CERTH will deploy the sensors and the service in a demo AV, perform internal tests and verifications with all the stakeholders. Subsequently, they will be connected with the PC installed into the vehicle, as well as with the AVENUE platform and the mobile application. Then, for a sufficient period, the installed sensors will measure the corresponding metrics such as $CO_2$, $NO_2$ humidity, temperature, fog, dust, smoke, etc. After data processing, and through a calibration process, it will be possible for the conclusions we are interested in to be extracted. Afterwards, based on the obtained information, the rates of successful measurements will be able to be determined. In this stage, minor modifications and fine-tuning may be applied on the final setup, mainly on design or algorithms of the service depending on the real operation conditions.

**Prototyping phase:** Amobility will in coordination with CERTH install cameras and sensors in the Amobility shuttles for testing and validation of the technologies, use cases etc.

## 2.7.5    Result analysis

At the final step, the evaluation of the whole system, from the sensor measurement process to the appropriate messages and measurements appeared onto the users' mobile application, will take place under real conditions, which can be performed on daily or controlled routes of Amoblity operator along with a safety driver inside the shuttle. Also, short-sessions for assessment with real passengers will be available. The results will also be cross validated with manual person notification by Amoblity for validation purposes of the service's accuracy.

The measurements from the sensors will be processed and maybe used in order to predict the condition indoor regarding the air quality, for the next couple of hours. These measurements, as well as the conclusions of the real-time assessment, will be sent both to passengers and operators. In cases of alarms, such as a possible fire, they will be sent to the suitable authorities. Moreover, this may be followed by appropriate notifications and/or instructions to the passengers, while the vehicle may also implement respective actions.

After the successful evaluation of the service in Amobility sites, the results will be used to fine tune the service if required and the service will be deployed and evaluated also in other sites of the involved operators in the AVENUE project towards a successful integration to the relevant pilot sites.

Not approved yet

# 2.8 Service: Smart feedback system

## 2.8.1 Concept of service

The service "Smart feedback system" aims at allowing the travellers inside the shuttle to give easy and effortless feedback to the operators of the shuttle - when the safety driver is no longer inside the shuttle. It is important for the operators to know if people are satisfied with the services and transportation. Currently the safety driver talks to the travellers and via his/hers presence also becomes the conversation channel between the operators and the travellers. When the safety driver is removed, knowing whether they are satisfied or disappointed can be even more important, as the safety driver is not there to support, hence knowing how to assist the travellers.
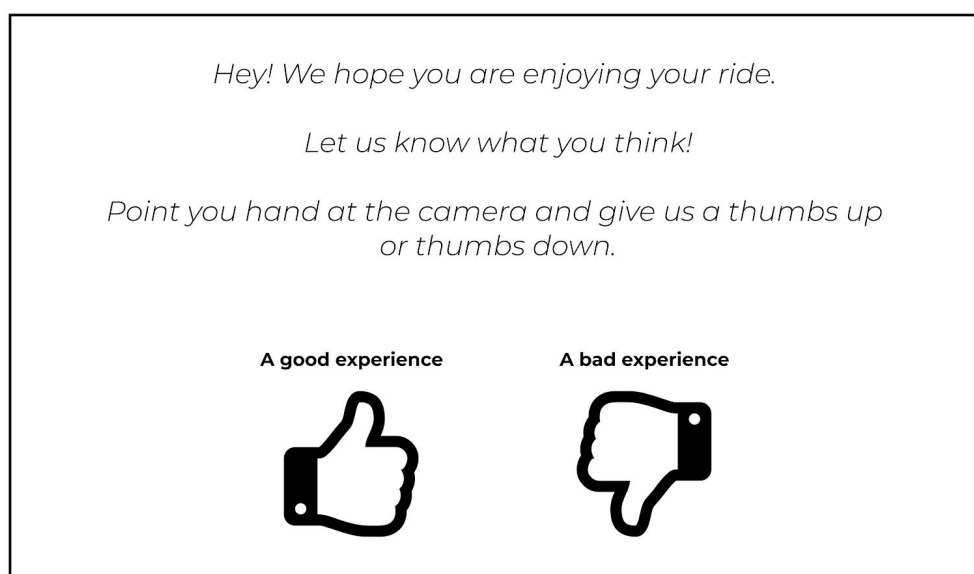
When removing the operator, automated services have to perform the same level of service and interaction as he/she did while being in the shuttle.

This service aims at allowing travellers to give their feedback about liking/disliking the service experience as easy as possible. This will be done by instructing the travellers to give a hand gesture to one of the cameras inside the shuttle. This will allow the travellers to effortlessly say "i like" or "i dont like" the experience with a thumbs up or a thumbs down.

The concept will be communicated to the travellers via stickers inside the shuttle. To begin with camera technology will be used to capture the thumbs up or thumbs down, but if possible sounds sensors will also be tested to capture the experience/feedback from the travellers.

### 2.8.1.1 Use case

The "smart feedback service" will allow travellers to give their feedback by a hand gesture to the cameras inside the shuttle. The service will be communicated to the travellers as follows:



*Hey! We hope you are enjoying your ride.*

*Let us know what you think!*

*Point you hand at the camera and give us a thumbs up or thumbs down.*

**A good experience**     **A bad experience**

81

**Use case 1: Giving a thumbs up/down in light settings**
- Mid day with sunlight
- Good visibility for the cameras

**Use case 2: Giving a thumbs up/down in dark settings**
- Early morning or night with no sunlight
- Low visibility for the cameras

**Use case 3: Giving a thumbs up/down in crowded settings**
- Many passengers inside the shuttle, both standing and seating
- Low visibility for cameras due to people standing close to the cameras

**Use case 4: Giving a thumbs up/down in empty settings (or few passengers)**
- Little or no passengers inside the shuttle
- Good visibility for the cameras, easy to see the hand gesture

# 2.8.2 Stakeholders (development/prototyping team)

For the "Smart feedback system" service the stakeholders are the same as in Stakeholders (development/prototyping team) (Section 2.4.2).

# 2.8.3 Technical requirements

Gesture recognition refers to the whole process of tracking human gestures to their representation and conversion to semantically meaningful commands. Research in hand gesture recognition aims to design and development of such systems than can identify explicit human gestures as input and process these gesture representations. In this service, we implement a hand recognition system capable of detecting hand gestures from a video frame. Although there are multiple methods for hand gesture recognition, we focus on simplicity and high performance using Single Shot Detectors and shallow CNN models for classification.

## 2.8.3.1 Technology

In the following sections, the research, involving algorithms and experiments conducted, is presented for the "Smart feedback system" service. As depicted in Figure 50, the first layer of sensors connects to the Hardware Abstraction Layer (HAL). The HAL implements the USB protocol supporting USB. The input data is converted and transformed in a compatible format and passed into the analytics algorithms. The output prediction is then transferred via the API endpoints into the cloud, available for further processing from other platforms.
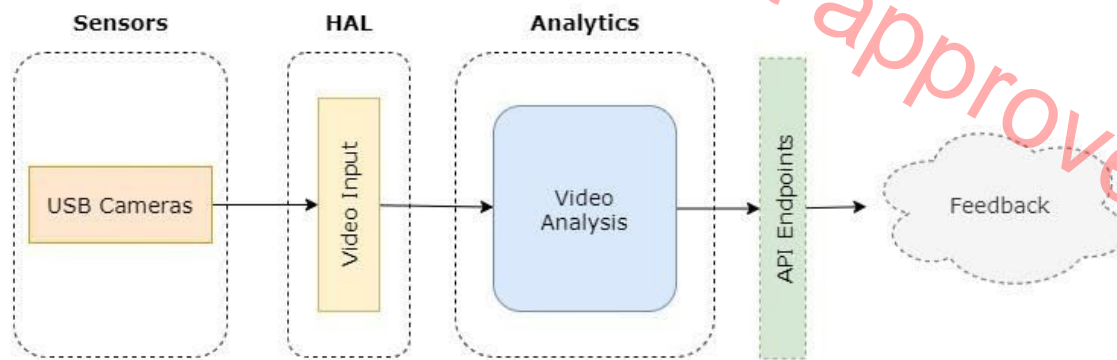
**Figure 50. High level overview of the "Smart feedback" service**

In the first stage, we use the Single Shot Detector (SSD) for the detection of a bounding box of where the hand(s) is and the corresponding cropped frame.  The SSD approach is based on a feed-forward convolutional network that produces a fixed-size collection of bounding boxes and scores for the presence of object class instances in those boxes, followed by a non-maximum suppression step to produce the final detections. The early network layers are based on a standard architecture used for high quality image classification (truncated before any classification layers), which are called as the base network2. Auxiliary structure is added to the network to produce detections with the following key features.

**Multi-scale feature maps for detection:** Convolutional feature layers to the end of the truncated base network. These layers decrease in size progressively and allow predictions of detections at multiple scales. The convolutional model for predicting detections is different for each feature layer.

**Convolutional predictors for detection:** Each added feature layer (or optionally an existing feature layer from the base network) can produce a fixed set of detection predictions using a set of convolutional filters. These are indicated on top of the SSD network architecture in Figure 47. For a feature layer of size $m \times n$ with p channels, the basic element for predicting parameters of a potential detection is a $3 \times 3 \times p$ small kernel that produces either a score for a category, or a shape offset relative to the default box coordinates. At each of the $m \times n$ locations where the kernel is applied, it produces an output value. The bounding box offset output values are measured relative to a default box position relative to each feature map location. Default boxes and aspect ratios. A set of default bounding boxes is associated with each feature map cell, for multiple feature maps at the top of the network. The default boxes tile the feature map in a convolutional manner, so that the position of each box relative to its corresponding cell is fixed. At each feature map cell, the offsets relative to the default box shapes in the cell are predicted, as well as the per-class scores that indicate the presence of a class instance in each of those boxes. Specifically, for each box out of k at a given location, c class scores and the 4 offsets relative to the original default box shape are computed.
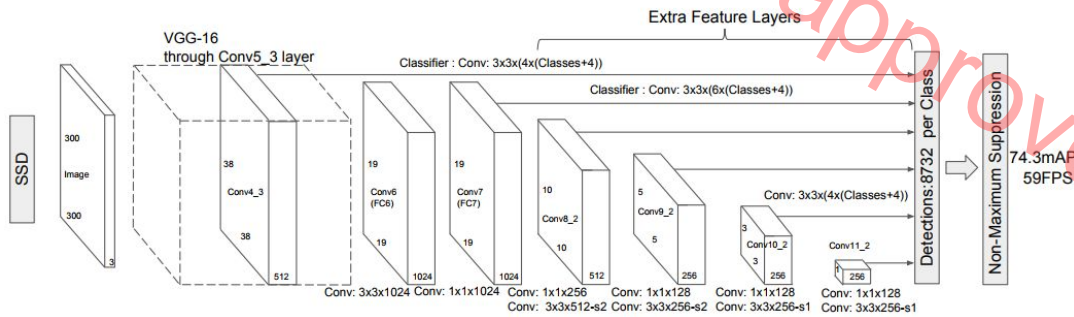
**Figure 47. SSD architecture using VGG16 for feature extraction**

This cropped frame of the hand is then passed to the CNN (Figure 48), that predicts a class vector output of values between 0 and 1. The values correspond to the probability of the frame to be one of the classes.
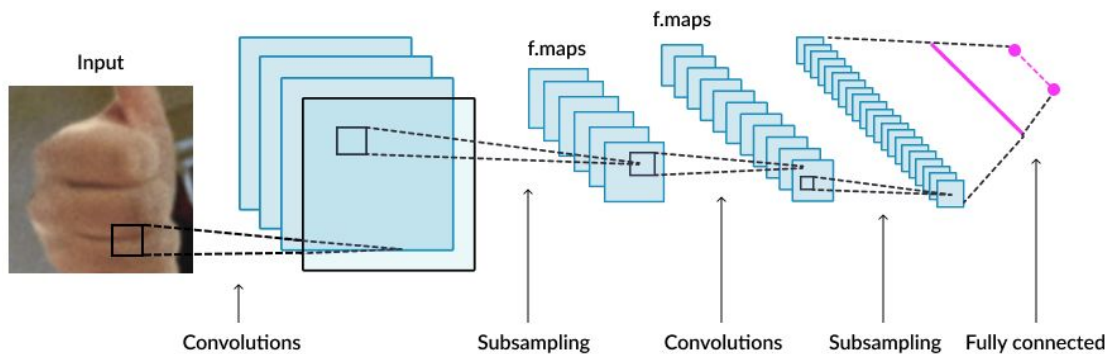


**Figure 48. CNN architecture**

### 2.8.3.1.1 Research results

The model is trained end-to-end and regularized so that it distils the most compact profile of the normal patterns of training data (Figure 49) and effectively detects the "thumbs-up and thumbs-down" gestures (Figure 50). The system will be extended and more gestures will be supported in the future.
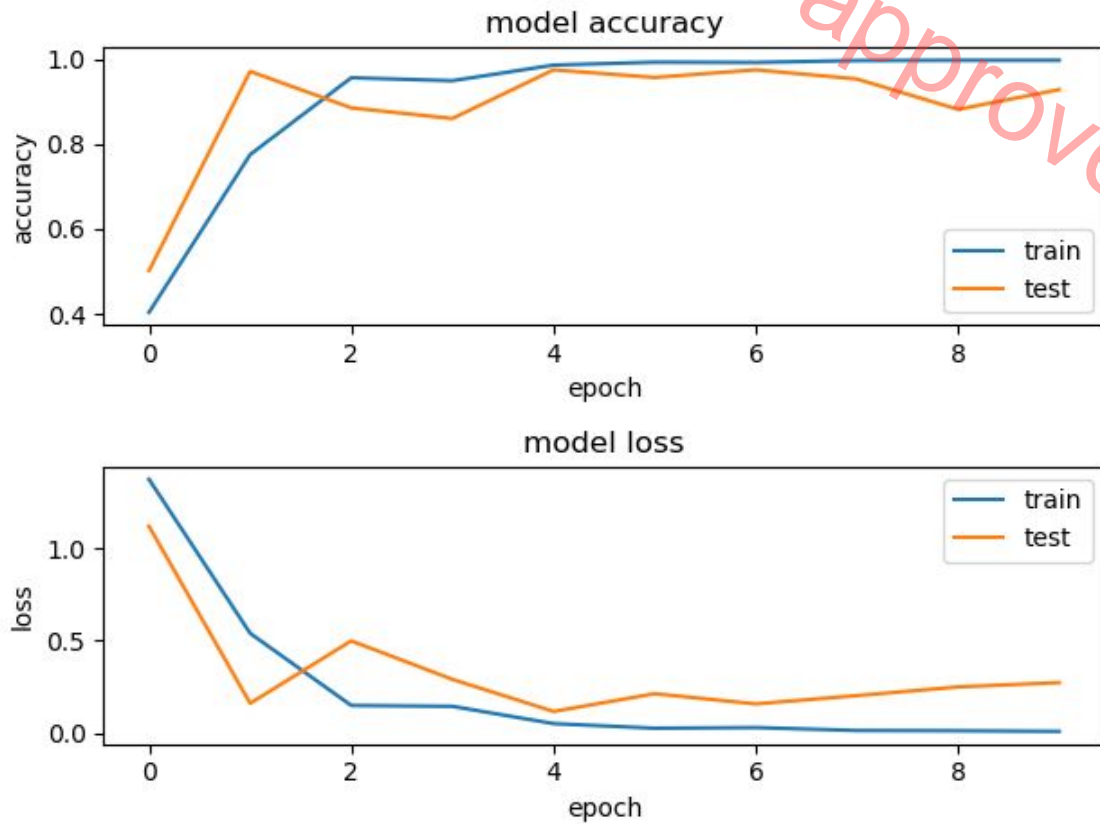
**Figure 49. Training accuracy and loss at 10 epochs**



**Figure 51. Inference results: (a) Detection of "thumbs-up" gesture. (b) Detection of "thumbs-down" gesture**

### 2.8.3.2    Equipment

The same equipment is utilised that is enlisted in Section 2.4.3.2 and 2.4.3 of the Security trust service without the need for any audio sensing equipment (such as microphones).

## 2.8.4     Prototyping plan

The prototyping plan of this service consists of the development phase and the deployment phase:

**Development phase:** In this phase, CERTH initially plans and analyses the requirements with inputs from all the stakeholders. Once the requirement analysis is done, the final software and hardware representation and documentation of the requirements are accepted from the project stakeholders. The next stage consists of designing the critical service components, the research and the development of the algorithms and the integration of the components across several platforms.

**Deployment phase**: In the deployment phase, the individual service components are unified into a complete system. The system has discrete input and output and is ready to be integrated with the other platforms. In the next months, CERTH will deploy the service in a demo AV, perform internal tests and verifications with all the stakeholders. In this stage, minor modifications and fine-tuning may be applied on the final setup, mainly on design or algorithms of the service depending on the real operation conditions.

**Prototyping phase:** Amobility will in coordination with CERTH install cameras and sensors in the Amobility shuttles for testing and validation of the technologies, use cases etc.

## 2.8.5     Result analysis

At the final step, an evaluation of the service under real conditions follows, that can be performed on daily or controlled routes of Amobility operator along with a safety driver inside the shuttle, depending on the GDPR permissions. Also, short-sessions for assessment with real passengers will be available. The results will also be cross validated with manual person feedback counting by the Amoblity operator. After the successful evaluation of the service in Amoblity sites, the results will be used to fine tune the service if required and the service will be deployed and evaluated also in other sites of the involved operators in the AVENUE project towards a successful integration to the relevant pilot sites.

# 3 Conclusions

Five in-vehicle services have been chosen to prototype and test in the coming 8-12 months in a collaboration between CERTH, Amoblity, Bestmile and MobileThinking. Each service targets different essential functions that the safety driver currently provides inside the shuttle. Once we he/she is removed, automated in-vehicle services have to replace the role of the safety driver. All the services are centered around camera and sensor systems inside the shuttle. The prototyping and testing will not lead to fully developed services, but feed the AVENUE project with recommendations that have to be taken into account in the future when the technology becomes more mature and when autonomous vehicles represent a large part of the transport means in public as well as private transport systems.

The prototyping and testing of the described in-vehicle services will start as soon as COVID-19 is under control, meaning that operations have started again and when flights can be taken to the site to integrate and prepare for the prototyping sessions.

Technical preparations that can be conducted without being present on the site, will start in M26, meaning digital architecture setup, GDPR preparations and so forth.

Once COVID-19 allows for resuming operations, technological implementations like installing ekstra sensors and cameras will be prepared, so that once flights are available CERTH and Amoblity can visit the sites together and install the necessary technology, and test the setup with a third computer etc.

# Appendix A: