POLITECNICO DI TORINO Repository ISTITUZIONALE

Dual-View Single-Shot Multibox Detector at Urban Intersections: Settings and Performance Evaluation

Original

Dual-View Single-Shot Multibox Detector at Urban Intersections: Settings and Performance Evaluation / Lenatti, Marta; Narteni, Sara; Paglialonga, Alessia; Rampa, Vittorio; Mongelli, Maurizio. - In: SENSORS. - ISSN 1424-8220. - 23:6(2023). [10.3390/s23063195]

Availability: This version is available at: 11583/2977584 since: 2023-05-31T14:35:13Z

Publisher: MDPI

Published DOI:10.3390/s23063195

Terms of use:

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

Publisher copyright

(Article begins on next page)



16

Dual-view Single Shot Multibox Detector at Urban Intersections: Settings and Performance Evaluation

Marta Lenatti^{1,‡}, Sara Narteni^{1,2,‡}, Alessia Paglialonga¹, Vittorio Rampa¹, Maurizio Mongelli^{1,*}

- ¹ National Research Council of Italy (CNR), Institute of Electronics and Information and Telecommunications Engineering (IEIIT);
- ² Politecnico di Torino Department of Control and Computer Engineering (DAUIN);
- * Correspondence: maurizio.mongelli@ieiit.cnr.it
- ‡ These authors contributed equally to this work.

Abstract: The explosion of Artificial Intelligence methods has paved the way to more sophisticated 1 smart mobility solutions. In this work we present a multi-camera Video Content Analysis (VCA) 2 system that exploits a Single Shot multibox Detector (SSD) network to detect vehicles, riders, and 3 pedestrians and triggers alerts to drivers of public transportation vehicles approaching the surveilled 4 area. The evaluation of the VCA system will address both detection and alert generation perfor-5 mances, by combining visual and quantitative approaches. Starting from a SSD model trained for a single camera, we added a second one, under a different field of view (FOV), to improve accuracy 7 and reliability of the system. Due to real-time constraints, the complexity of the VCA system must be limited, thus calling for a simple multi-view fusion method. According to the experimental test-bed, 9 the use of two cameras achieves a better balance between precision (68%) and recall (84%) with 10 respect to the use of a single camera (i.e., 62% precision and 86% recall). In addition, a system evalu-11 ation in temporal terms is provided, showing that missed alerts (false negatives) and wrong alerts 12 (false positives) are typically transitory events. Therefore, adding spatial and temporal redundancy 13 increases the overall reliability of the VCA system. 14

Keywords: Smart mobility; Object Detection; Video Content Analysis; Single Shot Multibox Detector. 15

1. Introduction

Nowadays, the smart city paradigm is changing the asset of the urban environment 17 thanks to the rapid growth of digital technologies and communication infrastructures. By 18 interconnecting people and things, smart cities scenarios provide more efficient, fast, ubiq-19 uitous, and accessible services to the citizens [1]. In this context, smart mobility applications 20 are empowered by the high speed and low latency properties of 5G networks [2], being 21 suitable for ensuring road safety [3] and monitoring dangerous situations [4]. The huge 22 amount of sensor data and the availability of fast computing resources at the edge of the 23 5G networks have paved the way to advanced Deep Learning (DL) models for real-time 24 Video Content Analysis (VCA) scenarios [5]. 25

Both real-time localization and object classification methods from video streams are 26 mandatory requirements for VCA solutions. To this aim, different DL architectures based 27 on Convolutional Neural Networks (CNNs) have been recently proposed [6]. However, 28 among the most widely exploited approaches, You-Only-Look-Once (YOLO) and Single 29 Shot Multibox Detectors (SSD) algorithms stand out for their performances and computing 30 efficiency [7]: the former is indeed one of the fastest and most accurate networks for real-31 time object detection [8], while the latter is a benchmark for real-time multi-class object 32 detection at different scales [9]. 33

In this paper, we consider a driver alert scenario, where an urban intersection is monitored by two cameras and an SSD-based object detection model is trained to identify, localize and, eventually, signal the presence of obstacles to public transportation vehicles

Citation: Lenatti M., Narteni S., Paglialonga A., Rampa V., Mongelli M. Dual-view Single Shot Multibox Detector at Urban Intersections: Settings and Performance Evaluation. *Sensors* **2023**, *1*, 0. https://doi.org/

Received: Accepted: Published:

Article

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Copyright: © 2023 by the authors. Submitted to *Sensors* for possible open access publication under the terms and conditions of the Creative Commons Attri-bution (CC BY) license (https:// creativecommons.org/licenses/by/ 4.0/). approaching the surveilled area. Particular focus will be addressed on investigating the 37 advantages of using two cameras instead of a single one, in terms of object detection and 38 alert generation performances. To this purpose, the VCA model will be evaluated using 39 qualitative and tailored quantitative approaches, exploiting both spatial and temporal 40 redundancy. 41

The paper is organized as follows. First, we discuss relevant literature on the topic. 42 Then, we recall the SSD-based method adopted, and we thoroughly describe the on-field 43 implementation. Finally, we present and discuss the results in terms of object detection 44 performance and the related alert generation performance.

2. Related Works

Object detection and/or tracking via multiple camera sensors is a widespread topic 47 in computer vision research. Multi-view 3D object recognition [10] consists in reducing 48 complex 3D object classification tasks to simpler 2D classification tasks, by rendering 49 3D objects into 2D images. Real objects are surrounded by cameras posed at different 50 viewpoints with configurations leading to multi-view proposals such as MVCNN [11], 51 GVCNN [12], View-GCN [13] and RotationNet [14] architectures. These methods use the 52 most successful image classification networks, i.e., VGG, GoogleNet, AlexNet, ResNet, 53 as backbone networks. Then, global 3D shape descriptors are obtained by aggregating 54 selected multi-view features through approaches that account for both content and spatial 55 relationships between the views. 56

Transfer learning approaches prove extremely useful, especially when dealing with 57 scarcely available data. To this end, several open source datasets for object detection 58 in urban traffic optimization and management have recently become available. These datasets focus either on pedestrian or vehicle tracking and detection, combining inputs 60 from multiple cameras and extending visual coverage (e.g., [15,16]).

An overview of recent multi-camera solutions for object detection is presented below. 62 In [17], a novel multi-view region proposal network that infers the vehicles position on 63 the ground plane by leveraging multi-view cross-camera scenarios is presented, whereas 64 an end-to-end DL method for multi-camera people detection is studied in [18]. In [19], 65 a vehicle detection method that applies transfer learning on two cameras with different 66 focal length is proposed. The processing consists of two steps: first, a mapping relationship 67 between input images from the cameras is calculated offline through a robust evolutionary 68 algorithm; then, CNN-based object detection is performed online. More specifically, after 69 a vehicle region is detected from one camera, it is transformed into a binary map. This 70 map is then used to filter CNN feature maps computed for the other camera's image. It 71 is important to outline that finding the relationship between the two cameras is crucial 72 to solve the problem of duplicated detection, as different cameras may focus on the same 73 vehicles. The same problem is raised in [20], where the Authors present a novel edge-AI solution for vehicles counting in a parking area monitored via multiple cameras. They 75 combine a CNN-based technique for objects localization with a geometric approach aimed 76 at analyzing the shared area between the cameras and merging data collected from them. 77 Multi-camera object detection is also investigated in [21], which presents an autonomous drone detection and tracking system exploiting a static wide-angle camera and a lower-79 angle camera mounted on a rotating turret. In order to save computational resources and time, the frame coming from the second camera is overlaid on the static camera's 81 frame. Then, a lightweight version of YOLOv3 detector is developed to perform the object 82 detection. Another recent work on multi-camera fusion for CNN-based object classification 83 [22] devised three fusion strategies: early, late and score fusion. A separate CNN was first 84 trained on each camera. Afterwards, feature maps were stacked together and processed 85 either from the initial layers (early fusion) or at the penultimate layers (late fusion). In addition, score fusion was performed, by aggregating the softmax classification scores 87 in three possible ways: by summing, or by multiplying, the scores across cameras, or by 88 taking the maximum score across them. Results showed that late and score fusion led 89

46

59

61

to an accuracy improvement, with respect to early fusion and single camera proposals. Multi-camera detection has gained increasing importance in several areas besides smart 91 mobility applications. For example, several solutions have been recently proposed in the 92 area of fall detection for remote monitoring of fragile patients. In [23], multi-camera 93 fusion is performed by combining models trained on single cameras together into a global ensemble model at the decision-making level, providing higher accuracy with respect to 95 local single-camera models and avoiding computationally expensive cameras calibration. 96 The Dual-Stream Fused Neural Network method, proposed in [24], first trains two deep 97 neural networks to detect falls by using two single cameras and then merges the results 98 through a weighted fusion of prediction scores. The obtained results overcome the existing 99 methods in this domain. 100

All these proposals deal with high-intensity computational methods while, on the 101 contrary, real-time field-deployable applications impose computational complexity con-102 straints as well. To solve this key issue, we propose here a simple but effective dual-view 103 fusion and detection method and compare its performances with real field experiments 104 [25]. In particular, our solution exploits a transfer learning approach, which consists in 105 training the object detection model on a single camera, in updating it through an additional 106 training by feeding the other camera's images, and then by fusing the single detection 107 signals to generate alerts at the decision level. This speeds up the overall training time and 108 saves computational resources, with respect to other existing decision-making level camera 109 fusion approaches such as [22,23]. 110

34567

3. Video Content Analysis system

3.1. Single Shot Multibox Detector Model

The SSD network is composed of a *backbone* stage for feature extraction and a *head* 114 stage for determining the output. The backbone is a Feature Pyramid Network (FPN)[26], 115 which is a CNN able to extract feature maps representing objects at different scales. It 116 comprises a bottom-up pathway connected to a top-down pathway via lateral connections. 117 The SSD head is a sequence of output maps, which determines the output of the network in 118 the form of bounding boxes coordinates and objects classes. Also, the SSD network exploits 119 the concept of *priors* (also known as *anchor boxes*), a special kind of boxes whose predefined 120 shape can guide the network to correctly detect objects of the desired class. 121

The SSD head is composed of multiple output maps (grids) with different sizes. Each grid decomposes the image into cells, and each cell expresses whether or not it belongs to a particular object, in terms of bounding box coordinates and object class. Lower resolution output maps (i.e., smaller size grids), having larger cells, can detect larger scale objects; in contrast, larger size output grids, having denser cells, are used to predict smaller objects. The use of multiple outputs improve the accuracy of the model significantly, while maintaining the ability to predict objects in real-time.

3.1.1. Loss Function

The training of the SSD model is based on the minimization of the following loss function \mathcal{L} :

$$\mathcal{L} = \mathcal{L}_{loc} + \mathcal{L}_{conf} + \mathcal{L}_{boxiness},\tag{1}$$

where \mathcal{L}_{loc} evaluates the object localization of the model, \mathcal{L}_{conf} evaluates the object classification ability and $\mathcal{L}_{boxiness}$ term refers to the *boxiness*, i.e. the ability of discriminating boxes from background throughout SSD output grids.

Considering object localization, we define $\mathbf{y}_{gt} = (x, y, w, h)$ as the ground truth box coordinates vector for a generic object, with x, y expressing box center coordinates, wthe box width and h the box height. Similarly, we denote with $\mathbf{y}_{pr} = (x_{pr}, y_{pr}, w_{pr}, h_{pr})$ the predicted box coordinates vector for that same object. A discrepancy between real

129

111

112

and predicted box positions is measured by the vector $\mathbf{a} \doteq |\mathbf{y}_{gt} - \mathbf{y}_{pr}|$, with coordinates (a_1, a_2, a_3, a_4) = $(|x - x_{pr}|, |y - y_{pr}|, |w - w_{pr}|, |h - h_{pr}|)$.

The \mathcal{L}_{loc} term is then computed through the Pseudo-Huber loss function [27]:

$$\mathcal{L}_{loc} = \sum_{i=1}^{4} \delta^2 \left(\sqrt{1 + \left(\frac{a_i}{\delta}\right)^2} - 1 \right), \tag{2}$$

with δ being a fixed quantity that controls the steepness of the function. The pseudo-Huber loss provides the best performances, with minimal computational costs w.r.t. the Huber and other type of loss functions [28]. In this study δ was set to 1.5, following preliminary training runs. .

Referring to object classification, let y_c be the true class label for each class c = 1, ..., N, where N is the number of classes. Also, let \hat{p}_c be the corresponding class probability estimates. The second loss term, \mathcal{L}_{conf} , is then a cross-entropy loss, computed as follows:

$$\mathcal{L}_{conf} = -\sum_{c=1}^{N} y_c \log(\hat{p}_c) \tag{3}$$

After prediction, the SSD model also outputs an estimate of the boxiness, expressed as a real value $b_{pr} \in [0, 1]$, which can be interpreted as the model confidence in recognizing whether any object is present in each cell of the network output grids. Consequently, the quantity $b_{bg} = 1 - b_{pr}$ defines the level of confidence of each cell to be part of the background.

The last term $\mathcal{L}_{boxiness}$ relies on a focal loss function [29], which is chosen for its ability to penalize the false positives, i.e. the background points wrongly detected as objects by the model. The boxiness loss $\mathcal{L}_{boxiness}$ is then computed as:

$$\mathcal{L}_{boxiness} = -\left[\alpha \, b_{bg}^{\gamma} log(b_{pr}) + (1 - \alpha) \, b_{pr}^{\gamma} \, log(b_{bg})\right],\tag{4}$$

where the parameter α acts as a weight for those cells being covered by a box and $1 - \alpha$ acts as weight for background cells; the parameter γ controls the shape of the function: higher values of γ require lower loss values to better distinguish boxes from background (i.e., to have $b_{pr} > 0.5$). The attention of the model is thus devoted to the harder-to-detect samples.

3.1.2. Network parameters, training and testing

Non-Maximum Suppression (NMS) [30] was performed to refine the predictions of the 158 model. Indeed, it may often occur that multiple boxes are predicted for the same ground 159 truth object. NMS algorithm filters out the predicted boxes based on the class confidence and the Intersection over Union (IoU) method [31] between them. In particular, for a given 161 SSD output grid and class, for each real object, the predicted box (if any) with the highest 162 class confidence is picked. This box is then chosen as a reference to compute the IoU 163 between itself and all the other predicted boxes, keeping only those with a value below a 164 threshold. In our case, we fixed this threshold at 0.1. Choosing such a low value allows 165 to filter out boxes characterized by even small overlaps with the reference one, therefore 166 reducing the presence of false positives. 167

Tab. 1 summarizes properties and parameters of the SSD model adopted in this work.168The choice of the SSD output grids dimensions was guided by a preliminary analysis on169a range of suitable values, performed to individuate a proper balance between model170accuracy and computational complexity. Also, the selection of Regions of Interest from the171foreground area, as better detailed in the scenario definition, required lower sized grids,172able to capture bigger foreground objects. The network is trained to recognize 3 classes of173objects: 'Vehicle', 'Rider', and 'Pedestrian'.174

Output grids	24×40, 12×20, 6×10 and 3×5
Priors	1×1 , 2×1 , 4×1 , 1×4 and 1×2
# trainable parameters	5000
Learning rate	10^{-4}
δ	1.5
α	0.85
γ	2
IoU threshold	0.1

Table 1. Parameters and properties of the adopted SSD.

3.2. VCA Architecture

We define here the main pipeline of the VCA system for alert generation, whose 176 inference and training/retraining flowcharts are sketched in Fig. 1. The first pipeline 177 (Fig. 1a) sketches the object detection blocks employed to generate alerts (inference phase) 178 by exploiting image fusion on both cameras. The second pipeline (Fig. 1b) focuses on 179 retraining the baseline SSD by adding *TLC2* images via transfer learning, thus obtaining a 180 final model, i.e., SSD_{ret}.More specifically, the inference block diagram shows the real-time 181 processing pipeline adopted to generate the alarm signal A_L by fusing together the single-182 view alerts A_{L1} and A_{L2} produced by the Alert Generation blocks AG_1 and AG_2 that are fed 183 by the output of the SSD_{ret} object detectors attached to the single camera TLC1 and TLC2, 184 respectively. The two cameras have a broad field of view, but in order to define the area 185 of potential danger to be monitored, a Region Of Interest (ROI) has been determined and 186 adapted for each camera. The alert A_L is then employed to alert the driver by activating 187 visual and acoustic alarms on the bus console. The final inference stage A_L is designed to 188 integrate the two independent outputs of the single alert generators related to each camera 189 view and to perform information fusion at the decision level with the aim of increasing the 190 overall reliability and accuracy of the system. 191

Fig. 1b shows the retraining procedure adopted to update the baseline SSD network 192 of the single-view system (that uses only *TLC*1 data) by including also images from the 193 *TLC2* camera. In fact, the baseline SSD model (i.e., the green block in Fig.1b) has been 194 preliminarily trained on a set of images extracted from three open-source dataset (Open 195 Images Dataset [32], ETH Pedestrian Dataset [33] and EuroCity Dataset [34]) that contain 196 annotated images of urban traffic scenes. Afterwards, the images captured by TLC1 were 197 added to these datasets to complete the training of the baseline SSD model. To further 198 improve the flexibility, reliability and, in particular, the detection accuracy of the VCA 199 system, the baseline SSD model was later retrained on a set of 10,000 additional images 200 acquired from the TLC2 camera. The term retraining refers to the procedure of updating 201 the parameters of a previously trained model based on the addition of new data by transfer 202 learning methods [35]. From now on, we will refer to the final retrained model as SSD_{ret} 203 (blue block in Fig. 1b). The generalizing capabilities of the baseline SSD and the retrained 204 SSD_{ret} models were assessed using a test dataset consisting of frames extracted from a 1-hour video, for both cameras. Both videos were first synchronized and cut to align the 206 start and end time stamps, then converted from the h263 format to the mp4 format using 207 the FFmpeg tool [36] (with compression factor 1.25). Finally, 1,000 frames were extracted 208 for each recording.

3.3. Data labeling

YOLOv5x[37], one of the state-of-the-art YOLO networks for object detection in realtime applications, was adopted to define ground truth boxes, i.e., to label the objects actually present in each image. For this purpose, YOLOv5x was applied on each image of the

175



(b) Retraining

Figure 1. Flowchart of the procedures exploited for the proposed VCA system, sketching the alert generation and fusion (Fig.1a) based on model SSD_{ret}, obtained via a retraining process (Fig.1b).ret.

training, retraining, and test datasets in order to recognize objects of the classes 'Car', 'Bus', 214 'Truck', 'Motorcycle', 'Bicycle', and 'Pedestrian'. Then, these classes were grouped into three 215 more generic classes, namely 'Vehicle', 'Rider', and 'Pedestrian'. Ground truth boxes were 216 provided in the YOLO format (x_{center}, y_{center}, width, height), and subsequently converted 217 in the SSD format (x_{min} , y_{min} , width, height). The results of this automatic labeling step 218 were then manually inspected to verify the presence of sufficiently accurate ground truth 219 boxes. In presence of detection errors inside the monitored area, the corresponding images 220 were removed from the dataset. Based on the ground truth boxes, we also defined the 221 number of ground truth alerts, which were raised anytime at least one ground truth box 222 was detected within the ROI. 223

4. Driver Alert Use Case

4.1. Scenario definition

Piazza Caricamento is one of the locations with the highest concentration of pedestrian and road traffic in the historic center of Genoa, Italy, as it connects the east and west 227 areas of the city and, above all, it is located nearby the main tourist attractions (e.g., the 228 aquarium, the pedestrian area on the harbor, and the most important architectural and 229 artistic sites of the city). The area monitored by the proposed VCA system is the intersection 230 between the pedestrian area of the harbor, the vehicular access to the parking lot, and the 231 access roads to the underground tunnel below Piazza Caricamento corresponding to the 232 latitude and longitude coordinates 44.4110720922656 N, 8.928069542327654 E (expressed 233 in decimal degrees). A dedicated public transportation bus lane which is characterized 234 by limited visibility interconnects with the monitored intersection. . The area is often 235 crowded with pedestrians and vehicles frequently passing through to access the car parking. 236 Hence, potential collisions with buses coming from their dedicated lane represent a real 237 risk scenario that makes Piazza Caricamento a suitable location where to implement a 238 VCA system. The proposed solution consists in an automatic system able to detect 239 the presence of pedestrians and/or vehicles inside the area via VCA processing, and to 240

generate an appropriate alert to the bus approaching the intersection. Real-time monitoring is performed via two Bosh DINION IP Bullet 6000I HD, 2,8 - 12 MM cameras, which are professional surveillance HD cameras compliant to the SMPTE 296M-2001 standard [38] and ONVIF profiles G and S [39] to guarantee the interoperability with the AI components. We will refer to these cameras as *TLC1* and *TLC2*.

As previously noted, only objects within each ROI of the cameras can generate an alert to be sent to the driver. As a result, the two ROIs strictly overlap. Since our SSD model involves multiple output grids, the ROI was resized for each of them based on their dimension. Fig. 2 displays the fields of view covered by the two cameras and reports the selected ROIs for each adopted grid.



(a) ROI on *TLC1* for output grid of size 3x5



(c) ROI on *TLC1* for output grid of size 6x10



(b) ROI on *TLC2* for output grid of size 3x5



(d) ROI on *TLC2* for output grid of size 6x10



(e) ROI on *TLC1* for output grid of size 12x20

(f) ROI on *TLC2* for output grid of size 12x20

Figure 2. Regions Of Interest (ROIs) inside the monitored area (green rectangles), for each considered SSD output grid on *TLC1* (left column) and *TLC2* (right column).

As it will further emphasized in the following sections, the main goal of our work is to understand to what extent the joint use of two cameras can represent an added value for the VCA task with respect to the use of a single camera (either *TLC1* or *TLC2*).

4.2. Performance evaluation

Two types of performance figures will be considered to evaluate the VCA monitoring 255 system, namely *object detection performance*, that is the ability of the system to correctly 256

identify different classes of objects inside the ROI, and *alert generation performance*, that is the ability of the system to trigger an alert if and only if at least one object is present in the monitored area.

For the sake of simplicity, the system performances were assessed considering only 3 grids (i.e., 12x20, 6x10 and 3x5) with priors of size 1x2 (more suitable for identifying people) and priors of size 2x1 (more suitable for identifying vehicles). 262

Finally, the VCA system performances were evaluated also in terms of computation time required for object detection and alert generation. The average inference time per frame was assessed locally on a host equipped with an Intel Core i5 dual-core processor at 2.6 GHz, 8GB RAM memory banks, and running the macOS 10.15.7 operating system.

4.2.1. Object detection performance

The ability of each component of SSD_{ret} (according to the aforementioned grids and 268 priors) to identify objects of different classes inside the ROI was evaluated by calculating the 269 average confusion matrix over the whole test dataset, for each camera, namely the average 270 number of correctly identified objects (TP_{obi}), the average number of undetected objects 271 (FN_{obi}), and the average number of objects detected but not actually present in the ground 272 truth image (FP_{obi}). The obtained values were then compared with the average number of 273 real objects per image. Then, in order to measure the object detection performance from 274 a comprehensive point of view, precision (PRE_{obi}) and recall (REC_{obi}) were assessed for 275 each considered frame, both individually for single grids and priors, and aggregating all 276 outputs. Precision measures the number of correctly identified objects to the total number 277 of detected objects, whereas recall measures the number of correctly detected objects to 278 the total number of ground truth objects. These metrics were then averaged across all the 279 frames in the test dataset (i.e., 1000 frames). 280

4.2.2. Alert generation performance

The ability of SSD_{ret} to generate alerts when an object is inside the ROI was assessed 282 by calculating the confusion matrix over the entire test dataset, considering two possible 283 outputs of the system, namely the presence of an alert (*alert=1*) or its absence (*alert=0*), 284 for each input image. The following elements of the confusion matrix were considered: 285 the total number of correctly generated alerts (TP_{alert}), the total number of ground truth 286 alerts not triggered by the system (FN_{alert}), the total number of alerts incorrectly triggered 287 by the system (FP_{alert}), and the total number of non-alert situations in which the alert is 288 correctly not triggered by the system (TN_{alert}). It is also important to underline that, in 289 light of the technological implementation of the alerts triggering system of each camera, 290 incorrect alerts (either FN_{alert} or FP_{alert}) were only triggered when no true positives had 291 already been generated for the same image. 292

As previously described, SSD models provide different outputs from output maps of different sizes. Therefore, system performance was first evaluated by considering alerts detected individually by each grid and prior and then by evaluating the total amount of alerts identified by the aggregation of all grids and all priors. Alert generation performance was evaluated both individually on the two cameras (*TLC1* and *TLC2*, separately) and then on their fusion. In the latter case, an alert is generated when at least one of the two cameras detects an object within the ROI.

Since the frames considered in our use case are temporally continuous, we also decided to evaluate if the presence of FN_{alert} and FP_{alert} could be considered as a transient phenomenon or not. Hence, we computed also the FN^*_{alert} and FP^*_{alert} , representing the false negatives and false positives occurred at least in two consecutive frames. Any FN_{alert} or FP_{alert} events present in just one frame were therefore considered spurious and avoided by waiting for the next frame before doing inference.

267

5. Results

5.1. Object detection performance

3.1.2A base model was trained on a set of images composed by *TLC1* images and 308 external images from open-source datasets on mobility scenarios. The base model was then 309 retrained on a dataset extracted from TLC2 recordings yielding SSD_{ret}. The procedure of 310 retraining (on TLC2 images only) an already pre-trained model offers several advantages 311 over training from scratch (using TLC1 and TLC2 images). Notably, retraining was faster 312 than the full training. Specifically, the time required to retrain the model was more than 313 10 times shorter than the original training time of the baseline SSD (i.e., 42 hours). Tab. 2 314 reports the obtained object detection performance for each camera, each grid, and each 315 prior separately in terms of mean confusion matrix over the entire test dataset. Average 316 precision and average recall were also computed. 317

Table 2. Mean and standard deviation (between parentheses) of TP_{obj}, FP_{obj}, FN_{obj} and percentage of PRE_{obj} and REC_{obj} for each camera, grid, and prior of the SSD_{ret} model.

TLC1				TLC2								
	#real object	TP _{obj} ts	FP _{obj}	FN _{obj}	PRE _{obj}	REC _{obj}	#real object	TP _{obj} ts	FP _{obj}	FN _{obj}	PRE _{obj}	REC _{obj}
Grid: 12x20 Prior: 1x2	0.19 (0.80)	0.10 (0.55)	0.04 (0.20)	0.08 (0.87)	55%	54%	0.71 (1.43)	0.24 (0.79)	0.22 (0.57)	0.40 (0.96)	43%	31%
Grid: 12x20 Prior: 2x1	0.02 (0.31)	0.02 (0.24)	0.08 (0.36)	0.005 (0.13)	11%	66%	0.34 (1.18)	0.24 (0.99)	0.28 (0.65)	0.10 (0.61)	24%	67%
Grid: 6x10 Prior: 1x2	0.05 (0.35)	0.05 (0.30)	0.15 (0.42)	0.02 (0.23)	19.76%	63.46%	0.13 (0.48)	0.07 (0.36)	0.07 (0.27)	0.06 (0.28)	37%	43%
Grid: 6x10 Prior: 2x1	0.005 (0.08)	0.005 (0.10)	0.18 (0.47)	0.00 (0.00)	1.6%	100%	0.08 (0.40)	0.01 (0.14)	0.08 (0.35)	0.07 (0.43)	15%	21%
Grid: 3x5 Prior: 1x2	0.01 (0.14)	0.00 (0.05)	0.07 (0.25)	0.01 (0.13)	4.11%	42.86%	-	-	1.70 (0.59)	-	0%	-
Grid: 3x5 Prior: 2x1	0.001 (0.03)	0.00 (0.00)	0.08 (0.28)	0.001 (0.03)	0%	0%	0.07 0.33	0.04 0.23	1.63 0.56	0.03 (0.19)	1.3%	61.44%

Table 3. Global object detection performance of SSD_{ret} for each camera by considering all the grids and priors as defined in Tab. 1. Precision: PRE_{obj} ; Recall: REC_{obj} .

	TLC1	TLC2
PRE _{obj}	17%	73%
REC _{obj}	90%	89%

According to Tab. 4, it appears that the *TLC1* images contain fewer ground truth 318 objects inside the ROI than the TLC2 ones. However, no ground truth events filmed by 319 TLC2 are captured by the 3x5 grid with 1x2 prior. Hence, it was not possible to calculate 320 TP_{obj}, FN_{obj} and recall in that case. Since the number of false positives is on average higher 321 than the number of false negatives, PRE_{obj} is lower than REC_{obj}, except when considering 322 a 12x20 grid with 1x2 prior. In addition, we can observe how grids with a larger number 323 of cells (i.e., 12x20 and 6x10) are generally able to detect more objects than the smallest 324 grid (i.e., 3x5). This may be due to the fact that objects within the ROI are typically in the 325

334

background and thus more easily detected by denser grids, characterized by smaller cell sizes.

The global object detection performances of SSD_{ret} on both cameras were then evaluated in terms of precision and recall and reported in Tab. 3. These values were obtained by considering all the grids and priors used to define the model's architecture (as defined in Tab. 1). *TLC1* yielded a low precision of about 17% and a satisfying recall, equal to about 90%. In contrast, *TLC2* yielded a much higher precision of about 73% and recall similar to *TLC1* (i.e., about 89%).

5.2. Alert generation performances

Alert generation performance was first evaluated separately on the two cameras and 335 then considering the fusion between the alerts generated by the two, as shown in Tab. 336 4. The results reported in Tab. 4 are consistent with those shown in Tab. 2, since grids 337 with a larger number of cells (i.e., 12x20 and 6x10) are able to generate more alerts than 338 the smallest grid (i.e., 3x5). In particular, with the exception of the 3x5 grid, that mostly 339 detects vehicles, most of the alerts seem to be raised by objects that correspond to the prior 340 of size 1x2 (i.e., pedestrians in the ROI). From these results, we can observe that both the 341 number of ground truth alerts and the number of correctly predicted alerts (TP_{alert}) increase when considering the data fusion of both cameras (*fusion*(*TLC1*,*TLC2*)), compared to the 343 individual *TLC1* and *TLC2*. Fig. 3 shows an example of an alert correctly detected by *TLC2*, 344 but not by *TLC1*. This image would therefore constitute a FN event considering only *TLC1*, 345 but it is correctly classified as a TP event when *fusion(TLC1,TLC2)* is considered.

Table 4. Number of ground t	truth alerts and TP _{alert} for	r each grid and	l prior using sing	le camera
processing (TLC1, TLC2) and d	lata fusion of both cameras	s (fusion(TLC1,T	TLC2))	

	TLC1		TLC2		fusion(TLC1,TLC2)		
	Ground truth alerts	TP _{alert}	Ground truth alerts	TP _{alert}	Ground truth alerts	TP _{alert}	
Grid: 12x20 Prior: 1x2	62	54	41	27	76	61	
Grid: 12x20 Prior: 2x1	9	6	29	25	34	28	
Grid: 6x10 Prior: 1x2	66	49	29	23	76	57	
Grid: 6x10 Prior: 2x1	8	8	3	0	11	8	
Grid: 3x5 Prior: 1x2	3	2	2	0	5	2	
Grid: 3x5 Prior: 2x1	7	4	3	0	10	4	

If we focus, for example, on the 12x20 grid and the 1x2 prior (Tab. 4), we can observe $_{347}$ that *TLC1* alone detects 62 ground truth alerts (54 TP_{alert}), while *TLC2* detects 41 ground $_{348}$ truth alerts (27 TP_{alert}) and *fusion(TLC1,TLC2)* detects 76 ground truth alerts (61 TP_{alert}). $_{349}$ These results confirm how different grids and priors are able to identify different objects, $_{350}$ and consequently generate different alerts. For this reason, we finally evaluated the global $_{351}$ alert generation performances, obtained by combining all the outputs provided by different $_{352}$ this global evaluation are reported in Tab. 5. $_{354}$

The estimated average elapsed time during the inference phase for the whole alert generation process on a single camera is about 0.46 seconds per frame, while the elapsed time of the decision fusion is about $1.8 \cdot 10^{-6}$ seconds and may be neglected. Thus, the total



Figure 3. Example of the same object correctly detected within the ROI (green area) by *TLC2* (right), but missed by *TLC1* (left). Ground truth boxes are shown in blue while predicted boxes are shown in red.

Table 5. Ground truth alerts, TP_{alert} , TN_{alert} , FP_{alert} , FN_{alert} , FP^*_{alert} and FN^*_{alert} obtained from all grids and priors, on single cameras (*TLC1*, *TLC2*) and their fusion (fusion(TLC1,TLC2))

	Ground truth alerts	TP _{alert}	TN _{alert}	FP _{alert}	FN _{alert}	FP* _{alert}	FN* _{alert}
TLC1	89	77	865	46	12	2	0
TLC2	74	59	908	18	15	1	2
fusion(TLC1,TLC2)	125	105	827	48	20	3	3

inference time of the multi-camera VCA system (not parallelized) is about 0.92 seconds per frame. 358

6. Discussion

A VCA monitoring system based on a SSD architecture has been implemented and evaluated in terms of its ability to detect objects in the surveilled area and its related ability to generate alerts. Specifically, the VCA system foresees possible dangerous situations inside a intersection through the use of a multi-camera deep learning-based object detection system. The choice to merge data at decision level was motivated by its simplicity that allows to operate within the time constraints dictated by a real-time application. In addition, the system built in this way can easily compensate for the lack of one of the two possible inputs, ensuring robustness against possible failures or damages to the system.

Comparing the *TLC1* and *TLC2* cases, it can be seen that the former has a rather low precision in detecting objects, . This result is further confirmed by the performances of 370 alert generation (Tab. 5). Provably, the precision of *TLC1* in terms of alert generation is lower than the corresponding TLC2 precision (i.e., 62% and 77%, respectively). As 372 a result, fusion(TLC1,TLC2) reaches a higher precision (i.e., 69%) with respect to TLC1 373 alone. In contrast, the recall of *TLC1* in terms of alert generation is slightly higher than the 374 corresponding TLC2 recall (i.e., 86% and 80%, respectively). As a result, fusion(TLC1,TLC2) 375 yields to a higher recall (i.e., 84%) with respect to TLC2 alone. In summary, by combining 376 the two cameras, there is a significant increase in precision with respect to TLC1 alone and 377 a slight improvement in recall compared to *TLC2*. The monitoring system based on SSD_{ret} 378 yields a quite satisfactory alert generation accuracy when considering a single camera (i.e., 379 about 94%). This means that the retraining phase did not erase what the model learned from *TLC1* images, i.e., there is no catastrophic forgetting [40]. Although accuracy remains 381 almost stable (93%) when considering *fusion*(*TLC1*,*TLC2*), the introduction of a second 382 camera TLC2 improves the overall safety by allowing the identification of a higher number 383 of real dangerous situations (i.e., 125 ground truth alerts) within the area of interest. In fact, the combination of *TLC1* and *TLC2* enables the triggering of 40% more ground truth 385 alerts than *TLC1* alone. The increase in the number of alerts is mainly due to the different framing of the two cameras, and thus the increased field of view of the object detection 387 system. Consequently, also the absolute number of TP_{alert} increases (from 77 to 105) after 388

12 of 15

the outputs of the two cameras are merged. Since we are dealing with a highly unbalanced dataset, where the number of dangerous situations is considerably lower than the number of safe situations, it could be useful to evaluate the F1-score. Specifically, it can be seen that the use of two cameras results in an F1-score of 75%, which is higher than that obtained by using *TLC1* alone (i.e., 73%).

By using not only spatial redundancy, i.e., the different views of the same monitored 394 area captured by TLC1 and TLC2, but also the temporal continuity of the frames, we can 395 design a post-processing algorithm that uses the information of two or more consecutive 396 video frames instead of a single one as assumed so far. In this case, the actual output alert 397 signal is generated if it is triggered by at least two consecutive frames. By exploiting the 398 temporal continuity, the amount of wrong predictions is reduced, as indicated in Tab. 5, 399 where FP*_{alert} and FN*_{alert} (i.e., the number of FPs and FNs persisting in at least two con-400 secutive frames) is consistently lower than FP_{alert} and FN_{alert}, respectively. This reduction 401 in the number of false and missed alarms proves that FPs and FNs are generally spurious 402 events that can be easily removed by considering a certain time window. However, it is 403 worth noting that this method introduces a one-frame delay in the alert signal generation 404 stage. 405

In addition, a local evaluation of the total inference time per frame was performed, demonstrating the ability of the proposed multi-camera VCA system to generate the alert in a sufficiently short time (less than 1 second), that is compatible with the system requirements to make a decision in real time. However, more precise evaluations will be needed following specific on-site deployment.

This study presents some limitations. First of all, the multi-camera system was evalu-411 ated using a single fusion technique directly applied at the decision level. In future studies, 412 different data fusion techniques including early and late fusion at different depths of the 413 network should be compared to evaluate possible further improvements in terms of the sys-414 tem reliability. Moreover, although the network was originally trained on a heterogeneous 415 set of images from the experimental test-bed (TLC1) and open source datasets, the dataset 416 used for retraining SSD_{ret} included only TLC2 frames captured in daytime. Therefore, it 417 will be necessary to evaluate the system's ability to generalize in different scenarios, such 418 as its robustness in different weather and light conditions (e.g., day/night, sunny/rainy 419 weather). Lastly, at the current stage, possible security issues following malicious attacks 420 on the main components of the system (e.g., cameras, onboard units, edge servers) have 421 not been considered yet. In particular, the alert generation system could be vulnerable to 422 adversarial attacks aimed at changing the output of the system, which could cause potential dangerous situations. In the future, it will be necessary to devise robust solutions to these 474 types of attacks, such as considering the introduction of a Bayesian layer in the vision 425 system [41]. 426

7. Conclusion

This work focuses on the development and evaluation of an Single Shot multibox 428 Detector-based object detection system applied to a urban scenario. In particular, we 429 evaluated the effectiveness of adding a second camera (TLC2) in terms of detecting potential hazardous situations within the region of interest. The introduction of a second camera, 431 in addition to the first one (TLC2), not only makes the Video Content Analysis system 432 more robust w.r.t. possible failures due to TLC1 malfunctions, but also leads to a higher 433 number of correctly detected alarms thanks to a wider coverage of the surveilled area. 434 Furthermore, the number of False Negatives (FN-type) events is reduced by considering 435 temporal continuity in successive frames. In the specific Smart Mobility use case, FN-type 436 errors were considered to be more important than False Positive (FP-type) errors. Indeed, 437 the number of negative events misclassified as positive (i.e., FP-type), will result in alarms 438 that do not correspond to the presence of objects or obstacles in the region of interest. Such 439 errors are considered less critical because they simply cause unnecessary alerts to be sent, if 440 few, without endangering the driver. However, in the long run, these redundant alarms 441

may make the driver less confident in the system's ability to correctly identify dangerous situations. Future studies will focus on further validation of the proposed solution. Finally, the formalization of an algorithm that can leverage the temporal continuity provided by videos, instead of relying on individual frames, could be investigated.

Author Contributions: Conceptualization, M.M., A.P. and V.R.; methodology, M.L., M.M., S.N.,A.P. and V.R.; software, M.L. and S.N.; validation, M.L., M.M., S.N., and V.R.; formal analysis, M.L.and S.N.; investigation, M.L., M.M., S.N., and V.R.; resources, M.M. and V.R.; data curation, M.L.and S.N.; writing—original draft preparation, M.L., M.M., S.N., A.P. and V.R.; writing—review andediting, M.L., M.M., S.N., A.P. and V.R.; visualization, M.L., S.N., V.R.; supervision, M.M.; projectadministration, M.M.; funding acquisition, M.M. and V.R. All authors have read and agreed to thepublished version of the manuscript.

Funding: The work was carried out within the Genova 5G project, a tender by the italian Ministry of
Economic Development for the acquisition of technologies aimed at the safety of road infrastructures
in the territorial area of Genova through experiments based on 5G technology. The project has just
ended in October 2022.453454

Informed Consent Statement: Not applicable.

Data Availability Statement: The data presented in this study may be available upon request to the corresponding author. 459

Acknowledgments: The authors would like to thank Vodafone, being the administrative and technical
coordinator of the Genova 5G project, as well as network operator and 5G technology enabler. The
authors are also grateful to Aitek S.p.A. (Vanessa Orani, Stefano Delucchi and Bernardo Pilarz) for
their assistance in the development of the VCA solution. The authors would also like to thank all
partners involved in the project: Azienda Mobilità e Trasporti SpA, Genova, Comune di Genova,
Leonardo, Start 4.0.460

Marta Lenatti is a PhD student enrolled in the National PhD in Artificial Intelligence, XXXVIII cycle, 4 course on Health and life sciences, organized by Università Campus Bio-Medico di Roma.

Conflicts of Interest: The authors declare no conflict of interest.

References

- Founoun, A.; Hayar, A. Evaluation of the concept of the smart city through local regulation and the importance of local initiative. In Proceedings of the 2018 IEEE International Smart Cities Conference (ISC2). IEEE, 2018, pp. 1–6.
- Savithramma, R.; Ashwini, B.; Sumathi, R. Smart Mobility Implementation in Smart Cities: A Comprehensive Review on State-of-art Technologies. In Proceedings of the 2022 4th International Conference on Smart Systems and Inventive Technology (ICSSIT). IEEE, 2022, pp. 10–17.
- Celidonio, M.; Di Zenobio, D.; Fionda, E.; Panea, G.G.; Grazzini, S.; Niemann, B.; Pulcini, L.; Scalise, S.; Sergio, E.; Titomanlio, S. 475 Safetrip: a bi-directional communication system operating in s-band for road safety and incident prevention. In Proceedings of the 2012 IEEE 75th Vehicular Technology Conference (VTC Spring). IEEE, 2012, pp. 1–6.
- Wen, J.; He, Z.; Yang, Y.; Cheng, Y. Study on the factors and management strategy of traffic block incident on Hangzhou Province Highway. In Proceedings of the 2020 International Conference on Intelligent Transportation, Big Data & Smart City (ICITBS).
 IEEE, 2020, pp. 67–71.
- Mauri, A.; Khemmar, R.; Decoux, B.; Ragot, N.; Rossi, R.; Trabelsi, R.; Boutteau, R.; Ertaud, J.Y.; Savatier, X. Deep learning for real-time 3D multi-object detection, localisation, and tracking: Application to smart mobility. Sensors 2020, 20, 532.
- Jiao, L.; Zhang, R.; Liu, F.; Yang, S.; Hou, B.; Li, L.; Tang, X. New generation deep learning for video object detection: A survey. *IEEE Transactions on Neural Networks and Learning Systems* 2021.
- Chen, Z.; Khemmar, R.; Decoux, B.; Atahouet, A.; Ertaud, J.Y. Real Time Object Detection, Tracking, and Distance and Motion Estimation based on Deep Learning: Application to Smart Mobility. In Proceedings of the 2019 Eighth International Conference on Emerging Security Technologies (EST), 2019, pp. 1–6. https://doi.org/10.1109/EST.2019.8806222.
- Wang, C.Y.; Bochkovskiy, A.; Liao, H.Y.M. YOLOv7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors. arXiv preprint arXiv:2207.02696 2022.
- Liu, W.; Anguelov, D.; Erhan, D.; Szegedy, C.; Reed, S.; Fu, C.Y.; Berg, A.C. Ssd: Single shot multibox detector. In Proceedings of the European conference on computer vision. Springer, 2016, pp. 21–37.
- Qi, S.; Ning, X.; Yang, G.; Zhang, L.; Long, P.; Cai, W.; Li, W. Review of multi-view 3D object recognition methods based on deep learning. *Displays* 2021, 69, 102053.
- Su, H.; Maji, S.; Kalogerakis, E.; Learned-Miller, E. Multi-view convolutional neural networks for 3d shape recognition. In Proceedings of the Proceedings of the IEEE international conference on computer vision, 2015, pp. 945–953.

468 469

- Feng, Y.; Zhang, Z.; Zhao, X.; Ji, R.; Gao, Y. Gvcnn: Group-view convolutional neural networks for 3d shape recognition. In Proceedings of the Proceedings of the IEEE conference on computer vision and pattern recognition, 2018, pp. 264–272.
- Wei, X.; Yu, R.; Sun, J. View-gcn: View-based graph convolutional network for 3d shape analysis. In Proceedings of the Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020, pp. 1850–1859.
- Kanezaki, A.; Matsushita, Y.; Nishida, Y. Rotationnet: Joint object categorization and pose estimation using multiviews from unsupervised viewpoints. In Proceedings of the Proceedings of the IEEE conference on computer vision and pattern recognition, 2018, pp. 5010–5019.
- Chavdarova, T.; Baqué, P.; Bouquet, S.; Maksai, A.; Jose, C.; Bagautdinov, T.; Lettry, L.; Fua, P.; Van Gool, L.; Fleuret, F.
 WILDTRACK: A Multi-camera HD Dataset for Dense Unscripted Pedestrian Detection. In Proceedings of the 2018 IEEE/CVF
 Conference on Computer Vision and Pattern Recognition, 2018, pp. 5030–5039. https://doi.org/10.1109/CVPR.2018.00528.
- Tang, Z.; Naphade, M.; Liu, M.Y.; Yang, X.; Birchfield, S.; Wang, S.; Kumar, R.; Anastasiu, D.; Hwang, J.N. CityFlow: A City-Scale Benchmark for Multi-Target Multi-Camera Vehicle Tracking and Re-Identification. In Proceedings of the 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2019, pp. 8789–8798. https://doi.org/10.1109/CVPR.2019.009 00.
- Wu, H.; Zhang, X.; Story, B.; Rajan, D. Accurate Vehicle Detection Using Multi-camera Data Fusion and Machine Learning. In Proceedings of the ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2019, pp. 3767–3771. https://doi.org/10.1109/ICASSP.2019.8683350.
- Chavdarova, T.; Fleuret, F. Deep multi-camera people detection. In Proceedings of the 2017 16th IEEE international conference on machine learning and applications (ICMLA). IEEE, 2017, pp. 848–853.
- Dinh, V.Q.; Munir, F.; Azam, S.; Yow, K.C.; Jeon, M. Transfer learning for vehicle detection using two cameras with different focal lengths. *Information Sciences* 2020, 514, 71–87.
- Ciampi, L.; Gennaro, C.; Carrara, F.; Falchi, F.; Vairo, C.; Amato, G. Multi-camera vehicle counting using edge-AI. *Expert Systems* with Applications 2022, 207, 117929.
- Unlu, E.; Zenou, E.; Riviere, N.; Dupouy, P.E. Deep learning-based strategies for the detection and tracking of drones using several cameras. *IPSJ Transactions on Computer Vision and Applications* 2019, 11, 1–13.
- 22. Seeland, M.; Mäder, P. Multi-view classification with convolutional neural networks. *Plos one* **2021**, *16*, e0245230.
- Ezatzadeh, S.; Keyvanpour, M.R.; Shojaedini, S.V. A human fall detection framework based on multi-camera fusion. *Journal of Experimental & Theoretical Artificial Intelligence* 2022, 34, 905–924.
- Saurav, S.; Saini, R.; Singh, S. A dual-stream fused neural network for fall detection in multi-camera and 360° videos. Neural Computing and Applications 2022, 34, 1455–1482.
- Narteni, S.; Lenatti, M.; Orani, V.; Rampa, V.; Paglialonga, A.; Ravazzani, P.; Mongelli, M. Technology transfer in smart mobility: the driver alert pilot of 5G Genova project. In Proceedings of the Proc. of the 11th World Conference on Information Systems and Technologies (WorldCIST'23), 1st Workshop on Artificial Intelligence for Technology Transfer (WAITT'23), 2023, pp. 1–4.
- Lin, T.Y.; Dollár, P.; Girshick, R.; He, K.; Hariharan, B.; Belongie, S. Feature pyramid networks for object detection. In Proceedings of the Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 2117–2125.
- Charbonnier, P.; Blanc-Feraud, L.; Aubert, G.; Barlaud, M. Deterministic edge-preserving regularization in computed imaging. IEEE Transactions on Image Processing 1997, 6, 298–311. https://doi.org/10.1109/83.551699.
- Barrow, D.; Kourentzes, N.; Sandberg, R.; Niklewski, J. Automatic robust estimation for exponential smoothing: Perspectives from statistics and machine learning. *Expert Systems with Applications* 2020, 160, 113637. https://doi.org/https://doi.org/10.101
 6/j.eswa.2020.113637.
- Lin, T.Y.; Goyal, P.; Girshick, R.; He, K.; Dollár, P. Focal loss for dense object detection. In Proceedings of the Proceedings of the IEEE international conference on computer vision, 2017, pp. 2980–2988.
- Bodla, N.; Singh, B.; Chellappa, R.; Davis, L.S. Soft-NMS-improving object detection with one line of code. In Proceedings of the Proceedings of the IEEE international conference on computer vision, 2017, pp. 5561–5569.
- Kowalczyk, P.; Izydorczyk, J.; Szelest, M. Evaluation Methodology for Object Detection and Tracking in Bounding Box Based Perception Modules. *Electronics* 2022, 11, 1182.
- Krasin, I.; Duerig, T.; Alldrin, N.; Ferrari, V.; Abu-El-Haija, S.; Kuznetsova, A.; Rom, H.; Uijlings, J.; Popov, S.; Kamali, S.; et al. OpenImages: A public dataset for large-scale multi-label and multi-class image classification. Dataset available from https://storage.googleapis.com/openimages/web/index.html 2017.
- Ess, A.; Leibe, B.; Schindler, K.; .; van Gool, L. A Mobile Vision System for Robust Multi-Person Tracking. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR'08). IEEE Press, 2008.
- Braun, M.; Krebs, S.; Flohr, F.B.; Gavrila, D.M. EuroCity Persons: A Novel Benchmark for Person Detection in Traffic Scenes. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 2019, pp. 1–1. https://doi.org/10.1109/TPAMI.2019.2897684.
- Dhillon, A.; Verma, G.K. Convolutional neural network: a review of models, methodologies and applications to object detection.
 Progress in Artificial Intelligence 2020, *9*, 85–112.
- 36. FFmpeg 5.0. https://ffmpeg.org/, accessed July 05, 2022.
- 37. Jocher, G. YOLOv5 by Ultralytics (Version 7.0)[Computer Software], 2020. https://doi.org/10.5281/zenodo.3908559.
- Informative, A.B.P.A. 1280× 720 Progressive Image Sample Structure—Analog and Digital Representation and Analog Interface 2011.

521

551

39.	ONVIF profiles [Available online]. https://www.onvif.org/profiles/, accessed March 01, 2023.	555
40.	French, R.M. Catastrophic forgetting in connectionist networks. Trends in cognitive sciences 1999, 3, 128–135.	556
41.	Y. Pang, S. Cheng, J.H.; Liu, Y. Evaluating the robustness of bayesian neural networks against different types of attacks. In	557
	Proceedings of the CVPR 2021 Workshop on Adversarial Machine Learning in Real-World Computer Vision Systems and Online	558
	Challenges (AML-CV), 2021.	559