# Journal Pre-proof

Multi-camera vehicle counting using edge-AI

Luca Ciampi, Claudio Gennaro, Fabio Carrara, Fabrizio Falchi,
Claudio Vairo, Giuseppe Amato

Please cite this article as: L. Ciampi, C. Gennaro, F. Carrara et al., Multi-camera vehicle counting using edge-AI. *Expert Systems With Applications* (2022), doi: https://doi.org/10.1016/j.eswa.2022.117929.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

# Highlights

- Smart mobility is crucial for smart cities and traffic-related issues.
- We introduce a multi-camera system able to count cars from images of parking areas.
- We combine a deep learning-based technique and a decentralized geometric-based approach.
- All the algorithms run on the edge devices reducing the traffic on the network.
- Our solution benefits from redundant information from different data sources.

# ORCID Information

Below the ORCID information about the authors:

- Luca Ciampi: https://orcid.org/0000-0002-6985-0439
- Claudio Gennaro: https://orcid.org/0000-0002-3715-149X
- Fabio Carrara: https://orcid.org/0000-0001-5014-5089
- Fabrizio Falchi: https://orcid.org/0000-0001-6258-5313
- Claudio Vairo: https://orcid.org/0000-0003-2740-4331
- Giuseppe Amato: https://orcid.org/0000-0003-0171-4315

# *<u>CRediT author statement</u>*

**Luca Ciampi: Writing -** Conceptualization, Software, Validation, Writing - Original Draft, Writing - Review & Editing

**Claudio Gennaro:** Conceptualization, Writing - Review & Editing, Funding acquisition

**Fabio Carrara:** Conceptualization, Investigation, Writing - Review & Editing

**Fabrizio Falchi:** Writing - Review & Editing, Supervision, Funding acquisition

**Claudio Vairo:** Writing - Review & Editing, Supervision

**Giuseppe Amato:** Writing - Review & Editing, Supervision, Funding acquisition

# Multi-Camera Vehicle Counting Using Edge-AI

Revised Manuscript (clean)

Luca Ciampi[a] (luca.ciampi@isti.cnr.it), Claudio Gennaro[a]
(claudio.gennaro@isti.cnr.it), Fabio Carrara[a] (fabio.carrara@isti.cnr.it),
Fabrizio Falchi[a] (fabrizio.falchi@isti.cnr.it), Claudio Vairo[a]
(claudio.vairo@isti.cnr.it), Giuseppe Amato[a] (giuseppe.amato@isti.cnr.it)

[a] Institute of Information Science and Technologies of the National Research Council
of Italy (ISTI-CNR), via G. Moruzzi 1 - 56124, Pisa, Italy

**Corresponding Author:**

Luca Ciampi

Institute of Information Science and Technologies of the National Research Council
of Italy (ISTI-CNR), via G. Moruzzi 1 - 56124, Pisa, Italy

Tel: +39 050 6213054

Email: luca.ciampi@isti.cnr.it

**Abstract**

This paper presents a novel solution to automatically count vehicles in a parking lot using images captured by smart cameras. Unlike most of the literature on this task, which focuses on the analysis of *single* images, this paper proposes the use of multiple visual sources to monitor a wider parking area from different perspectives. The proposed multi-camera system is capable of automatically estimating the number of cars present in the *entire* parking lot directly on board the edge devices. It comprises an on-device deep learning-based detector that locates and counts the vehicles from the captured images and a decentralized geometric-based approach that can analyze the inter-camera shared areas and merge the data acquired by all the devices. We conducted the experimental evaluation on an extended version of the *CNRPark-EXT* dataset, a collection of images taken from the parking lot on the campus of the National Research Council (CNR) in Pisa, Italy. We show that our system is robust and takes advantage of the redundant information deriving from the different cameras, improving the overall performance without requiring any extra geometrical information of the monitored scene.

*Keywords:* Smart Parking, Counting Objects, Edge AI, Counting Vehicles, Smart Mobility, Deep Learning

## 1. Introduction

Traffic-related issues are constantly increasing, and tomorrow's cities cannot be considered intelligent if they do not enable smart mobility. Smart mobility applications, such as smart parking and road traffic management, are nowadays widely employed worldwide, making our cities more livable and bringing benefits to the cities and, consequently, to our lives.

Images are perhaps the best sensing modality to perceive and assess the flow of vehicles in large areas. Like no other sensing mechanism, city camera net-

9  works can monitor large areas while simultaneously providing visual data to AI
10 systems to extract relevant information from this deluge of data. However, this
11 application is often hampered by the massive flow of data that must be sent to
12 central servers or the cloud for processing. On the other hand, edge computing
13 is a recent paradigm that promotes the decentralization of data processing to
14 the border, i.e., where the data are generated, thus reducing the traffic on the
15 network and the pressure on central servers. No wonder that combination of
16 recent Computer Vision deep learning-based techniques and the edge comput-
17 ing paradigm is an emerging trend, as witnessed, for example, by Khan et al.
18 (2019) that tackles the face recognition task or by Amato et al. (2019b); Ciampi
19 et al. (2020a) that instead can detect people directly onboard surveillance cam-
20 eras. Nonetheless, this promising paradigm brings along with it also some new
21 challenges related to the limited computational resources on the disposable edge
22 devices and also concerning security inside IoT networks (Ujjan et al., 2020).

23    In this work, we tackle the problem of estimating the number of vehicles
24 present in a parking lot using images captured by smart cameras. Whereas
25 classic car counting solutions are sensor-based (e.g., entrance-level photocells,
26 per-space ground sensors), vision-based solutions provide several advantages,
27 such as a) flexibility, as cameras can adapt to more challenging configurations
28 of parking spaces (e.g., undelimited parking lots with non-fixed spaces), b) lower
29 hardware and maintenance cost, as smart cameras can cost few tens of dollars
30 while each monitoring multiple parking spaces, and c) being multi-purpose, as
31 the same hardware can be used to perform additional tasks (e.g., surveillance).
32 However, this vision-based counting task is challenging as the process of un-
33 derstanding the captured images faces many problems, such as shadows, light
34 variation, weather conditions, and inter-object occlusions. Although most of
35 the existing works concerning the vehicle counting task focus on the analysis of
36 *single* images, in many real-world scenarios, one can benefit from using multiple
37 cameras to monitor the same parking lot from different perspectives and view-
38 points. Furthermore, multiple neighboring cameras can also help cover a wider
39 area. At the same time, such an approach introduces issues related to merg-

2

ing the knowledge extracted from the single cameras with partially overlapping fields of views (FOVs), as shown in Figure 1.

In this paper, we propose a novel solution to improve car counting when scaled up with multi-camera setups. Specifically, we introduce a multi-camera system that estimates the number of cars present in the *entire* parking lot by combining a state-of-the-art Convolutional Neural Network (CNN), which can locate and count vehicles present in images belonging to individual cameras, along with a decentralized geometry-based approach that is responsible for aggregating the data gathered from all the devices. Our solution performs the task directly on the edge devices (i.e., the smart cameras) without using a central server or cloud, consequently reducing the communication overhead. The total count is built exploiting the partial results computed in parallel by the single cameras and propagated through messages. Hence, our system scales better when the number of monitored parking spaces increases. Moreover, our solution does not require any manual intervention or any extra information about the monitored parking area, such as the location of the parking spaces, nor any geometric information about the camera positions in the parking lot. In short, it is a flexible and ready-to-use solution that allows a simple "plug-and-play" insertion of new cameras into the system.

To validate our multi-camera solution, we employed the *CNRPark-EXT* dataset (Amato et al., 2017), a collection of images taken from the parking lot on the campus of the National Research Council (CNR) in Pisa, Italy. The pictures are acquired by multiple cameras having partially overlapping fields of view and describing challenging scenarios with different perspectives, illuminations, weather conditions, and many occlusions. Since the annotations of this dataset concern single images, we extended it by manually labeling a part of it to be consistent with our algorithm that instead considers the entire parking area. We conducted extensive experiments testing the generalization capabilities of the CNN-based technique responsible for detecting vehicles in single images and the effectiveness of our multi-camera algorithm, demonstrating that our system is robust and benefits from the redundant information deriving from

3

Figure 1: An example of two cameras monitoring the same parking area with partially over-lapping fields of view. This redundancy provides robustness and fault-tolerance but also raises the problem of aggregating knowledge extracted from the individual cameras.

the different cameras improving the overall performance.

To summarize, the main contributions of this work are the followings:

- We introduce a novel multi-camera system able to automatically estimate the number of cars present in the *entire* monitored parking area. It runs directly on the edge devices and combines a deep learning-based detector together with a decentralized technique that exploits the geometry of the captured images.

- We specifically extend the *CNRPark-EXT* dataset (Amato et al., 2017), a collection of images acquired by multiple cameras having partially over-lapping fields of views and describing various parking lots. We manually label a subset of it, making it suitable for our considered scenario in which we consider the whole parking area.

- We conduct an experimental evaluation showing that our system is ro-bust, flexible, and can benefit from redundant information from different

4

85 cameras while improving overall performance.

86 We organize the rest of the paper as follows. Section 2 reports other works
87 present in the literature related to our topic. Section 3 describes our multi-
88 camera counting algorithm. Section 4 states the experimental setup, describing
89 the dataset, the metrics, and the implementation details. Section 5 presents and
90 discusses the experiments and the obtained results. Finally, Section 6 concludes
91 the paper with some insights on future directions.

## 2. Related Work

93 This section overviews some works related to our, organizing them into two
94 categories. The first one concerns the counting task, while the second regards
95 multi-camera parking lot monitoring systems.

### 2.1. The counting task

97 The counting task estimates the number of object instances in still images
98 or video frames (Lempitsky & Zisserman, 2010). This topic has recently at-
99 tracted much attention due to its inter-disciplinary and widespread applicability
100 and paramount importance for many real-world applications. Examples include
101 counting bacterial cells from microscopic images (Xie et al., 2016; Ciampi et al.,
102 2022), estimating the number of people present at an event (Boominathan et al.,
103 2016; Benedetto et al., 2022), counting animals in ecological surveys to moni-
104 tor the population of a specific region (Arteta et al., 2016) and evaluating the
105 number of vehicles on a highway or in a car park (Amato et al., 2019a).

106 Several machine learning-based solutions (especially supervised) have been
107 suggested in the last years. Following the taxonomy adopted in Sindagi & Patel
108 (2018), we can broadly classify existing counting approaches into two categories:
109 counting by regression and counting by detection. Counting by *regression* is
110 a supervised method that tries to establish a direct mapping (linear or not)
111 from the image features to the number of objects present in the scene or a
112 corresponding density map (i.e., a continuous-valued function), skipping the

5

challenging task of detecting instances of the objects (Zhang et al., 2016, 2017; Oñoro-Rubio & López-Sastre, 2016; Ciampi et al., 2020b, 2021). Counting by *detection* is, instead, a supervised approach where we localize instances of the objects, and then we count them (Amato et al., 2018; Ciampi et al., 2018). While regression-based techniques work very well in very crowded scenarios where the single object instances are not well defined due to inter-class and intra-class occlusions, they perform poorly in images with a large perspective and oversized objects. Another remarkable drawback of the regression-based approaches is that they cannot precisely localize the objects present in the scene, eventually providing only a coarse position of the area in which they are distributed.

In this work, we estimate the number of vehicles present in a park area from images collected by smart cameras having large perspectives. The cars close to the cameras are much larger than those far away from them. Therefore, we employ a detection-based method. Furthermore, another reason which led us to discard counting by regression approaches is that we need to know the precise localization (with boundaries) of the detected vehicles. Most of the existing counting solutions do not directly deal with edge computing devices and the consequent constraints due to the limited available computing resources. They use deep learning-based approaches that typically require the use of a GPU and that are computationally expensive. Moreover, they consider the images as single entities. They do not account for the possible benefits of monitoring the same lots from different perspectives or covering a wider parking area with multiple cameras. Instead, our solution runs directly on the edge devices and can estimate the number of vehicles present in the entire parking lot.

### 2.2. Multi-camera parking lot monitoring

Only a few works addressed parking lot monitoring considering a multi-camera scenario. In Nieto et al. (2019), the authors applied a homography to project the detected vehicles from the plane of each camera to a common plane, where they performed a perspective correction to correct matching between the vehicle detections and the parking spots. Also, the authors in Vítek &

6

Melničuk (2017) proposed a multi-camera system to classify parking spaces as vacant or occupied. In this solution, the acquired images are processed onboard Raspberry Pi devices. The extracted information about the status of parking spaces is then transmitted to a central server, which evaluates the parking spaces in the overlapping areas. Their algorithm is based on the histogram of oriented gradients (HOG)(Dalal & Triggs, 2005) feature descriptor and support vector machine (SVM) classifier. Since the HOG feature descriptor cannot adequately describe rotated vehicles, the authors have provided a descriptor with additional information about rotation to increase the system accuracy.

However, these solutions rely on prior knowledge of the monitored scene, such as the position of the parking spaces or some geometric information concerning the parking area. For instance, the proposed system in Nieto et al. (2019) requires manually annotating the corners of the parking area and the number of spots. In essence, a preliminary annotation of the new areas and a new training phase of the algorithm are often mandatory operations. Consequently, these techniques are not very flexible. On the other hand, we propose a simple yet effective solution that does not need any extra information about the monitored scene. The smart cameras can automatically localize and count the vehicles present in their field of view, propagating the single results to the other edge devices through messages. A decentralized technique, again running directly on the edge devices, is instead in charge of analyzing and merging these results, exploiting the captured images geometry, and automatically outputs the number of cars present in the entire parking area.

## 3. Proposed approach

### 3.1. Overview

In this section, we describe our multi-camera counting algorithm. We based our system on the parallel processing of each of the smart cameras followed by the fusion of their results to estimate the number of vehicles present in the *entire* parking area.

7

<sup>172</sup> Figure 2 shows an example of our multi-camera counting system, together
<sup>173</sup> with its graphical representation. We model our system as a graph $G$, comprised
<sup>174</sup> of $n$ nodes $\nu_i$ and one Sink node $S$, $V = \{\nu_1, \nu_2, \cdots, \nu_n, S\}$. Each node $\nu_i$
<sup>175</sup> represents an independent edge device, i.e., a smart camera in our case. Two
<sup>176</sup> nodes $\nu_i$ and $\nu_j$ are considered neighbors if their FOVs overlap. In this case,
<sup>177</sup> a directed edge of the graph connects them. Each edge device $\nu_i$ can capture
<sup>178</sup> images, localize and count the vehicles present in its FOV exploiting a deep
<sup>179</sup> learning-based detector, and communicate with its neighboring nodes through
<sup>180</sup> messages $m_i$ containing the cars detections. Furthermore, each node $\nu_i$ can also
<sup>181</sup> run a local counting algorithm in charge of computing partial counting results
<sup>182</sup> concerning the estimation of the number of vehicles present in overlapped areas
<sup>183</sup> between its FOV and the ones belonging to its neighbors.

<sup>184</sup> The fusion of the partial results is performed by the Sink node $S$, which is
<sup>185</sup> also in charge of providing the final result and synchronizing all the algorithm
<sup>186</sup> steps through synchronization signals headed towards the other nodes $\nu_i$. On
<sup>187</sup> the other hand, the nodes $\nu_i$ can also communicate through messages with the
<sup>188</sup> Sink node. Messages can be of two types: i) messages $\eta_i$ containing the number
<sup>189</sup> of cars captured by the node $\nu_i$ in its FOV, and ii) messages $\mu_{j,i}$ representing
<sup>190</sup> the partial counting estimation related to the overlapping area between two
<sup>191</sup> neighboring nodes $\nu_i$ and $\nu_j$.

<sup>192</sup> In the following sections, we describe all the steps of our algorithm in detail.
<sup>193</sup> First, in Section 3.2, we outline the automatic system initialization performed by
<sup>194</sup> the smart cameras themselves, in which they compute the homographic trans-
<sup>195</sup> formations between the scene they are monitoring and the scene observed by the
<sup>196</sup> neighboring cameras. Then, in Section 3.3, we describe the CNN-based local
<sup>197</sup> counting algorithm that runs on each of the smart cameras and the geometric-
<sup>198</sup> based technique helpful for the overlapped areas. Finally, in Section 3.4, we
<sup>199</sup> depict the global counting algorithm responsible for the fusion of these individ-
<sup>200</sup> ual and partial results, and that finally outputs the number of cars present in
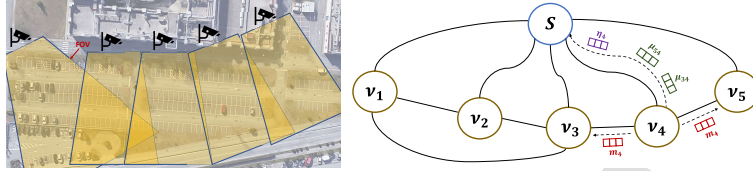<sup>201</sup> the *entire* parking area.

8

Figure 2: An example of our multi-camera counting system, with $n = 5$ smart cameras. We model it as a graph $G$, comprised of $n$ nodes $\nu_i$ (one for each camera) and one Sink node $S$, $V = \{\nu_1, \nu_2, \cdots, \nu_n, S\}$. Each node $\nu_i$ can capture images, localize and count the vehicles present in its FOV, and communicate with its neighboring nodes through messages $m_i$ containing these detections. Moreover, each node $\nu_i$ can run a local counting algorithm in charge of computing partial counting results concerning the overlapped areas between its FOV and the ones belonging to its neighbors, exploiting images geometry. These partial results are sent through messages to the Sink node $S$, which is responsible for their fusion and provides the final result. Messages to $S$ can be of two types: i) $\eta_i$ containing the number of cars captured by the node $\nu_i$ in its FOV, and ii) $\mu_{j,i}$ representing the partial counting estimation related to the overlapping area between two neighboring nodes $\nu_i$ and $\nu_j$.

## 3.2. Initialization

This step is aimed at *automatically* initializing the system, estimating the geometric relationship between each node (i.e., each scene monitored by a smart camera) and its neighbors. The only hypotheses we impose are i) each smart camera is aware of the IP addresses of its neighbors, i.e., the cameras having the field of view overlapped with its own; ii) the Sink node $S$ is aware of the IP addresses of all the smart cameras belonging to the system.

The Sink node $S$ starts the initialization phase, sending a synchronization signal to the other nodes. Once received, each smart camera captures an image of the scene it monitors and sends it to all its neighbors. Once a smart camera $i$ receives an image from a neighboring camera $j$, it computes a homographic transformation $H_{j,i}$ between the image $j$ and the image $i$ describing its monitored scene. This allows us to establish a correspondence between the points belonging to the pair of images taken by the two cameras, which will be used subsequently in the algorithm. We formalized the system initialization for a

9

<sup>217</sup> generic node $\nu_i$ in the Algorithm 1.

<sup>218</sup>    However, finding this homography can be challenging because neighboring
<sup>219</sup> cameras can have different angles of view, leading to a perspective distortion be-
<sup>220</sup> tween the images captured by them. Given a pair of neighboring nodes $\nu_i, \nu_j$, we
<sup>221</sup> employ a procedure that starts with finding the SIFT (Lowe, 2004) key-points
<sup>222</sup> and feature descriptors of the images $i, j$ captured by the two nodes. Then, we
<sup>223</sup> match the two sets of feature descriptors by performing David Lowe's ratio test
<sup>224</sup> (Lowe, 2004), and we further filter the matched feature descriptors by keeping
<sup>225</sup> only the pairs whose euclidean distance is below a given threshold. Finally, we
<sup>226</sup> obtain the homographic transformation by applying the random sample con-
<sup>227</sup> sensus (RANSAC (Fischler & Bolles, 1981)) algorithm to the filtered feature
<sup>228</sup> descriptors. All these computations are performed *automatically* without the
<sup>229</sup> need of any extra geometric information about the monitored scene, and no
<sup>230</sup> manual intervention is needed. Figure 3 shows the concatenation of two neigh-
<sup>231</sup> boring images $i$ and $j$ in which we apply the found homographic matrix to the
<sup>232</sup> image $i$, to have the same perspective as the image $j$.

---

**Algorithm 1 : Initialization**

At each Initialization Signal by $S$, each node $\nu_i$ performs the following steps:

1:  RECEIVEINITSIGNAL()                    ▷ waits the initialization signal from $S$
2:  image$_i \leftarrow$ CAMERACAPTURE()
3:  **for each** $j \in J$ **do**          ▷ $J$ is the set of neighboring nodes of node $\nu_i$
4:     SENDIMAGE(image$_i$,$\nu_j$)                    ▷ sends image$_i$ to node $\nu_j$
5:     image$_j \leftarrow$ RECEIVEIMAGE()          ▷ receives image$_j$ from node $\nu_j$
6:     $H_{j,i} =$ COMPUTEHOMOGRAPHY(image$_j$, image$_i$)

---

<sup>233</sup> *3.3. Local Counting Algorithm*

<sup>234</sup>    This section describes the local counting algorithm that runs directly on-
<sup>235</sup> board the edge devices. It combines a CNN-based counting technique in charge
<sup>236</sup> of the localization and the estimation of the number of vehicles present in the
<sup>237</sup> acquired single images, i.e., the contents of the messages $m_i$ and the quantities

10

Figure 3: Example of concatenation of two images using a homographic transformation, where it is also visible the overlapping area between them.

238    $\eta_i$ shown in Figure 2, together with a geometric-based approach responsible of
239    estimating the number of vehicles present in the overlapping areas between the
240    nodes and their neighbors, i.e., the quantities $\mu_{j,i}$.

241    *A vehicle counting CNN on the Edge.* Each smart camera needs to indepen-
242    dently detect and count vehicles from its captured frame. For this step, every
243    approach providing precise localization of the detected vehicles in the pixel
244    space is suitable, and the choice of a particular approach should be guided by
245    resource constraints, e.g., available memory, prediction frequency, or energy con-
246    sumption, if any. Here, we base our vehicle counting technique on *Mask R-CNN*
247    (He et al., 2017), a popular deep CNN for instance segmentation that operates
248    within the 'recognition using regions paradigm' (Gu et al., 2009). In particular,
249    it extends the *Faster R-CNN* detector (Ren et al., 2017) by adding a branch
250    that outputs a binary mask saying whether or not a given pixel is part of an
251    object. Briefly, a CNN acts a backbone in the first stage, extracting the input
252    image features. Starting from this feature space, another CNN named Region
253    Proposal Network (RPN) generates region proposals that might contain objects.
254    RPN slices pre-defined region boxes (called anchors) over this space and ranks
255    them, suggesting those most likely containing objects. Once RPN produces the
256    Regions Of Interests (ROIs), they might be of different sizes. Since it is hard
257    to work on features having different sizes, RPN reduces them into the same di-
258    mension using the Region of Interest Pooling algorithm. Finally, these fixed-size

11

proposals are processed by two parallel CNN-based branches: one is responsible for classifying and localizing the objects inside them with bounding boxes; the second produces a binary mask that says whether or not a given pixel is part of an object. In the end, given an input image, the network produces per-pixel masks localizing the detected objects together with the associated labels classifying them.

To make our counting solution able to run efficiently directly on the edge devices, we employ, as a backbone, the *ResNet50* architecture, a lighter version of the popular *ResNet101* (He et al., 2016). This simplification is also justified because the more powerful version of Mask R-CNN based on the ResNet101 model was designed for more complicated visual detection tasks than ours. Originally, Mask R-CNN was trained on the *COCO* dataset (Lin et al., 2014) to detect and recognize 80 different classes of everyday objects. In our case, we have to localize and identify objects belonging to just one category (i.e., the *vehicle* category). To this end, we further simplify the model by reducing the number of the final fully convolutional layers responsible for the classification of the detected objects, making the model lighter. Once we have localized the instances of the objects, we count them estimating the number of vehicles present in the scene.

*Local counting.* The Sink node $S$ starts this phase, sending a synchronization signal to all the smart cameras belonging to the system. Once received the synchronization signal, each node $\nu_i$ captures an image belonging to its underlying FOV and feeds it to the previously described CNN-based counting technique obtaining a set of masks masks$_i$ localizing the vehicles present in the scene. The cardinality of this set of masks corresponds to the number of cars present in the image, i.e., the quantity $\eta_i$, that is sent with a message to the Sink node $S$. Then, the node $\nu_i$ packs this set of masks masks$_i$ in a message $m_i$, sends it to all its neighboring nodes $\nu_j$, and receives from them their corresponding set of masks masks$_j$ packed in a message $m_j$. Once received a message $m_j$, the node $\nu_i$ is responsible for analyzing the potential vehicles present in the overlapped

12

289 area between its FOV and the one of the node $\nu_j$. To this end, it employs the
290 homographic transformation $H_{j,i}$ computed during the system initialization, as
291 described in Section 3.2. Specifically, it projects the masks belonging to the set
292 masks$_j$ into its image plane, filtering them and discarding the ones that overlap
293 with the masks belonging to the set masks$_i$ having a value of Intersection over
294 Union (IoU) greater than a threshold that we empirically found to be optimal
295 at 0.2. These masks indeed localize vehicles already detected, which should not
296 be considered a second time. On the other hand, the cars left after this filtering
297 are vehicles that were not detected in the FOV underlying the node $\nu_i$, but
298 instead found by the node $\nu_j$, probably because of having a better view of this
299 object. Referring to our graph modeling the system and reported in Figure 2,
300 the number of the discarded cars after this filtering operation corresponds to
301 the message $\mu_{j,i}$, that is sent to the Sink node $S$. We detail all the described
302 steps in the Algorithm 2 and in the Procedure 3.

---

**Algorithm 2 : Local Counting**

At each Computational Signal by $S$, each node $\nu_i$ performs the following steps:

1: RECEIVECOMPUTSIGNAL() ▷ waits the computational signal from $S$

2: image$_i \leftarrow$ CAMERACAPTURE()

3: masks$_i \leftarrow$ MASKRCNN(image$_i$)

4: $\eta_i \leftarrow |$masks$_i|$

5: SENDMESSAGE($\eta_i, S$) ▷ sends $\eta_i$ to Sink node $S$

6: $m_i \leftarrow$ PACKMESSAGE(masks$_i$) ▷ builds message $m_i$ containing masks$_i$

7: **for each** $j \in J$ **do** ▷ $J$ is the set of neighboring nodes of node $\nu_i$

8:     SENDMESSAGE($m_i, \nu_j$) ▷ sends $m_i$ to node $\nu_j$

9:     $m_j \leftarrow$ RECEIVEMESSAGE() ▷ receives message $m_j$ from node $\nu_j$

10:     masks$_j \leftarrow$ UNPACKMESSAGE($m_j$) ▷ unpacks $m_j$ containing masks$_j$

11:     $\mu_{j,i} \leftarrow$ COMPUTE_$\mu$(masks$_i$, masks$_j$, $H_{j,i}$)

12:     SENDMESSAGE($\mu_{j,i}, S$) ▷ sends $\mu_{j,i}$ to Sink node $S$

---

13

---

**Algorithm 3** : **Computation of** $\mu$

$\mu$ represents the num of cars detected by $\nu_j$ and already detected by $\nu_i$

Each node $\nu_i$ performs the following procedure:

---

1: **procedure** COMPUTE_$\mu$(masks$_i$, masks$_j$, $H_{j,i}$)

2:　　n_cars_already_detected $\leftarrow$ 0

3:　　**for each** mask $\in$ masks$_j$ **do**

4:　　　　mask$_h$ $\leftarrow$ PROJECT($H_{j,i}$, mask)　　▷ projects mask points on plane $i$

5:　　　　**if** mask$_h$ falls within image$_i$ **then**

6:　　　　　　mask$_{max}$ $\leftarrow$ arg max$_{m \in masks_i}$ IoU(mask$_h$, $m$)

7:　　　　　　**if** IoU(mask$_h$, mask$_{max}$) $> \tau$ **then**

8:　　　　　　　　n_cars_already_detected $++$

9:　　**return** n_cars_already_detected

---

### 3.4. Global Counting Algorithm

In this section, we describe the global counting algorithm that runs on the Sink node $S$, responsible for the fusion of the partial results coming from all the other nodes, and that finally outputs the number of cars present in the *entire* monitored parking area.

This phase starts when $S$ receives all the $\eta_i$ and the $\mu_{j,i}$ messages, i.e., the number of vehicles estimated in the single FOVs and the estimation of the number of cars already considered in the overlapping areas between neighboring cameras, from all the nodes belonging to the system. Specifically, for each overlapped area shared between a pair of nodes $\nu_i, \nu_j$, the node $S$ receives two messages $\mu_{j,i}$ and $\mu_{i,j}$, the contents of which are computed by the two nodes employing two homographic transformations $H_{j,i}$ and $H_{i,j}$, respectively. These two quantities can be potentially different. We choose the best value by aggregating them, choosing between three different functions - max, min and mean, finding that the latter is the best one. Finally, the node $S$ builds the final result, i.e., the estimation of the number of vehicles present in the *entire* parking lot, by summing up the content of all the $\eta_i$ messages and subtracting the computed aggregated values. We detail all these steps in the Algorithm 4.

14

---

**Algorithm 4 : Global Counting**

The Sink node $S$ performs the following steps:

1: **for each** $(\mu_{i,j}, \mu_{j,i})$ **do**

2: $\quad \overline{\mu_k} \leftarrow \text{AGGREGATE}(\mu_{i,j}, \mu_{j,i})$

3: global_cars_count $\leftarrow \sum_{n=1}^{N} \eta_n - \sum_{k=1}^{K} \overline{\mu_k}$

$\quad \triangleright N$ is the set of nodes, $K$ is the set of aggregations

---

## 4. Experimental Setup

In this section, we describe the simulated scenario that we exploited for our experiments. In particular, we extended the *CNRPark-EXT* dataset (Amato et al., 2017), adapting it to be suitable for the counting task so that it was usable for training the vehicles counting CNN running on the smart cameras and applicable to validate our multi-camera algorithm. Furthermore, we briefly describe the *PKLot* dataset (de Almeida et al., 2015), a public dataset comprising parking lot scenes that we exploited for further assessing the generalization capabilities of the local vehicles counting network. Then, we illustrate the employed evaluation metrics, and, finally, we report some implementation details.

### 4.1. The CNRPark-EXT Dataset

In this work, we exploit the *CNRPark-EXT* public dataset introduced in Amato et al. (2017), a collection of annotated images of vacant and occupied parking spaces on the campus of the National Research Council (CNR) in Pisa, Italy. This dataset represents most of the challenging situations that can be found in a real scenario: nine different cameras capture the images under various weather conditions, angles of view, light conditions, and many occlusions. Furthermore, the cameras have their fields of view partially overlapped. Since this dataset is specifically designed for parking lot occupancy detection, it is not directly usable for the counting task. Indeed, each image, called *patch*, contains one parking space labeled according to its occupancy status - 0 for vacant and 1 for occupied. Since this work aims at counting the cars present in the parking

15

343 area, we extended it by considering the full images and adapting the ground
344 truth to our purposes.

345 Specifically, we created a suitable label set to train and evaluate the local ve-
346 hicles counting based on Mask R-CNN. In this case, labels correspond to *binary*
347 *masks*, i.e., binary images identifying the polygons surrounding the vehicles we
348 want to detect. Since mask creation is a very time-consuming operation, dif-
349 ferently from our previous work (Ciampi et al., 2018), we considered the *raw*
350 masks obtained directly from the bounding boxes localizing the occupied park-
351 ing spaces. The idea is that we do not need precise polygons that identify the
352 vehicles we want to detect. Still, we can use the region within the delimiters
353 that identify the occupied parking spaces and the underlying part of the car.

354 On the other hand, to validate our multi-camera algorithm, we built a simu-
355 lated scenario considering some sequences of images belonging to different cam-
356 eras captured simultaneously. In other words, a sequence is defined as the set of
357 images captured by the different smart cameras that are monitoring the parking
358 area at the same moment. Hence, a sequence represents a snapshot of the *entire*
359 parking lot at a given timestamp, and it takes into account all the spaces from
360 the available different views. We manually annotated these sequences to obtain
361 the ground truth car counts. Specifically, we considered the single images com-
362 posing a sequence, counting the vehicles present in the scenes, but taking care of
363 accounting for them just once if they appear in more than one view, i.e., discard-
364 ing the cars from the global count if they were located in the overlapping areas.
365 We labeled six different sequences, two for each weather condition, considering
366 the images belonging from $camera_2$ to $camera_9$. We did not consider $camera_1$
367 since it has small and particularly skewed field-of-view overlaps with the other
368 cameras, hindering the automatic homography estimation and the subsequent
369 projections.

370 *4.2. The PKLot Dataset*

371 To further validate the generalization capabilities of the CNN-based local
372 vehicles counting algorithm, we exploited an additional public dataset, named

16

<sup>373</sup> *PKLot* (de Almeida et al., 2015). In particular, this dataset is composed by
<sup>374</sup> three different scenarios describing three different parking lot scenes - *UFPR04*,
<sup>375</sup> *UFPR05* and *PUC*. We considered only the first two subsets since the third one
<sup>376</sup> contains images captured from a fixed camera located at the height of the 10th
<sup>377</sup> floor of a building, which provides a slanted view of the parking lot and results
<sup>378</sup> in a different setting without intra-vehicle occlusions. Since also the *PKLot*
<sup>379</sup> dataset, like the *CNRPark-EXT* one, is specifically designed for the parking
<sup>380</sup> lot occupancy detection task, we manually re-labeled the ground truth for our
<sup>381</sup> purposes as already described in Section 4.1, obtaining a simulation scenario
<sup>382</sup> suitable for measure the performance of our solution for the counting task.

<sup>383</sup> *4.3. Evaluation Metrics*

<sup>384</sup> Following other counting benchmarks, we exploited Mean Absolute Error
<sup>385</sup> (*MAE*), Mean Square Error (*MSE*), and Mean Relative Error (*MRE*) as the
<sup>386</sup> metrics for the performance evaluation, defined as follows:

$$MAE = \frac{1}{N} \sum_{n=1}^{N} |c_n^{gt} - c_n^{pred}|, \tag{1}$$

$$MSE = \frac{1}{N} \sum_{n=1}^{N} (c_n^{gt} - c_n^{pred})^2, \tag{2}$$

$$MRE = \frac{1}{N} \sum_{n=1}^{N} \frac{|c_n^{gt} - c_n^{pred}|}{\text{num\_spaces}_n}, \tag{3}$$

<sup>387</sup> where $N$ is the total number of the images, $c_{gt}$, $c_{pred}$ and $num\_spaces_n$ are
<sup>388</sup> the actual count, the predicted count, and the total number of parking spaces
<sup>389</sup> of the n-th image, respectively. Note that as a result of the squaring of each
<sup>390</sup> difference, MSE effectively penalizes large errors more heavily than small ones
<sup>391</sup> and thus should be more useful when large errors are particularly undesirable.
<sup>392</sup> On the other hand, MRE also considers the relation between the error and the
<sup>393</sup> total number of objects present in the image.

*4.4. Implementation Details*

We report in this section some implementation details concerning the Mask R-CNN-based algorithm responsible for the prediction of the number of vehicles in the single images. In particular, we trained the modified Mask R-CNN initializing the weights of the ResNet50 backbone with the ones of a pre-trained model on *ImageNet* (Deng et al., 2009), a popular dataset for classification tasks, and the remaining ones at random. We froze the backbone for the first 10 epochs, and then we trained the whole network for 20 additional epochs. We used Stochastic Gradient Descent (SGD) to perform the CNN parameters update. Concerning the Region Proposal Network, explained in Section 3.3, we exploited a set of five anchors of sizes 16, 32, 64, 128, and 256 pixels. To prevent overfitting, we applied some standard augmentation techniques to the training data: images are horizontally flipped with a 0.5 probability, then their pixels are multiplied by a random value between 0.8 and 1.5, and finally, they are blurred using a Gaussian kernel with a standard deviation of a random value between 0 and 5. Then, to support training multiple images per batch, we resized all pictures to the same size. If an image was not square, we padded it with zeros to preserve the aspect ratio. In the end, we obtained images of size $1024 \times 1024$. At inference time, images were resized and padded with zeros to get a square picture of size $1024 \times 1024$, and no other augmentations took place.

## 5. Experiments and Results

In this section, we report the experiments and the obtained results. First, we evaluate the performance against other state-of-the-art solutions of the CNN-based technique responsible for estimating the vehicles in the single images directly onboard the smart cameras, also stressing its generalization capabilities. Then, we validate the effectiveness of our multi-camera algorithm by testing it in the simulated scenario previously described. We demonstrate that our system can benefit from the redundant information deriving from the different cameras, obtaining performance improvements in all the considered counting metrics.

18

### 5.1. Experiments on the CNN-based counting solution on the edge

### 5.1.1. State-of-the-art comparison

We compared our solution with the results obtained in our previous work Ciampi et al. (2018), where we presented a centralized counting approach based on the original version of Mask R-CNN having the ResNet101 model as a features extractor, which has been fine-tuned on a very small manually annotated subset of the CNRPark-EXT dataset, starting from the model pre-trained on the *COCO* dataset (Lin et al., 2014). We filtered the detections considering only the predictions related to the car class, and we counted them. Although this solution is very computationally expensive and unsuitable for edge devices, it represents a direct comparison in terms of counting on the same dataset. We also compared our technique against the method proposed in Amato et al. (2017), an approach for car parking occupancy detection based on *mAlexNet*, a deep CNN designed explicitly for smart cameras. This work represents an indirect method for counting cars in a parking lot, as the counting problem is cast as a classification problem: if a parking space is occupied, we increment the total number of cars; otherwise, we do not. We illustrate the results in Table 1, where we also report the performance obtained using the Mask R-CNN network without a preliminary fine-tuning on the CNRPark-EXT dataset. Our solution performs better than the other considered methods, considering all three counting metrics. In particular, our approach outperforms the solution introduced in Ciampi et al. (2018), despite the latter employing a more deep and powerful CNN, and it is designed to be used as a centralized-server solution. This is explained by the fact that in Ciampi et al. (2018) the authors fine-tuned the CNN using a tiny dataset. Consequently, the algorithm overfits on the training data, and it cannot generalize over the test subset. It is also worthy of notice that our CNN also outperforms the mAlexNet network, even though the latter knows the exact location of the parking spaces. Figure 4 shows some examples of images belonging to different cameras and different weather conditions together with the masks localizing them computed by our counting solution.

19

| Method | CNRPark-EXT | | | PKLot | | |
|---|---|---|---|---|---|---|
| | MAE | MSE | MRE | MAE | MSE | MRE |
| (Amato et al., 2017) | 1.34 | 8.00 | 0.04 | | - | |
| (Ciampi et al., 2018) | 1.05 | 4.41 | 0.03 | | - | |
| ResNet50 Mask R-CNN | 11.20 | 247.40 | 0.30 | 16.90 | 522.40 | 0.48 |
| Our solution | **0.49** | **1.04** | **0.01** | **4.56** | **33.88** | **0.13** |

Table 1: Local Counting: Left-side: results obtained using our counting solution on the edge compared with other state-of-the-art approaches; we get the best results on all the three considered counting metrics. Right-side: evaluation of the generalization capabilities on the *PKLot* dataset (de Almeida et al., 2015), using the model trained on the *CNRPark-EXT* dataset; we achieved an error that is approximately four times lower than the one obtained with the COCO pre-trained model.



(a) Image from Camera$_2$

(b) Image from Camera$_8$

Figure 4: Two examples of the output of our counting method. Images are taken from the CNRPark-EXT dataset. We report the predictions and the estimate of the number of vehicles present in the scene.

### 5.1.2. Generalization capabilities

Errors in vehicle detection and counting are due to many reasons, but critical points are different light conditions and diverse perspectives. Weather conditions might produce significant illumination changes since puddles and wet floors create a textural pattern that may lead to an error, and sunbeams can create reflections on the car windscreen, covering the majority of the images with saturated patterns. When a CNN does not generalize well, it works well only in the conditions where it was trained.

To measure the robustness of our approach to these scenarios, we performed two types of experiments exploiting the *CNRPark-EXT* dataset: i) *inter-weather* and ii) *inter-camera* experiments. In the former, we trained our CNN with images taken in one particular weather condition, and we computed the performance metrics obtained on images having different weather conditions. In particular, we performed three experiments, training respectively on the *Sunny*, *Overcast* and *Rainy* subsets of the CNRPark-EXT dataset. In the latter, we trained our algorithm employing images from one camera, and then we computed the performance metrics on pictures captured by another camera. In particular, we performed two experiments, training with images coming respectively from $camera_1$ and $camera_8$. We chose these two cameras because they are particularly representative since one has a side view of the parking lot while the other has a pure front view.

We report the results of the two experiments in Table 2 and Table 3, respectively. We achieve a good generalization in both the considered scenarios. We experienced a larger amount of error when the CNN was trained and tested on two opposite weather conditions, for instance, *Sunny* and *Rainy*, while the more accurate model was the one trained on *Overcast* weather conditions. However, the performance difference is quite small. On the other hand, in *inter-camera* experiments, the model trained on $camera_8$ is the best, and it has a slight drop in performance only when tested on the $camera_1$ subset. The model trained on the $camera_1$ dataset performs in general worse. This is probably due to a bias

21

| Train Set | Sunny | | | Overcast | | | Rainy | | |
|---|---|---|---|---|---|---|---|---|---|
| | MAE | MSE | MRE | MAE | MSE | MRE | MAE | MSE | MRE |
| Sunny | - | - | - | 0.29 | 0.34 | 0.009 | 0.96 | 2.78 | 0.02 |
| Overcast | 0.62 | 1.09 | 0.02 | - | - | - | 0.56 | 1.26 | 0.01 |
| Rainy | 0.84 | 1.65 | 0.02 | 0.49 | 0.65 | 0.01 | - | - | - |

Table 2: CNRPark-EXT: Results of inter-weather experiments in terms of counting metrics obtained when training on sunny, overcast, or rainy weather.

| Metric | Train Set | Test Set | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | C1 | C2 | C3 | C4 | C5 | C6 | C7 | C8 | C9 |
| MAE | C1 | - | 0.77 | 1.21 | 2.53 | 3.26 | 2.57 | 2.88 | 2.88 | 1.54 |
| | C8 | 3.87 | 0.85 | 0.76 | 0.45 | 0.48 | 0.71 | 1.07 | - | 0.41 |
| MRE | C1 | - | 0.08 | 0.05 | 0.06 | 0.07 | 0.05 | 0.06 | 0.05 | 0.05 |
| | C8 | 0.11 | 0.09 | 0.03 | 0.01 | 0.01 | 0.01 | 0.02 | - | 0.01 |
| MSE | C1 | - | 1.48 | 2.91 | 10.61 | 20.24 | 13.50 | 19.82 | 17.30 | 7.19 |
| | C8 | 22.60 | 1.78 | 1.36 | 0.57 | 0.74 | 0.95 | 4.97 | - | 2.13 |

Table 3: CNRPark-EXT: Results of inter-camera experiments in terms of counting metrics obtained when training on camera 1 and camera 8.

⁴⁸³ in the CNRPark-EXT dataset, where the majority of the images are captured
⁴⁸⁴ from a frontal viewpoint.

⁴⁸⁵ Moreover, to further validate the generalization capabilities of our approach,
⁴⁸⁶ we considered our counting network trained on the entire training set of the
⁴⁸⁷ *CNRPark-EXT* dataset, and we tested it over a different dataset, the *PKLot*
⁴⁸⁸ dataset (de Almeida et al., 2015). Results are shown in Table 1 where we also
⁴⁸⁹ report the performance obtained using the Mask R-CNN network without a
⁴⁹⁰ preliminary fine-tuning on the *CNRPark-EXT* dataset. As we can see, using
⁴⁹¹ our solution, we achieve an error that is approximately four times lower than
⁴⁹² the one obtained with the COCO pre-trained model.

22

### 5.2. Experiments on the Multi-Camera Scenario

To the best of our knowledge, there are no annotated datasets in the literature suitable for evaluating counting algorithms operating on multiple FOV-overlapping cameras. The most relevant work in this context is Nieto et al. (2019), in which there are only two overlapping cameras facing each other with an extreme perspective transformation between the two; this makes any automatic perspective computation nearly impossible without manual intervention, and this is a mandatory assumption for our proposed method. Hence, we performed our experiments on the extended version of the CNRPark-EXT dataset created on purpose in this work, which we hope will become a new benchmark for this task. Furthermore, to demonstrate that our algorithm can benefit from the redundant information deriving from the different cameras, we compared the obtained results against a baseline and a simplified version of our algorithm.

Specifically, we compared our solution against a system that is not aware of the other cameras' overlapped areas, and so it just sums up all the vehicles detected by all the cameras belonging to a sequence (Naïve Counting $\mathbf{N}$). Then, we considered a more conservative approach, where the nodes employ the homographic transformations only with the purpose of black-masking the overlapped areas (Overlap Masking $\mathbf{M}$). This latter baseline then loses the ability to take advantage of monitoring the same from different views. However, it is still aware of the locations of the overlapping areas, and it considers the vehicles inside them only once.

Results are shown in Table 4. Our solution obtains the best results compared to the considered baselines in all the three counting metrics and all the employed scenarios. We report the errors concerning the considered six sequences of the CNRPark-EXT dataset, together with the MAE, MSE, and MRE, which summarize 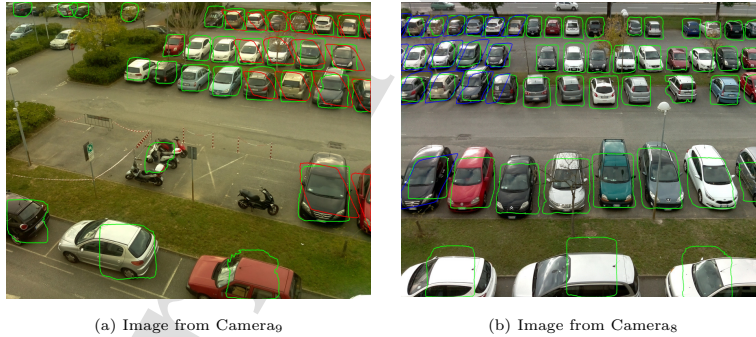the mean results regarding all the scenarios. As an example, in Figure 5 we also report the output of our multi-camera algorithm for a pair of images belonging to two different cameras having a shared area in their field of view, where we highlight in red and blue the masks projected from one camera to the other, using the previously computed homographic transformations.

23

| | Error | | | Absolute Err. | | | Squared Err. | | | Relative Err. (%) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | **N** | **M** | **O** | **N** | **M** | **O** | **N** | **M** | **O** | **N** | **M** | **O** |
| Overcast-1 | 124 | -33 | **2** | 124 | 33 | **2** | 15,376 | 1,089 | **4** | 71.6 | 19.0 | **1.2** |
| Overcast-2 | 131 | -26 | **1** | 131 | 26 | **1** | 17,161 | 676 | **1** | 76.1 | 15.1 | **0.6** |
| Rainy-1 | 80 | -39 | **-5** | 80 | 39 | **5** | 6,400 | 1,521 | **25** | 47.6 | 23.2 | **2.9** |
| Rainy-2 | 105 | -44 | **-5** | 105 | 44 | **5** | 11,025 | 1,936 | **25** | 54.4 | 22.8 | **2.6** |
| Sunny-1 | 117 | -38 | **2** | 117 | 38 | **2** | 13,689 | 1,444 | **4** | 68.0 | 22.1 | **1.2** |
| Sunny-2 | 113 | -37 | **2** | 113 | 38 | **2** | 12,769 | 1,444 | **4** | 66.1 | 22.2 | **1.2** |
| Mean | 111.6 | -36.1 | **-0.5** | 111.6 | 36.3 | **2.8** | 12,736.6 | 1,351.6 | **10.5** | 63.9 | 20.7 | **1.6** |

**N**: Naïve Counting; **M**: Overlap Masking; **O**: Ours (mean aggr., IoU Threshold $\tau = 0.2$)

Table 4: Results using our multi-camera counting algorithm, considering the *entire* parking lot. We compare our solution against a baseline and a simplified version of our algorithm. We report the errors obtained on the six considered sequences (two for each weather condition) of the CNRPark-EXT dataset that we extend on purpose.



(a) Image from Camera$_9$

(b) Image from Camera$_8$

Figure 5: Example of the output of our multi-camera algorithm for a pair of images belonging to two different cameras $i, j$ having a shared area in their FOV. We report in green the masks localizing the vehicles detected by a camera in its own FOV, while in red and blue, the masks projected from camera j to camera i and vice-versa, employing the homographic transformations pre-computed during the system initialization.

## 6. Conclusion

This paper presented a distributed artificial intelligence-based system that automatically counts the vehicles present in a parking lot using images taken by multiple smart cameras. Unlike most of the works in literature, we introduced a multi-camera approach that can estimate the number of cars present in the *entire* parking area and not only in the single captured images. The main peculiarities of this approach are that all the computation is performed in a distributed manner at the edge of the network and that there is no need for any extra information about the monitored parking area, such as the location of the parking spaces, nor any geometric information about the position of the cameras in the parking lot. We modeled our system as a graph. The nodes, i.e., the smart cameras, are responsible for estimating the number of cars present in their view and merging data from nearby devices with an overlapping field of view. Our solution is simple but effective, combining a deep-learning technique with a distributed geometry-based approach. We evaluated our algorithm on the CNRPark-EXT dataset, which we specifically extended and which we hope will become a new benchmark for counting vehicles in multi-camera parking area scenarios. Through an experimental evaluation, we showed how we benefit from redundant information from different cameras while improving overall performance.

There are multiple lines of future development that can help improve the proposed system. Although our multi-camera algorithm is flexible, one limitation relies on computing the homographic matrix between images captured by cameras placed in completely different locations, such as facing each other. In such cases, the two perspectives are totally different, and manual intervention is required to avoid the generation of an inaccurate geometric transformation.

25

FSE 2014-2020 AI-MAP (CNR4C program, CUP B15J19001040004).

**References**

de Almeida, P. R., Oliveira, L. S., Britto, A. S., Silva, E. J., & Koerich, A. L. (2015). PKLot – a robust dataset for parking lot classification. *Expert Systems with Applications*, *42*, 4937–4949. URL: https://doi.org/10. 1016%2Fj.eswa.2015.02.009. doi:10.1016/j.eswa.2015.02.009.

Amato, G., Bolettieri, P., Moroni, D., Carrara, F., Ciampi, L., Pieri, G., Gennaro, C., Leone, G. R., & Vairo, C. (2018). A wireless smart camera network for parking monitoring. In *2018 IEEE Globecom Workshops (GC Wkshps)*. IEEE. URL: https://doi.org/10.1109%2Fglocomw.2018. 8644226. doi:10.1109/glocomw.2018.8644226.

Amato, G., Carrara, F., Falchi, F., Gennaro, C., Meghini, C., & Vairo, C. (2017). Deep learning for decentralized parking lot occupancy detection. *Expert Systems with Applications*, *72*, 327–334. URL: https://doi.org/ 10.1016%2Fj.eswa.2016.10.055. doi:10.1016/j.eswa.2016.10.055.

Amato, G., Ciampi, L., Falchi, F., & Gennaro, C. (2019a). Counting vehicles with deep learning in onboard UAV imagery. In *2019 IEEE Symposium on Computers and Communications (ISCC)*. IEEE. URL: https://doi.org/10.1109%2Fiscc47284.2019.8969620. doi:10. 1109/iscc47284.2019.8969620.

Amato, G., Ciampi, L., Falchi, F., Gennaro, C., & Messina, N. (2019b). Learning pedestrian detection from virtual worlds. In *Lecture Notes in Computer Science* (pp. 302–312). Springer International Publishing. URL: https://doi.org/10.1007%2F978-3-030-30642-7_27. doi:10.1007/978-3-030-30642-7_27.

Arteta, C., Lempitsky, V., & Zisserman, A. (2016). Counting in the wild. In *Computer Vision – ECCV 2016* (pp. 483–498). Springer International

26

580 Publishing. URL: https://doi.org/10.1007%2F978-3-319-46478-7_30.
581 doi:10.1007/978-3-319-46478-7_30.

582 Benedetto, M. D., Carrara, F., Ciampi, L., Falchi, F., Gennaro, C., & Am-
583 ato, G. (2022). An embedded toolset for human activity monitoring in
584 critical environments. *Expert Systems with Applications*, *199*, 117125.
585 URL: https://doi.org/10.1016%2Fj.eswa.2022.117125. doi:10.1016/
586 j.eswa.2022.117125.

587 Boominathan, L., Kruthiventi, S. S. S., & Babu, R. V. (2016). Crowd-
588 Net. In *Proceedings of the 24th ACM international conference on Mul-
589 timedia*. ACM. URL: https://doi.org/10.1145%2F2964284.2967300.
590 doi:10.1145/2964284.2967300.

591 Ciampi, L., Amato, G., Falchi, F., Gennaro, C., & Rabitti, F. (2018). Counting
592 vehicles with cameras. In S. Bergamaschi, T. D. Noia, & A. Maurino (Eds.),
593 *Proceedings of the 26th Italian Symposium on Advanced Database Systems,
594 Castellaneta Marina (Taranto), Italy, June 24-27, 2018*. CEUR-WS.org
595 volume 2161 of *CEUR Workshop Proceedings*. URL: http://ceur-ws.
596 org/Vol-2161/paper12.pdf.

597 Ciampi, L., Carrara, F., Amato G., & Gennaro, C. (2022). Counting or localiz-
598 ing? evaluating cell counting and detection in microscopy images. In *Pro-
599 ceedings of the 17th International Joint Conference on Computer Vision,
600 Imaging and Computer Graphics Theory and Applications*. SCITEPRESS
601 - Science and Technology Publications. URL: https://doi.org/10.5220%
602 2F0010923000003124. doi:10.5220/0010923000003124.

603 Ciampi, L., Messina, N., Falchi, F., Gennaro, C., & Amato, G. (2020a). Virtual
604 to real adaptation of pedestrian detectors. *Sensors*, *20*, 5250. URL: https:
605 //doi.org/10.3390%2Fs20185250. doi:10.3390/s20185250.

606 Ciampi, L., Santiago, C., Costeira, J., Gennaro, C., & Amato, G. (2021). Do-
607 main adaptation for traffic density estimation. In *Proceedings of the 16th In-
608 ternational Joint Conference on Computer Vision, Imaging and Computer*

609 *Graphics Theory and Applications*. SCITEPRESS - Science and Technology

610 Publications. URL: `https://doi.org/10.5220%2F0010303401850195`.

611 doi:10.5220/0010303401850195.

612 Ciampi, L., Santiago, C., Costeira, J. P., Gennaro, C., & Amato, G. (2020b).

613 Unsupervised vehicle counting via multiple camera domain adaptation. In

614 A. Saffiotti, L. Serafini, & P. Lukowicz (Eds.), *Proceedings of the First*

615 *International Workshop on New Foundations for Human-Centered AI (Ne-*

616 *HuAI) co-located with 24th European Conference on Artificial Intelligence*

617 *(ECAI 2020), Santiago de Compostella, Spain, September 4, 2020* (pp. 82–

618 85). CEUR-WS.org volume 2659 of *CEUR Workshop Proceedings*. URL:

619 `http://ceur-ws.org/Vol-2659/ciampi.pdf`.

620 Dalal, N., & Triggs, B. (2005). Histograms of oriented gradients for human

621 detection. In *2005 IEEE Computer Society Conference on Computer Vision*

622 *and Pattern Recognition (CVPR'05)*. IEEE. URL: `https://doi.org/10.`

623 `1109%2Fcvpr.2005.177`. doi:10.1109/cvpr.2005.177.

624 Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., & Fei-Fei, L. (2009). ImageNet:

625 A large-scale hierarchical image database. In *2009 IEEE Conference on*

626 *Computer Vision and Pattern Recognition*. IEEE. URL: `https://doi.`

627 `org/10.1109%2Fcvpr.2009.5206848`. doi:10.1109/cvpr.2009.5206848.

628 Fischler, M. A., & Bolles, R. C. (1981). Random sample consensus. *Com-*

629 *munications of the ACM*, *24*, 381–395. URL: `https://doi.org/10.1145%`

630 `2F358669.358692`. doi:10.1145/358669.358692.

631 Gu, C., Lim, J. J., Arbelaez, P., & Malik, J. (2009). Recognition using regions.

632 In *2009 IEEE Conference on Computer Vision and Pattern Recognition*.

633 IEEE. URL: `https://doi.org/10.1109%2Fcvpr.2009.5206727`. doi:10.

634 `1109/cvpr.2009.5206727`.

635 He, K., Gkioxari, G., Dollar, P., & Girshick, R. (2017). Mask r-CNN. In

636 *2017 IEEE International Conference on Computer Vision (ICCV)*. IEEE.

637   URL: https://doi.org/10.1109%2Ficcv.2017.322. doi:10.1109/iccv.
638   2017.322.

639   He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image
640   recognition. In *2016 IEEE Conference on Computer Vision and Pattern*
641   *Recognition (CVPR)*. IEEE. URL: https://doi.org/10.1109%2Fcvpr.
642   2016.90. doi:10.1109/cvpr.2016.90.

643   Khan, M. Z., Harous, S., Hassan, S. U., Khan, M. U. G., Iqbal, R.,
644   & Mumtaz, S. (2019). Deep unified model for face recognition
645   based on convolution neural network and edge computing. *IEEE Ac-*
646   *cess*, *7*, 72622–72633. URL: https://doi.org/10.1109%2Faccess.2019.
647   2918275. doi:10.1109/access.2019.2918275.

648   Lempitsky, V. S., & Zisserman, A. (2010). Learning to count objects in
649   images. In J. D. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R. S.
650   Zemel, & A. Culotta (Eds.), *Advances in Neural Information Process-*
651   *ing Systems 23: 24th Annual Conference on Neural Information Pro-*
652   *cessing Systems 2010. Proceedings of a meeting held 6-9 December 2010,*
653   *Vancouver, British Columbia, Canada* (pp. 1324–1332). Curran Asso-
654   ciates, Inc. URL: https://proceedings.neurips.cc/paper/2010/hash/
655   fe73f687e5bc5280214e0486b273a5f9-Abstract.html.

656   Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár,
657   P., & Zitnick, C. L. (2014). Microsoft COCO: Common objects in context.
658   In *Computer Vision – ECCV 2014* (pp. 740–755). Springer International
659   Publishing. URL: https://doi.org/10.1007%2F978-3-319-10602-1_48.
660   doi:10.1007/978-3-319-10602-1_48.

661   Lowe, D. G. (2004). Distinctive image features from scale-invariant key-
662   points. *International Journal of Computer Vision*, *60*, 91–110. URL:
663   https://doi.org/10.1023%2Fb%3Avisi.0000029664.99615.94. doi:10.
664   1023/b:visi.0000029664.99615.94.

29

Nieto, R. M., Garcia-Martin, A., Hauptmann, A. G., & Martinez, J. M. (2019). Automatic vacant parking places management system using multicamera vehicle detection. *IEEE Transactions on Intelligent Transportation Systems*, *20*, 1069–1080. URL: `https://doi.org/10.1109%2Ftits.2018.2838128`. doi:`10.1109/tits.2018.2838128`.

Oñoro-Rubio, D., & López-Sastre, R. J. (2016). Towards perspective-free object counting with deep learning. In *Computer Vision – ECCV 2016* (pp. 615–629). Springer International Publishing. URL: `https://doi.org/10.1007%2F978-3-319-46478-7_38`. doi:`10.1007/978-3-319-46478-7_38`.

Ren, S., He, K., Girshick, R., & Sun, J. (2017). Faster r-CNN: Towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, *39*, 1137–1149. URL: `https://doi.org/10.1109%2Ftpami.2016.2577031`. doi:`10.1109/tpami.2016.2577031`.

Sindagi, V. A., & Patel, V. M. (2018). A survey of recent advances in CNN-based single image crowd counting and density estimation. *Pattern Recognition Letters*, *107*, 3–16. URL: `https://doi.org/10.1016%2Fj.patrec.2017.07.007`. doi:`10.1016/j.patrec.2017.07.007`.

Ujjan, R. M. A., Pervez, Z., Dahal, K., Bashir, A. K., Mumtaz, R., & González, J. (2020). Towards sFlow and adaptive polling sampling for deep learning based DDoS detection in SDN. *Future Generation Computer Systems*, *111*, 763–779. URL: `https://doi.org/10.1016%2Fj.future.2019.10.015`. doi:`10.1016/j.future.2019.10.015`.

Vítek, S., & Melničuk, P. (2017). A distributed wireless camera system for the management of parking spaces. *Sensors*, *18*, 69. URL: `https://doi.org/10.3390%2Fs18010069`. doi:`10.3390/s18010069`.

Xie, W., Noble, J. A., & Zisserman, A. (2016). Microscopy cell counting and detection with fully convolutional regression networks. *Computer Meth-*

30

693     *ods in Biomechanics and Biomedical Engineering: Imaging & Visualiza-*

694     *tion*, *6*, 283–292. URL: `https://doi.org/10.1080%2F21681163.2016.`

695     `1149104`. doi:10.1080/21681163.2016.1149104.

696 Zhang, S., Wu, G., Costeira, J. P., & Moura, J. M. F. (2017). Understand-

697     ing traffic density from large-scale web camera data. In *2017 IEEE Con-*

698     *ference on Computer Vision and Pattern Recognition (CVPR)*. IEEE.

699     URL: `https://doi.org/10.1109%2Fcvpr.2017.454`. doi:10.1109/cvpr.

700     2017.454.

701 Zhang, Y., Zhou, D., Chen, S., Gao, S., & Ma, Y. (2016). Single-image

702     crowd counting via multi-column convolutional neural network. In *2016*

703     *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

704     IEEE. URL: `https://doi.org/10.1109%2Fcvpr.2016.70`. doi:10.1109/

705     cvpr.2016.70.

**Declaration of interests**

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