# UAV-based Crowd Surveillance in Post COVID-19 Era

**NIZAR MASMOUDI[1], WAEL JAAFAR[2], (Senior Member, IEEE), SAFA CHERIF[3], JIHENE BEN ABDERRAZAK[4] and HALIM YANIKOMEROGLU[5], (Fellow, IEEE)**

[1]ESPRIT School of Engineering, Tunis, Tunisia (e-mail: nizar.masmoudi@esprit.tn)
[2]Department of Systems and Computer Engineering, Carleton University, Ottawa, Canada (e-mail: waeljaafar@sce.carleton.ca)
[3]ESPRIT School of Engineering, Tunis, Tunisia (e-mail: safa.zhiouacherif@esprit.tn)
[4]ESPRIT School of Engineering, Tunis, Tunisia (e-mail: jihene.benabderrazek@esprit.tn)
[5]Department of Systems and Computer Engineering, Carleton University, Ottawa, Canada (e-mail: halim@sce.carleton.ca)

Corresponding author: Wael Jaafar (e-mail: waeljaafar@sce.carleton.ca).

This work is accepted for publication in IEEE Access. Copyright will be transferred to IEEE at the final publishing stage.

**ABSTRACT** To cope with the current pandemic situation and reinstate pseudo-normal daily life, several measures have been deployed and maintained, such as mask wearing, social distancing, hands sanitizing, etc. Since outdoor cultural events, concerts, and picnics, are gradually allowed, a close monitoring of the crowd activity is needed to avoid undesired contact and disease transmission. In this context, intelligent unmanned aerial vehicles (UAVs) can be occasionally deployed to ensure the surveillance of these activities, that health restriction measures are applied, and to trigger alerts when the latter are not respected. Consequently, we propose in this paper a complete UAV framework for intelligent monitoring of post COVID-19 outdoor activities. Specifically, we propose a three steps approach. In the first step, captured images by a UAV are analyzed using machine learning to detect and locate individuals. The second step consists of a novel coordinates mapping approach to evaluate distances among individuals, then cluster them, while the third step provides an energy-efficient and/or reliable UAV trajectory to inspect clusters for restrictions violation such as mask wearing. Obtained results provide the following insights: 1) Efficient detection of individuals depends on the angle from which the image was captured, 2) coordinates mapping is very sensitive to the estimation error in individuals' bounding boxes, and 3) UAV trajectory design algorithm 2-Opt is recommended for practical real-time deployments due to its low-complexity and near-optimal performance.

**INDEX TERMS** Object detection, clustering, Unmanned Aerial Vehicle, computer vision, image coordinates mapping.

## I. INTRODUCTION

SINCE early 2020, the world has been living a multi-wave COVID-19 pandemic, which forced populations into quarantine and limited movement episodes. However, such a situation cannot last forever for the sake of mental health. Also, with the progress of vaccination, health restrictions are being loosened precociously. Indeed, the World Health Organization started a huge campaign to spread awareness about COVID-19 and share health measures to minimize infections. Besides the wearing of masks and constant use of sanitizers to reduce viral transmission, social distancing rules were enforced. Specifically, individuals should always keep a safety distance of more than 1 or 2 meters (depending on the country) away from each other. This sparked a new interest in crowd surveillance not only to overcome the current situation but also to prevent future pandemics.

Several intelligent solutions were proposed to evaluate social distancing between people. For instance, Yang *et al.* described in [1] a system that sends a multimedia message to alert the crowd when social distancing is disrespected. Fed with a surveillance camera video-stream, their solution localizes pedestrians using a pre-trained conventional object detector. Then, inverse homography transformation is used to map image pedestrians' coordinates with real-world locations and to evaluate distances among them. Punn *et al.* adopted in [2] a similar setting, where they trained a YOLOv3 model alongside DeepSORT to localize detected people in a surveillance camera video-stream. Then, real-world coordinates were estimated using bounding boxes around detected individuals.

In addition to works relying on captured images/videos by fixed surveillance cameras, unmanned aerial vehicle (UAV) based surveillance systems have been also considered. Indeed, a large number of organizations and industries are leveraging the UAV technology to automate complicated tasks. Specifically, UAVs are capable of flying at high altitudes to expand their range of view and take advantage of obstacle-free spaces. For instance, modern agriculture uses UAVs for soil health scans, fertilization, and crop monitoring. Also, law enforcement organizations use UAVs for crime investigation and crowd control, especially during large public events. In the context of COVID-19, Somaldo *et al.* leveraged UAVs for social distance monitoring [3]. Using the UAV's ventral camera, the system feeds aerial images into a trained YOLOv3-based object detector. Based on the area expansion principle, the authors developed a coordinates mapping algorithm and evaluated social distancing conditions between detected pedestrians. Although Yang *et al.* managed in [1] to achieve accurate results, their system is sensitive to the environment setting. Furthermore, the use of fixed cameras to monitor wide open areas [1], [2] is not efficient due to the small covered area by a single camera. This issue is bypassed in [3] through the use of flexibly flying UAVs. Nevertheless, the UAV's camera has generally a limited vision range, thus requiring either flying over long distances to monitor a wide area, which consumes a significant amount of energy, or deploying several UAVs simultaneously.

Motivated by the aforementioned issues, we propose in this paper a complete UAV-based framework for outdoor crowd surveillance in the post COVID-19 era. The latter follows three steps as follows: 1) Given a single UAV that hovers and captures images on a crowded open area, the images stream is run through an object detector to fit bounding boxes around people. Due to the low-efficiency of detection model with heterogeneous individuals, i.e., standing or seated adults, kids, etc., we propose a bounding box correction process. 2) Corrected bounding boxes are then used for coordinates mapping into a 1:1 scale coordinates system. Subsequently, individuals are clustered with respect to the social distancing conditions. 3) Each cluster is attributed an infection risk score, calculated based on distances between individuals within the same cluster. Then, the UAV is deployed closer to clusters in order to check other conditions, such as mask wearing, identity, or vaccination pass. Several trajectory design algorithms are evaluated, with the objective trading-off between prioritizing clusters with high risks and reducing energy consumption. To the best of our knowledge, this is the first work that provides a complete UAV-based surveillance framework for post COVID-19 circumstances. The contributions of the paper are given as follows:

1) We present a complete UAV based surveillance framework that combines three aspects, namely individuals detection and localization, coordinates mapping, distance evaluation, individuals clustering, and UAV trajectory design.

2) Differently from previous works, we propose an adaptation of the individuals localization algorithm to support more efficiently heterogeneous cases, including standing/seated adults and kids. Specifically, since low altitude wide-angle captured images raise bounding box errors with seated adults and kids detection, we propose a bounding box correction technique to reduce its effect on coordinates mapping.

3) Unlike most state-of-the-art works that rely on additional information, e.g., camera characteristics, capturing angle, GPS information, landmarks, etc., in order to map image coordinates to real-world ones, our work proposes a coordinates mapping solution based solely on the set of images, without involving any other information.

4) Finally, we propose an efficient UAV trajectory design to further investigate the crowd, for instance, to verify mask wearing restrictions. The UAV path is composed of two components: The first defines the order of visit for each individuals cluster, while the second determines how the individuals within a cluster are inspected.

The remaining of the paper is organized as follows. Section II reviews the related works. Section III describes the people detection phase. Section IV details how to map detected individuals to real-world locations, and the clustering approach. Section V formulates and solves the UAV trajectory planning problem. Then, Section VI presents and discusses the results. Finally, Section VII concludes the paper.

## II. RELATED WORK

We present in this section the most relevant works to our framework's steps, namely human detection, social distance monitoring, and UAV trajectory optimization. A summary is presented in Table 1.

### A. HUMAN DETECTION

Human detection algorithms are mainly classified into two categories: Two-stage and one-stage object detectors.

Two-stage object detectors split the detection process into two steps. The first step is a convolutional region network, a.k.a, region based convolutional neural networks (R-CNN). It extracts feature maps from the input image and provides a certain number of regions called regions of interest (ROI). ROIs are then classified in the second step via a dense neural network to provide the final bounding boxes that locates individuals. The R-CNN is the foundation of most two-stage object detectors [4]. Faster R-CNN, proposed in [5], managed to reach a mean average precision (mAP) of 69% on the PASCAL visual object classes (VOC) 2007 dataset in [6].

In contrast, one-stage object detectors have a fully convolutional architecture. The latter behaves as a single regression model predicting bounding boxes and class probabilities from the entire image, which provides training and inference speed, and the capability to generalize through different scenarios. Redmon *et al.* were the first to introduce this architec-

ture in YOLOv1 [7]. Although YOLOv1 achieved only 63% mAP on PASCAL VOC 2007 dataset, its inference speed was significantly higher than that of two-stage object detectors, while the recently proposed YOLOv4 in [8] achieved a higher mAP with respect to the fully convolutional architecture.

Human detection can be seen as part of a bigger object detection problem. For instance, authors of [2] suggested to combine a YOLOv3-based model with DeepSORT. Their algorithm managed to reach a mAP of 84.6% when trained and evaluated on the Open Image Dataset. Several other works focused on object detection in aerial images [9]–[12]. They relied on the VisDrone2019 dataset for training and testing, but using different machine learning models. For instance, authors of [9] trained a CenterNet algorithm with a large input image resolution (2048 x 2048) in order to avoid the loss of discriminatory features for small objects. Their model reached mAP of 65% on the validation subset on human classes, i.e., *people* and *pedestrian*. Also, the authors of [10] proposed a novel network structure called PENet (Points Estimated Network) that achieved 41% mAP on the validation subset on human classes, while authors in [11] combined Cascaded R-CNN with CenterNet through weighted box ensemble to design a new network called SyNet. Their proposed model detected humans with mAP of 43%. In contrast, Yu *et al.* proposed in [12] an algorithm based on Faster R-CNN, called DSHNet. Between *people* and *pedestrian* classes, their model recorded a mAP of only 19.5%. The aforementioned works demonstrate the difficulty to efficiently detect humans through aerial imagery. Indeed, VisDrone2019, amongst only few other UAV-based image datasets, offers different perspectives, scales, and angles of capture, which complicates the detection algorithm generalization.

### B. SOCIAL DISTANCE MONITORING

Social distance monitoring is a complex task that relies essentially on mapping image coordinates to real-world coordinates. In this context, several works investigated mapping functions. In [1], authors used the inverse homography transformation matrix, while [2] relied on the convex lens geometric properties and the dimensions of bounding boxes to estimate the depth of each detected individual. Moreover, authors of [13] focused on implementing a camera calibration method to estimate intrinsic and extrinsic camera parameters. They developed a corner detection algorithm to establish the real-world coordinates system. In [3], individuals' real coordinates were estimated using the characteristics of the used UAV-mounted camera and a calibration constant defined based on the area expansion principle, while homography based estimation was realized in [14] by relaying on capturing at least four landmarks.

Although the aforementioned works successfully mapped image coordinates to real-world ones, most of them were built under several assumptions and camera parameters knowledge. For instance, [1], [13] and [14] require the presence of essential landmarks in order to estimate their
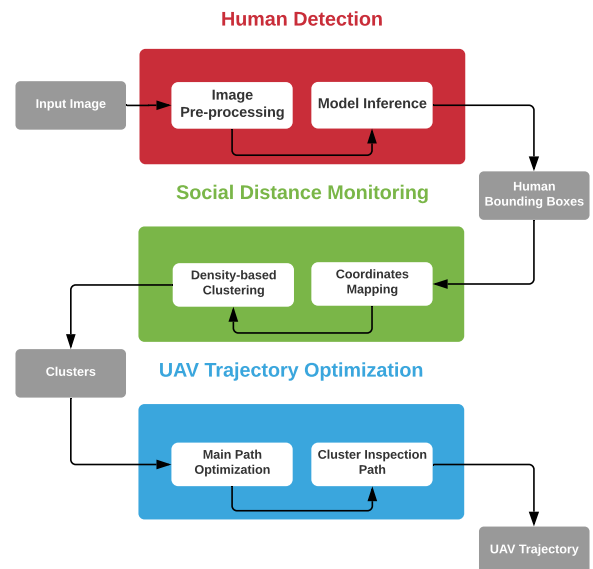


**FIGURE 1.** Proposed UAV framework for intelligent monitoring of post COVID-19 activities.

mapping function. These assumptions cannot be supported with others datasets, such as VisDrone2019.

### C. UAV TRAJECTORY OPTIMIZATION

With the growing interest in leveraging UAVs for automating complex tasks, UAV trajectory optimization is one of the most critical issues. In [15], the authors proposed UAV path design aiming to minimize energy consumption, data transmission, and coverage fairness, while authors of [16] solved optimally the UAV trajectory problem, targeting minimum cellular data offloading costs, using Hamiltonian-Jacobi equations (H&J eq.). In [17], the authors investigated the UAV path planning problem to maximize the total score collected from specific regions within a mission time constraint. Their solution inspired from the nearest neighbor (NN) algorithm built a trajectory by progressively adding nodes that maximize the ration reward by distance. Finally, authors of [18] proposed several algorithms in the offline and online settings to maximize UAV data collection from ground sensors and within time deadlines. Their offline solutions included Tabu search, simulated annealing (SA), and guided local search (GLS), while the online solutions involved reinforcement learning (RL) algorithms.

The following Sections III, IV, and V, expose our crowd surveillance framework, as presented in Fig. 1. For the sake of clarity, we summarize in Table 2 the symbols used in the remaining of the paper, along their descriptions.

### III. HUMAN DETECTION

In this section, we present and explain the setup for human detection.

**TABLE 1.** Summary of related works.

| Paper | Focus | | | Features | | | |
|---|---|---|---|---|---|---|---|
| | Human detection | Coordinates mapping | UAV trajectory optimization | UAV-based | Detection model | Coords. mapping method | UAV trajectory algorithm |
| Punn *et al.* [2] | ✓ | ✓ | – | – | YOLOv3 & Deepsort | Bounding box based | – |
| Pailla *et al.* [9] | ✓ | – | – | ✓ | CenterNet | – | – |
| Tang *et al.* [10] | ✓ | – | – | ✓ | PENet | – | – |
| Albaba *et al.* [11] | ✓ | – | – | ✓ | SyNet | – | – |
| Yu *et al.* [12] | ✓ | – | – | ✓ | DSHNet | – | – |
| Yang *et al.* [1] | – | ✓ | – | – | – | Inverse homography transformation | – |
| Somaldo *et al.* [3] | – | ✓ | – | ✓ | – | – | – |
| Babinec *et al.* [14] | – | ✓ | – | ✓ | – | Inverse homography transformation | – |
| Siswantoro *et al.* [13] | – | ✓ | – | – | – | Camera calibration | – |
| Zhang *et al.* [15] | – | – | ✓ | ✓ | – | – | RL |
| Coupechoux *et al.* [16] | – | – | ✓ | ✓ | – | – | H&J eq. |
| Fountoulakis *et al.* [17] | – | – | ✓ | ✓ | – | – | NN-based |
| Ghdiri *et al.* [18] | – | – | ✓ | ✓ | – | – | Tabu,SA,GLS,RL |
| This work | ✓ | ✓ | ✓ | ✓ | Scaled YOLOv4 | Bounding box based | 2-Opt, GA, ACO |

**TABLE 2.** List of symbols.

| Symbol | Description |
|---|---|
| $(x, y)$ | Image plane coordinates system (in pixels) |
| $(\bar{x}, \bar{y})$ | Image plane coordinates system (in meters) |
| $(\hat{x}, \hat{y}, \hat{z})$ | Real-world coordinates system (in meters) |
| $h$ | Bounding box height (in pixels) |
| $w$ | Bounding box width (in pixels) |
| $h_c$ | Corrected bounding box height (in pixels) |
| $\bar{h}$ | Corrected bounding box height (in meters) |
| $f$ | Focal length of camera lens |
| $\hat{h}$ | Assumed height of a perfectly standing adult |
| $\mathcal{G}$ | Complete directed graph of clusters |
| $\mathcal{C}$ | Set of clusters |
| $\mathcal{E}$ | Set of edges connecting clusters of $\mathcal{C}$ |
| $\mathbf{F}$ | Cost matrix of edges in $\mathcal{E}$ |
| $f_{ij}$ | Total cost of edge $(i, j) \in \mathcal{E}$ |
| $f_{ij}^p$ | Priority cost of edge $(i, j) \in \mathcal{E}$ |
| $f_{ij}^e$ | Energy consumption cost of edge $(i, j) \in \mathcal{E}$ |
| $\alpha$ | Weight of cost $f_{ij}^p$ |
| $\mathbf{c}_i$ | Coordinates of cluster $i$ in system $(\hat{x}, \hat{y})$ |
| $\lambda_i$ | Infection risk score of cluster $i$ |
| $N_i$ | Number of individuals in cluster $i$ |
| $t^*$ | Optimized UAV trajectory |
| $\mathcal{W}$ | Set of convex hull points related to a cluster |
| $d_S$ | Safety distance |

### A. AERIAL IMAGES DATASET

Amongst only few public datasets captured by camera-equipped UAVs, VisDrone2019 [19] stands out as one of the largest, most diverse and carefully annotated benchmark datasets. VisDrone2019 was provided by AISKEYE team at the Lab of Machine Learning and Data Mining in Tianjin University, China and was designed for multiple tasks with over 10 different object classes (people, pedestrian, car, van, truck, . . . ). With over 8,000 images, VisDrone2019 images were captured from different altitudes and angles and under diverse circumstances, namely night, rain, lens flare, and motion blur. Besides, the captured images vary in size from 480 x 360 up to 2000 x 1500.

VisDrone2019 was designed for machine learning, hence the images were carefully split into a training set (6,471 images), a validation set (1,610 images) and a test set (548 images).

AISKEYE team provided comprehensive annotations for each image presenting useful information, namely the bounding boxes' locations as well as their dimensions, the encoded object category and the truncation and occlusion ratios of each object. Given the annotation complexity of some regions due to a low resolution and/or a dense crowded area, AISKEYE team identified portions of images to be ignored while training or evaluation in order to avoid misleading the model. Annotated samples from the VisDrone2019 dataset are presented in Fig. 2.
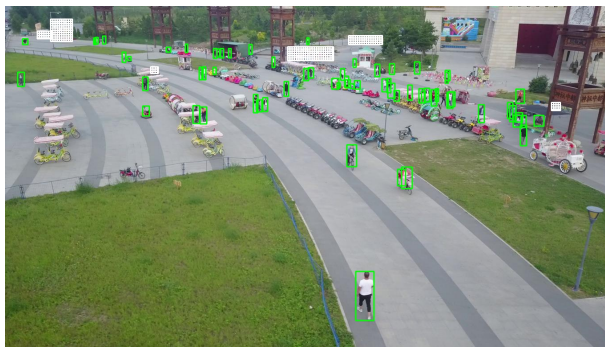
### B. PRE-PROCESSING

First of all, ignored regions that were likely to disturb the model's development were replaced with white noise as new pixels were drawn from a Gaussian distribution.

As mentioned previously, VisDrone2019 images can be very large (up to 2000 x 1500 pixels). In order to efficiently extract features from small objects (5 x 15 pixels), it is imperative that the network's input dimensions at least match the dimensions of these large images in order to avoid input compression and the loss of features crucial for detection. Evidently, large networks require more computational power. Due to the inaccessibility to high-end hardware, we were obligated to split the images into 4 portions each lowering, therefore, the required size of the network. Each portion was, then, resized to 736 x 736 pixels.

The splitting process generates a heavy amount of negative samples. Although these images may present an addition to the training set, their influence on the model's precision is very minimal compared to the amount of training time they cost. In fact, object detectors consider each group of pixels that does not belong to the annotated objects as negatives. Given that the objective of this task is to only detect
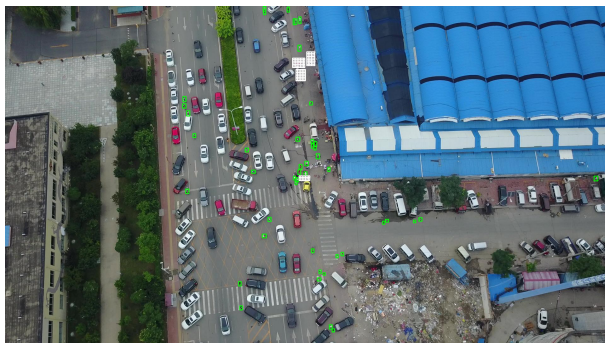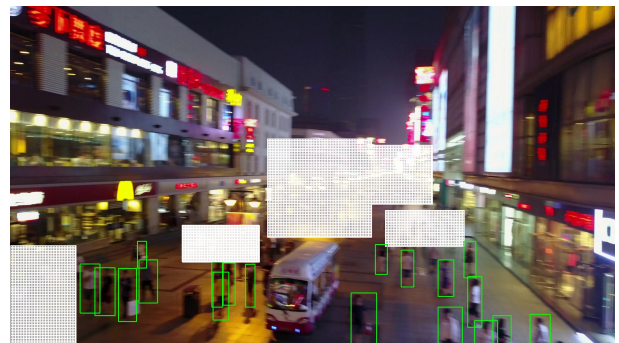
(a) Image 0000011_00234_d_0000001.

(b) Image 0000031_02000_d_0000041.

(c) Image 0000129_02411_d_0000138.

(d) Image 0000074_08777_d_0000017.

**FIGURE 2.** Annotated samples from VisDrone2019 dataset (Green rectangles identify humans; White grids are ignored regions).

humans, *pedestrian* and *people* classes were merged into one class while other objects were omitted from the dataset.

### C. SCALED YOLOV4

For human detection, we opted in this work for the recent Scaled YOLOv4 model [20]. Authors of [20] proposed initially the YOLOv4 model, which inherits the fully convolutional architecture of its ancestors. Although this object detection model achieved state-of-the-art results, Bochkovskiy *et al.* further enhanced their architecture by readjusting the backbone, adding more Cross Stage Partial (CSP) connections and scaling up the model to get the best speed/accuracy trade-off. This novel architecture, called Scaled YOLOv4, has been proven to achieve high detection results. Indeed, as shown in [20], ground-breaking results of 55.4% in average precision (AP) on the Microsoft Common Objects in Context (MS COCO) dataset with a relatively high inference speed is realized.

### D. MODEL TRAINING

In order to train the model for our dataset, an implementation of scaled YOLOv4 on GitHub is cloned to a Google Colaboratory environment [21], [22]. The GPU provided by the environment is NVIDIA Tesla T4.

For our work, we initialize scaled YOLOv4 with the pre-trained weights on the MS COCO dataset. The learning rate starts with a very low value and gradually increases during a short warm up phase until it reaches 0.001. The image input size is set to $736 \times 736$ pixels, and a preset of 9 bounding boxes were calculated with the K-means algorithm on the training subset [23]. The algorithm uses the complete intersection over union (CIoU) loss [24] and Nesterov accelerated gradient for back-propagation [25]. The implemented Scaled YOLOv4 model is trained over a total of 47,224 batches of 8 images each.

## IV. IMAGE ANALYSIS AND CLUSTERING

This section details how an image is analyzed to extract detected humans coordinates, then to proceed with coordinates mapping and clustering.

### A. PINHOLE CAMERA MODEL

Due to the low-altitude wide-angle captured images in the VisDrone2019 dataset, close objects to the camera would appear larger/longer than further ones. To correctly represent this optical phenomenon, we leverage the pinhole camera model in the following analysis [26]. The latter is commonly used in computer vision to mimic the geometrical projection of light on the image plane of the camera.

### B. RELATION BETWEEN THE COORDINATES SYSTEMS

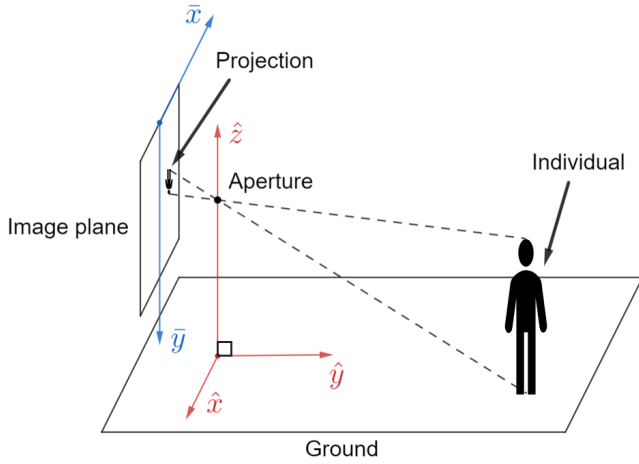In order to efficiently design a mapping function, two orthogonal coordinates systems are established, the image plane

**FIGURE 3.** Pinhole camera model illustration and relation between the image plane and real-world coordinates systems.
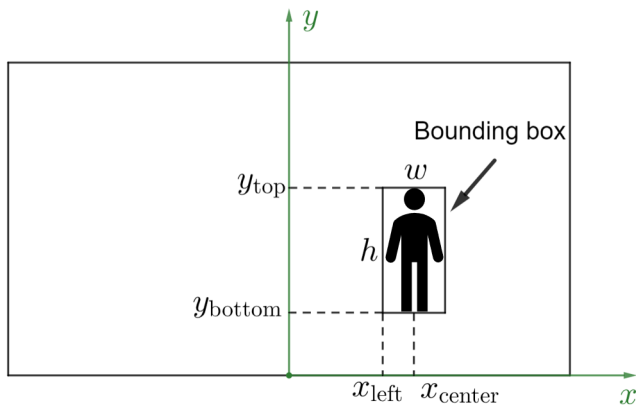


**FIGURE 4.** Bounding box definition in the image coordinates system.

coordinates system $(\bar{x}, \bar{y})$ that localizes individuals' projection within the image plane of the camera and a real-world coordinates system $(\hat{x}, \hat{y}, \hat{z})$ that localizes individuals at a distance from the camera lens, illustrated in Fig. 3. The ground is assumed perfectly flat, thus, any individual Z-coordinate is considered null in the real-world coordinates system.

Any detected individual is delimited by a bounding box in the image plane. Bounding boxes are described using specific annotations. Indeed, the detection model outputs a vector for each bounding box presenting multiple features, including the bounding box location $(x_{\text{left}}, y_{\text{top}})$, width $w$, and height $h$, in pixels, as shown in Fig. 4. For more accurate localization, we define new coordinates as follows:

$$x_{\text{center}} = x_{\text{left}} + \frac{w}{2} \tag{1a}$$

$$y_{\text{bottom}} = y_{\text{top}} - h. \tag{1b}$$

## C. BOUNDING BOX CORRECTION

As mentioned formerly, our intention is to leverage the object as well as its projection's dimensions in order to estimate its

location in the real-world coordinates system. Nonetheless, the objects in our setting are humans, where knowledge about the height of each individual cannot be determined, given the variable heights between adults, children, men, and women. Furthermore, bounding boxes fit only around pixels deemed to be part of a human's features and do not take into consideration posture and stance, i.e., standing, crouching, or sitting. As such, two major assumptions are made:

- People captured in an image are assumed to be perfectly standing adults.
- The height of a perfectly standing adult, either a man or a woman, is approximated to 1.75 m.

These assumptions may seem naive since stumbling upon a group of 1.75 m tall standing adults is very unlikely. Accordingly, children and seated/crouching adults will be perceived further than their actual real-world location as their bounding boxes are significantly smaller than that of a 1.75 m tall standing adult. In order to attenuate this effect, we introduce a bounding box correction mechanism. The idea is to bring all bounding boxes of detected individuals into an average state. Assuming that the majority of detected individuals follow the previous assumptions, the outliers, i.e., not standing adults or children, will see their bounding boxes adjusted to their average, with consideration to the pinhole camera model.

For a set of $L$ original bounding boxes $(y_{\text{bottom}}^{(i)}, h^{(i)})$, $\forall i = 1, \ldots, L$, drawn inside an image, a second degree polynomial function can be used to model the average height of each bounding box as a function of its Y-coordinate in the image plane, as follows:

$$h^{(i)} = \alpha_0 + \alpha_1 y_{\text{bottom}}^{(i)} + \alpha_2 (y_{\text{bottom}}^{(i)})^2 + \epsilon^{(i)} = g(y_{\text{bottom}}^{(i)}) + \epsilon^{(i)}, \tag{2}$$

where $\epsilon^{(i)}$ is the deviation of the $i^{th}$ bounding box size from the standard size and $g(\cdot)$ is a second degree polynomial function with parameters $\alpha_0$, $\alpha_1$, and $\alpha_2$.

In order to determine these parameters, we formulate the associated regression problem with the following least-square objective function:

$$\min_{\alpha_0, \alpha_1, \alpha_2} \sum_{i=1}^{L} (h^{(i)} - g(y_{\text{bottom}}^{(i)}))^2. \tag{3}$$

Subsequently, $\alpha_0$, $\alpha_1$, and $\alpha_2$ are obtained by solving the equality of corresponding partial derivatives to 0. Eventually, the heights of bounding boxes are corrected by subtracting the deviation from their original heights as follows:

$$h_c^{(i)} = h^{(i)} - \epsilon^{(i)} = g(y_{\text{bottom}}^{(i)}), \tag{4}$$

where $h_c^{(i)}$ is the corrected height of the $i^{th}$ bounding box.

## D. COORDINATES MAPPING

In order to obtain the metric dimensions of the objects projected on the image plane of the camera, we multiply

(a) Projection on plane $(\hat{y}, \hat{z})$.



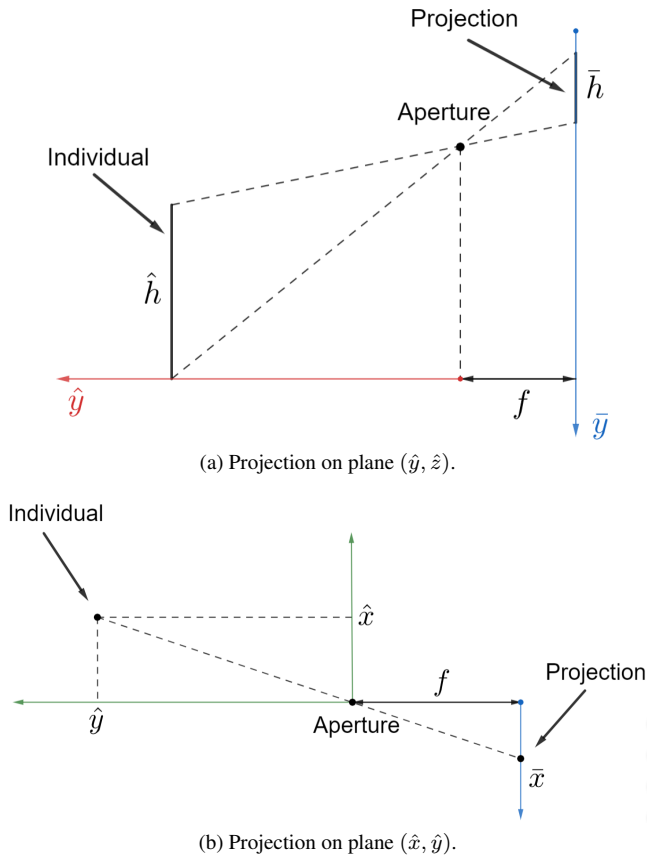(b) Projection on plane $(\hat{x}, \hat{y})$.

**FIGURE 5.** Projection of the pinhole camera model illustration on real-world system planes.

coordinates and dimensions by the pixel size of the camera sensor, denoted $p$. Hence, we obtain

$$\bar{x}_{center} = p \times x_{\text{center}} \tag{5a}$$

$$\bar{y}_{bottom} = p \times y_{bottom} \tag{5b}$$

$$\bar{h} = p \times h_c. \tag{5c}$$

For coordinates mapping purposes, we assume that the camera axis is perfectly parallel to the ground plane. Also, let $f$ be the focal length of the camera, i.e., distance between aperture and image plane, and $(\hat{x}_{center}, \hat{y}_{bottom})$ the estimated real-world ground coordinates of an individual in meters. Using different projection perspectives of the pinhole camera model, as illustrated in Fig. 5, the mapping function is designed using the homothety rules, as follows:

$$\hat{y} = \frac{\hat{h}f}{\bar{h}} \tag{6a}$$

$$\hat{x} = \frac{\bar{x}_{center}\,\hat{y}}{f} = \bar{x}_{center}\frac{\hat{h}}{\bar{h}}, \tag{6b}$$

where $\hat{h}$ is the height of the real-world object, assumed $\hat{h} = 1.75$ m.

### E. DENSITY-BASED CLUSTERING

Given the randomness of the distribution of people captured by the UAV camera, partition-based and hierarchical cluster-ing techniques are not adequate to identify groups with dense numbers of individuals. In contrast, density-based clustering is known to be efficient in recognizing randomly shaped clusters and in dealing with outliers. Hence, we opt here for density-based clustering to identify situations where social distancing is potentially not respected. Specifically, we use the popular Density-based Spatial Clustering of Applications with Noise (DBSCAN) algorithm for its intuitive and simple approach [27].

## V. UAV TRAJECTORY OPTIMIZATION

Within our crowd surveillance framework, we assume that the UAV acts in two phases. In the first, it flies at a relatively high altitude to detect people, determine their locations, and cluster them with respect to social distancing restrictions (Sections III and IV). In the second phase, it flies closer to the clusters with high risk of people closeness in order to verify, for instance, that mask wearing is in effect (Section V). The UAV trajectory for the second phase should be carefully designed in order to tradeoff between UAV energy consumption, i.e., shorter flying distance, and cluster visiting prioritization, i.e., clusters with higher risks should be visited first. Moreover, how to inspect each cluster (through hovering or slow flying) is a challenging task as it also needs to trade-off between energy consumption and efficacy of inspection. Given the complexity of the aforementioned problem, the entire process is divided to two separate sub-problems. In the first, we optimize the main UAV trajectory, i.e., to fly from a cluster to another, while in the second, we design an efficient cluster inspection path.

### A. MAIN TRAJECTORY OPTIMIZATION

The goal of the main trajectory design problem is to follow an energy-efficient path to visit all clusters and return to its initial location, while prioritizing high-risk clusters over low-risk ones. For the sake of simplicity, we assume that the aerial environment is obstacle-free, and that the UAV flies sufficiently high to avoid collisions with people or objects. This problem can be modeled using a complete directed graph $\mathcal{G} = (\mathcal{C}, \mathcal{E}, \mathbf{F})$ where $\mathcal{C} = \{0, \dots, |\mathcal{C}|\}$ is the set of clusters, with $|\mathcal{C}|$ is the cardinality of $\mathcal{C}$ and cluster $0$ identifies the UAV's initial location. $\mathcal{E}$ is the set of edges connecting the clusters, and $\mathbf{F} = [f_{ij}]_{|\mathcal{C}| \times |\mathcal{C}|}$ is the cost matrix of the edges, where $f_{ij}$ corresponds to the cost of directed edge $(i, j)$. Inspired by the Dantzig–Fulkerson–Johnson formulation [28], our problem can be written as

$$\min_{\mathbf{B}} \quad \sum_{i=0}^{|\mathcal{C}|-1} \sum_{j \neq i, j=0}^{|\mathcal{C}|-1} b_{ij} f_{ij} \tag{7}$$

$$\text{s.t.} \quad \sum_{\substack{i=0 \\ i \neq j}}^{|\mathcal{C}|-1} b_{ij} = 1 \tag{7a}$$

$$\sum_{\substack{j=0 \\ j \neq i}}^{|\mathcal{C}|-1} b_{ij} = 1 \tag{7b}$$

$$\sum_{i=0}^{|\mathcal{C}|-1} \sum_{\substack{j=0 \\ j \neq i}}^{|\mathcal{C}|-1} b_{ij} = |\mathcal{C}|, \tag{7c}$$

where $\mathbf{B} = [b_{ij}]_{|\mathcal{C}| \times |\mathcal{C}|}$ is a binary decision matrix that maps edges into the UAV trajectory, defined with

$$b_{ij} = \begin{cases} 1, & \text{if path goes from cluster } i \text{ to cluster } j \\ 0, & \text{otherwise.} \end{cases} \tag{8}$$

Constraints (7a) and (7b) guarantee that all clusters are visited only once, while (7c) ensures that a single tour covers all clusters. Clearly, problem (7) can be assimilated to a traveling salesman problem (TSP), which is known to be NP-hard. Consequently, by reduction, this problem is also NP-hard.

Since our goal trades off between energy consumption and cluster priority adherence, we define the cost function of an edge as the weighted sum of two functions as follows:

$$f_{ij} = \alpha f_{ij}^{p} + (1 - \alpha) f_{ij}^{e}, \tag{9}$$

where $f_{ij}^{p}$ and $f_{ij}^{e}$ are the priority cost and energy consumption cost of edge $(i, j)$, respectively, and $\alpha \in [0, 1]$ is the weight.

Assuming the UAV's velocity is uniform from one node to another, its energy consumption is linearly related to the traveled distance. Thus, the energy cost of an edge $(i, j)$ is equivalent to the distance between vertices $i$ and $j$. Subsequently, the energy consumption cost of edge $(i, j)$ can be defined by

$$f_{ij}^{e} = ||\mathbf{c}_{j} - \mathbf{c}_{i}||_{2} \tag{10}$$

where $\mathbf{c}_i$ and $\mathbf{c}_j$ are the coordinates of clusters $i$ and $j$, respectively, defined as the barycenters of their clusters, and $|| \cdot ||$ is the Euclidean norm. Moreover, the priority cost of an edge $(i, j)$ is defined as

$$f_{ij}^{p} = \begin{cases} \max(0, \ \lambda_j - \lambda_i), & \text{if } i \neq 0 \\ 0, & \text{otherwise} \end{cases} \tag{11}$$

where $\lambda_i$ and $\lambda_j$ are the infection risk scores of clusters $i$ and $j$, respectively, where the infection risk of a cluster $i$ is given by

---

**Algorithm 1:** 2-Opt Algorithm

**Input** : Set of clusters $\mathcal{C}$ and corresponding cost matrix $\mathbf{F}$.

**Output:** Optimized UAV trajectory $t^*$.

$t \leftarrow$ random trajectory ;
$improved \leftarrow True$;
**while** $improved$ **do**
  $\quad improved \leftarrow False$;
  $\quad$**for** $i \leftarrow 0$ **to** $|\mathcal{C}| - 2$ **by** 1 **do**
    $\quad\quad$**for** $j \leftarrow i + 2$ **to** $|\mathcal{C}|$ **by** 1 **do**
      $\quad\quad\quad$**if** $f_{i(i+1)} + f_{j(j+1)} > f_{ij} + f_{(i+1)(j+1)}$
      $\quad\quad\quad$**then**
        $\quad\quad\quad\quad$Swap edges $i$ and $j$ of trajectory $t$;
        $\quad\quad\quad\quad improved \leftarrow True$;

$t^* \leftarrow t$;

---

$$\lambda_i = \frac{2}{N_i(N_i - 1)} \sum_{l=1}^{N_i} \sum_{\substack{k=l+1 \\ k>l}}^{N_i} \mathbb{1}_{\{d_{kl} < 2\}} + N_i, \tag{12}$$

where $k$ and $l$ are the indexes of detected individuals in cluster $i$, $d_{kl}$ the distance in meters separating them, and $N_i$ is the number of individuals in cluster $i$. Finally, $\mathbb{1}_x$ assigns value 1 when condition $x$ is satisfied. To be noted that $\lambda_i$ is the sum of two components, the first is the social distancing disrespect ratio, while the second is the size of the cluster.

Several approaches can be used to solve problem (7). Intuitively, one would explore all possible combination in the graph $\mathcal{G}$ to obtain the optimal solution. This exhaustive search may work optimally for small-sized systems, however, it scales poorly with the number of clusters as $\mathcal{O}(|\mathcal{C}|!)$. For instance, with only 10 clusters more than 3 millions trajectory combinations need to be evaluated, which is very resource and time consuming, two factors that are scarce in crowd surveillance operations. Nevertheless, it is considered as a benchmark (for small-sized systems) in the next Section. Alternatively, we opt for heuristic and metaheursitic approaches as follows:

- **2-Opt [29]** is a lightweight iterative parameter-free TSP algorithm that consists of systematically swapping paths between vertices until a certain optimum is reached. We summarize our implementation of 2-Opt to solve problem (7) in Algorithm 1.

- **Genetic Algorithm (GA) [30]** reflects the process of natural selection where the fittest solutions are selected for breeding in order to produce better offspring for the upcoming generation. It evaluates full trajectories using the fitness function. Given a trajectory or route $r$, the fitness value $\phi$ is given by

$$\phi(r) = \frac{1}{\sum_{(i,j) \in \mathcal{E}} b_{ij} f_{ij}}. \tag{13}$$

---

**Algorithm 2:** Genetic Algorithm

**Input** : Set of clusters $\mathcal{C}$ and corresponding cost matrix $\mathbf{F}$.

**Output:** Optimized UAV trajectory $t^*$.

$\mathcal{P} \leftarrow$ Set of random trajectories (population);
`// `$N_g$`: number of generations`
**for** $i \leftarrow 1$ **to** $N_g$ **by** 1 **do**
  Initialize an empty set $\widetilde{\mathcal{P}}$;
  Calculate fitness of individuals in $\mathcal{P}$ using (13);
  Add the two fittest solutions of $\mathcal{P}$ to the new set $\widetilde{\mathcal{P}}$;
  Select mating pool using rank selection ($\mathcal{N} \subset \mathcal{P}$);
  Create offspring from mating pool $\mathcal{N}$ using ordered crossover;
  Use swap mutation on offspring;
  Add mutated offspring to $\widetilde{\mathcal{P}}$;
  $\mathcal{P} \leftarrow \widetilde{\mathcal{P}}$;

$t^* \leftarrow$ solution with best fitness from $\mathcal{P}$;

---

**Algorithm 3:** ACO Algorithm

**Input** : Set of clusters $\mathcal{C}$, corresponding cost matrix $\mathbf{F}$, colony size $s$, information elicitation factor $\delta$, $\beta$, pheromone intensity $Q$, and pheromone evaporation coefficient $\rho$.

**Output:** Optimized UAV trajectory $t^*$.

$\mathbf{H} \leftarrow [1/f_{ij}]_{|\mathcal{C}| \times |\mathcal{C}|}$;
Initialize pheromone matrix $\mathbf{T}$;
**for** $i \leftarrow 1$ **to** $N$ **by** 1 **do**
  Place a set $\mathcal{K}$ of $s$ ants in cluster 0;
  Calculate node selection probability matrix;
  Advance ants through a Hamiltonian cycle using predetermined probabilities;
  Update pheromone matrix $\mathbf{T}$;

$t^* \leftarrow$ best route found by the last colony $\mathcal{K}$;

---

For problem (7), the GA algorithm is presented in Algorithm 2. GA population $\mathcal{P}$ is a set of UAV trajectories. A gene is a cluster to visit. An individual is a UAV trajectory satisfying constraints (7a)-(7c). The parents are combined solutions. The mating pool is a collection of parents that creates the next generation. Mutations introduce variations in the population by randomly swapping clusters for a trajectory. Finally, elitism carries the best individuals to the next generation. The rank selection method is preferred over the roulette wheel approach to minimize the election of defective solutions for the mating pool.

Also, given the order constraint of the TSP, the ordered-crossover operator is used to generate offspring.

• **Ant Colony Optimization (ACO) [31]** emulates the swarm behavior of ants when looking for the shortest route from their shelter to a food source. In fact, randomly placed artificial ants inside a graph gradually advance through a Hamiltonian cycle using a probabilistic approach based on the pheromone intensity of the edge as well as the visibility of the target node. Such process is simulated through multiple ant colonies until a convergence criteria is met. ACO is similar to GA in the way it considers weights of edges as rewards rather than penalties. Thus, we define its visibility matrix as $\mathbf{H} = [\eta_{ij}] = [1/f_{ij}]_{|\mathcal{C}| \times |\mathcal{C}|}$. We implement ACO to solve our problem as described in Algorithm 3.

### B. CLUSTER INSPECTION PATH DESIGN

Once clusters and formed and main UAV trajectory is designed, the UAV has to adopt an inspection strategy when getting closer to clusters. Since individuals may be standing in different angles with respect to the UAV's perspective, we opt for a circular inspection path that would capture individuals from the front. Nevertheless, a safety distance must be kept to avoid any harm to individuals or to the UAV. Hence, we design the inspection path in two steps, namely convex hull design and safety distance design.

Since the size of a cluster is small, we propose to use the Jarvis March algorithm, which determines the convex hull surrounding a cluster's individuals [32]. The latter is time-efficient, achieving a complexity of $\mathcal{O}(N_i \eta_i)$, where $\eta_i$ is the number of points of the convex hull, linked to cluster $i$. The convex hull approach implies that the UAV has to fly dangerously and directly over some of the individuals in a cluster (the ones being included in the convex hull). This might compromise both safety and the quality of inspection for those individuals. Consequently, we expend the convex hull method with a safety constraint as follows.

First, let $\mathcal{W} = \{\omega_1, \omega_2, \ldots, \omega_{|\mathcal{W}|}\}$ be the set of convex hull points related to a cluster, and the closed ordered chain of points $\mathcal{W}' = \{\omega_1, \omega_2, \ldots, \omega_{|\mathcal{W}|}, \omega_1\}$ is oriented clockwise. Consequently, the following convex hull property holds:

$$(\overrightarrow{\omega_i \omega_j} \wedge \overrightarrow{\omega_i o}) \cdot \hat{z} < 0, \ \forall \ 1 \leq i < j \leq |\mathcal{W}|, \quad (14)$$

where $o$ is the barycenter of the shape constructed using $\mathcal{W}'$'s points, while $\wedge$ and $\cdot$ are the cross-product and dot-product of two vectors, respectively. Property (14) allows to identify the outbound of a cluster, from one convex hull edge perspective. Subsequently, we propose to geometrically translate each convex hull edge away from the cluster by a safety distance $d_S$ along one of its normal vectors $\overrightarrow{v}$, i.e.,

$$(\overrightarrow{\omega_i \omega_j} \wedge d_S \overrightarrow{v}) \cdot \hat{z} > 0, \ \forall \ 1 \leq i < j \leq |\mathcal{W}|. \quad (15)$$

Afterwards, to complete the novel path design, we fill the gaps between the translated convex hull edges using arcs with centers as the points in $\mathcal{W}$ and radius equal to $d_S$. The use of arcs ensures that safety distance is respected along the whole inspection path.

**TABLE 3.** Performance of the proposed human detection model.

| AP | 65% |
|---|---|
| Precision | 71% |
| Recall | 61% |
| Average IoU | 51.68% |

## VI. RESULTS AND DISCUSSION

### A. HUMAN DETECTION

The trained scaled YOLOv4 model is evaluated on 548 pre-processed images of the validation set of VisDrone2019 using the AP performance metric defined as follows:

$$AP = \int_0^1 \nu(r)dr \qquad (16)$$

where $\nu(r)$ is the precision-recall curve, obtained using the precision ($\gamma_p$) and recall ($\gamma_r$) metrics. The latter are defined as

$$\gamma_p = \frac{T_p}{T_p + F_p}, \qquad (17a)$$

$$\gamma_r = \frac{T_p}{T_p + F_n}, \qquad (17b)$$

where $T_p$ is the number of true positives, $F_p$ of false positives, and $F_n$ of false negatives. These parameters are calculated by comparing the Intersection over Union (IoU) measure of each predicted bounding box with a threshold set to 0.5. The IoU of a predicted bounding box is given by
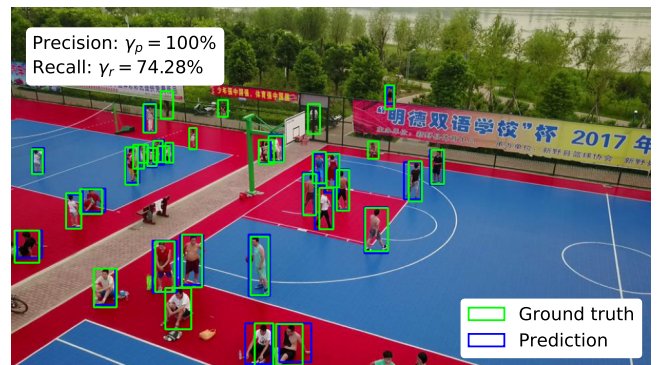
$$IoU = \frac{A_o}{A_u}, \qquad (18)$$

where $A_o$ is the area of overlap between the ground truth (in annotated images) and the prediction, and $A_u$ is the area of union of the ground truth and the prediction. Moreover, the average IoU of predicted bounding boxes is calculated to assess how well the model fits boxes around individuals.

In Table 3, we present the performances of our human detection model. First, it achieves AP of 65%, which is caused by a degraded recall of 61%. The latter means that only 61% of individuals in the evaluation set were identified. This degraded performance is issued from the difficulty encountered by the model to detect partially or highly occluded individuals, as shown in Fig. 6 below. Besides, the model fails to accurately replicate annotated bounding boxes as the average IoU is significantly low (around 52%). Such a defect can be caused by the dataset's annotation errors, which may disturb both the model's training and evaluation processes. For instance, we notice in Fig. 6 that our model correctly fitted some bounding boxes around individuals, as opposed to the erred ground truth, while in other occasions, it failed to do so.

Clearly, the dataset's heterogeneity, specifically the variation in scale and perspective of captured images, prevents the model from generalizing efficiently. Nevertheless, when compared to state-of-the-art methods in Table 4, we see that



(a) Image 0000022_00500_d_0000005.



(b) Image 0000086_00000_d_0000001.

**FIGURE 6.** Inference results on a sample from VisDrone2019 images.

**TABLE 4.** Comparison of different human detection models.

| Method | AP |
|---|---|
| CenterNet [9] | 65% |
| PENet [10] | 41% |
| SyNet [11] | 43% |
| DSHNet [12] | 19.5% |
| Our model | 65% |

our model outperforms the approaches of [10]–[12], while achieves the same AP as the CenterNet model [9]. Finally, our model is more interesting than CenterNet since it trains fast, while the CenterNet is known for its slow convergence and time-inefficiency as demonstrated in [33, Table 9].

### B. IMAGE ANALYSIS AND CLUSTERING

The evaluation of image analysis and clustering is split into three phases, namely bounding box correction, coordinates mapping, and individuals clustering. Due to the absence of a ground truth for assessment, these phases are evaluated based on the plain human observation. Moreover, in order to prevent error propagation from the human detection model, original annotations of bounding boxes are used for coordinates mapping and clustering. For the sake of clarity, the results are presented for the same VisDrone2019 sample image, which is shown in Fig. 6a.

Bounding box correction aims to bring original bounding boxes of detected individuals into an average state, under the
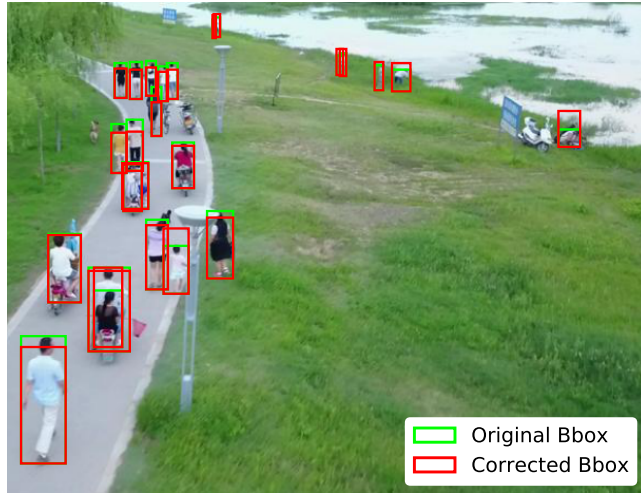
**FIGURE 7.** Bounding box (Bbox) correction performed on VisDrone2019 image 0000022_00500_d_0000005.



**FIGURE 8.** Coordinates mapping applied on corrected bounding boxes of VisDrone2019 image 0000022_00500_d_0000005.

assumption that the height of all individuals is around 1.75 m. Consequently, the bounding boxes of outliers, i.e., not standing adults and children, are averaged to their average value, with respect to the pinhole camera model. The results of this task are illustrated in Fig. 7, where green rectangles are the original bounding boxes and red rectangles the corrected ones. As it can be seen, the correction may consist on either reducing the height of the bounding box, e.g., in the case of adults, or on increasing its height, e.g., for children.

Next, given the corrected bounding boxes, the coordinates mapping task is evaluated. The proposed mapping function relies on two camera parameters, namely the focal length of the lens $f$ and the pixel size of the sensor $p$. These parameters were not provided with the VisDrone2019 dataset and the AISKEYE team did not affirm that all images were captured using the same camera specifications. To accurately execute the coordinates mapping task, we focus at first on determining the best pair $(f, p)$ that delivers most accurate results. Specifically, we apply coordinates mapping for several combinations of $(f, p)$ values over a random set of images. Then, based on human observation, we decide on which combination that provides the most accurate coordinates mapping output with respect to the original image. We found that the most accurate parameters that give acceptable representations for the selected random set of images are $(f, p) = (10 \text{ mm}, 18 \text{ } \mu\text{m})$. Subsequently, coordinates mapping is executed with these parameters and the output is shown in Fig. 8. Indeed, the estimated locations of detected individuals are plotted in the real-world coordinates system as blue dots.

In Fig. 9, we investigate the impact of human detection errors (i.e., misplaced bounding boxes) on the coordinates mapping performance. To do so, we establish the relation between the height of the annotated bounding boxes $h$, height of estimated bounding boxes $\hat{h}$, and the Y-axis original and estimated coordinates $y$ and $\hat{y}$ (calculated using (6a)). As it
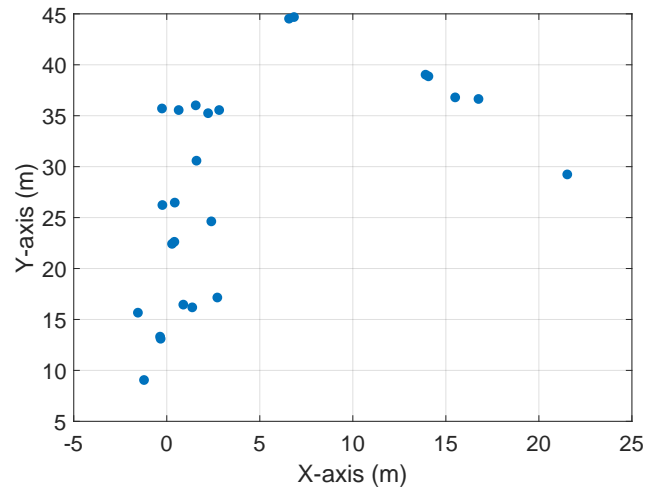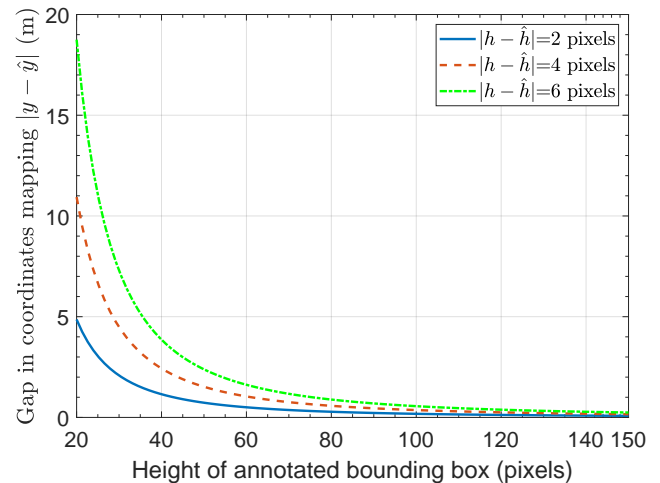


**FIGURE 9.** Impact of bounding box estimation error on coordinates mapping.

can be seen, when $h$ is small, the gap in coordinates mapping $|y - \hat{y}|$ is very high for any height gap $|h - \hat{h}|$. For instance, if $h = 40$ pixels (a typical value encountered with the VisDrone 2019 dataset) and the estimated bounding box is only 2 pixels away from the original one, the individual's location is shifted by 1 m, while a gap of 6 pixels in $|h - \hat{h}|$ causes a location estimation error of about 4 m. Although these values may seem small, they are significant when monitoring social distancing, which typically requires 2 m distance between individuals, in addition to leading to incorrect clustering.

Following the coordinates mapping phase, DBSCAN algorithm is applied for clustering. We set the DBSCAN parameters as follows: The social distancing threshold $\epsilon = 2$ m, while the minimum number of individuals within a cluster is defined as $m = 3$. In other words, groups of a single or two people are considered as outliers. The reason behind it is that groups of two people are likely to be from the same household, and thus do not need to respect social distancing. As

**FIGURE 10.** Clustering of detected individuals on VisDrone2019 image 0000022_00500_d_0000005.

**TABLE 5.** ACO parameters.

| Colony size | $s = \lceil (|\mathcal{C}| - 1)/2 \rceil$ |
|---|---|
| Information elicitation factor | $\delta = 1$ |
| Expected heuristic factor | $\beta = 5$ |
| Pheromone intensity | $Q = 10$ |
| Pheromone evaporation coefficient | $\rho = 0.5$ |

shown in Fig. 10, clusters were correctly identified, including children and seated individuals.

### C. UAV TRAJECTORY OPTIMIZATION

The UAV trajectory has two components, namely the main path and the cluster inspection path.

For the main path optimization, we choose our baseline as the optimal solution, which is obtained through exhaustive search for small-scale scenarios (i.e., images with less than 10 clusters) due to its high time complexity. The latter is compared to the proposed heuristic and meta-heuristic based ones, i.e., based on Algorithms 1 (2-Opt), 2 (GA), and 3 (ACO). We set the following parameters for GA and ACO:

- **GA:** For mutations, a random swap of nodes is implemented with rate 10%. To maintain the evolutionary aspect of the algorithm, only top-two solutions are considered for elitism. Due to the variability in number of clusters identified in each image sample of the VisDrone2019 dataset, the population size is set to twice the number of clusters. This allows to avoid both over-population and under-population, improve convergence, and explore a larger search space.
- **ACO:** The colony size is set dynamically to cope with the heterogeneity of data. The remaining of its parameters are set as in Table 5.

In Fig. 11, we present the performances of the trajectory optimization algorithms in terms of average total cost (7) and execution time, as functions of the number of clusters $|\mathcal{C}|$, and for different cost weights $\alpha$. Due to the small number of images with high number of clusters, averaged results (circles) may not present a smooth trend. Hence, we opted to add the results of third order polynomial curve fitting (solid lines) to better show the trend with the increasing number of clusters. When $\alpha = 0$ (Fig. 11a), the trajectory is optimized with regards to energy-efficiency only, reflected here through the traveled distance of the UAV. Obviously, the total cost increases rapidly with the number of clusters. This is expected since a higher $|\mathcal{C}|$ means that the UAV has to fly for longer distances to visit all clusters. Also, we notice that 2-Opt performs best compared to GA and ACO for any $|\mathcal{C}|$, and it achieves optimality for $|\mathcal{C}| < 10$. When $\alpha = 1$ (Fig. 11b), the UAV path is determined to adhere to the cluster priority order, i.e., prioritizing visiting clusters with high risk of social distancing violation. We notice here that the total cost increases slowly with $|\mathcal{C}|$, with preference for the 2-Opt based approach. However, 2-Opt is now sub-optimal for $|\mathcal{C}| < 10$. This is predictable since 2-Opt is a heuristic approach that typically finds near-optimal solutions. In addition, we remark that the priority adherence cost is approximately two orders of magnitude smaller than that of the traveled distance (compared to Fig. 11a). Thus, to accurately balance between them in the cost function, we set $\alpha = 0.99$ for the simulations of Fig. 11c. In this scenario, 2-Opt is near-optimal for $|\mathcal{C}| < 10$ while it outperforms GA and ACO solutions for any $|\mathcal{C}|$ value. Finally, the related execution times are evaluated in Fig 11d. In addition to its near-optimal performances, 2-Opt demonstrates the fastest execution time (in the range of micro and milliseconds), compared to exhaustive search, GA, and ACO. Hence, 2-Opt is seen as a practical solution for deployment in real-time environments.

For cluster inspection, we illustrate in Fig. 12 the construct of the inspection trajectory around one cluster. As it can be seen, the Jarvis March algorithm designs a non-smooth curve (dashed red line) around the cluster. Due to security concerns, a horizontal safety distance $d_S = 2$ m between any individual and the UAV is needed. To take this into account and at the same time design a practical (and thus a smooth) inspection path, we leverage our proposed approach that relies on (14)–(15).

Finally, the overall UAV trajectory is a combination of the main and cluster inspection paths. In Fig. 13, we present the complete UAV trajectory to inspect clusters of people, obtained through processing with our proposed framework (i.e., human detection, boundig box correction, coordinates mapping, clustering, and UAV trajectory design) image 0000011_00234_d_0000001 of the VisDrone2019 dataset, shown here as Fig. 2a. In Fig. 13, we mention the infection risk score of each cluster. Finally, given $\alpha = 1$, the UAV trajectory is designed to prioritize clusters with high risk scores.

### VII. CONCLUSION

In this paper, we investigated UAV-based crowd monitoring for post COVID-19 outdoor activities. Based on captured im-

(a) Total cost vs. number of clusters ($\alpha = 0$).

(b) Total cost vs. number of clusters ($\alpha = 1$).

(c) Total cost vs. number of clusters ($\alpha = 0.99$).

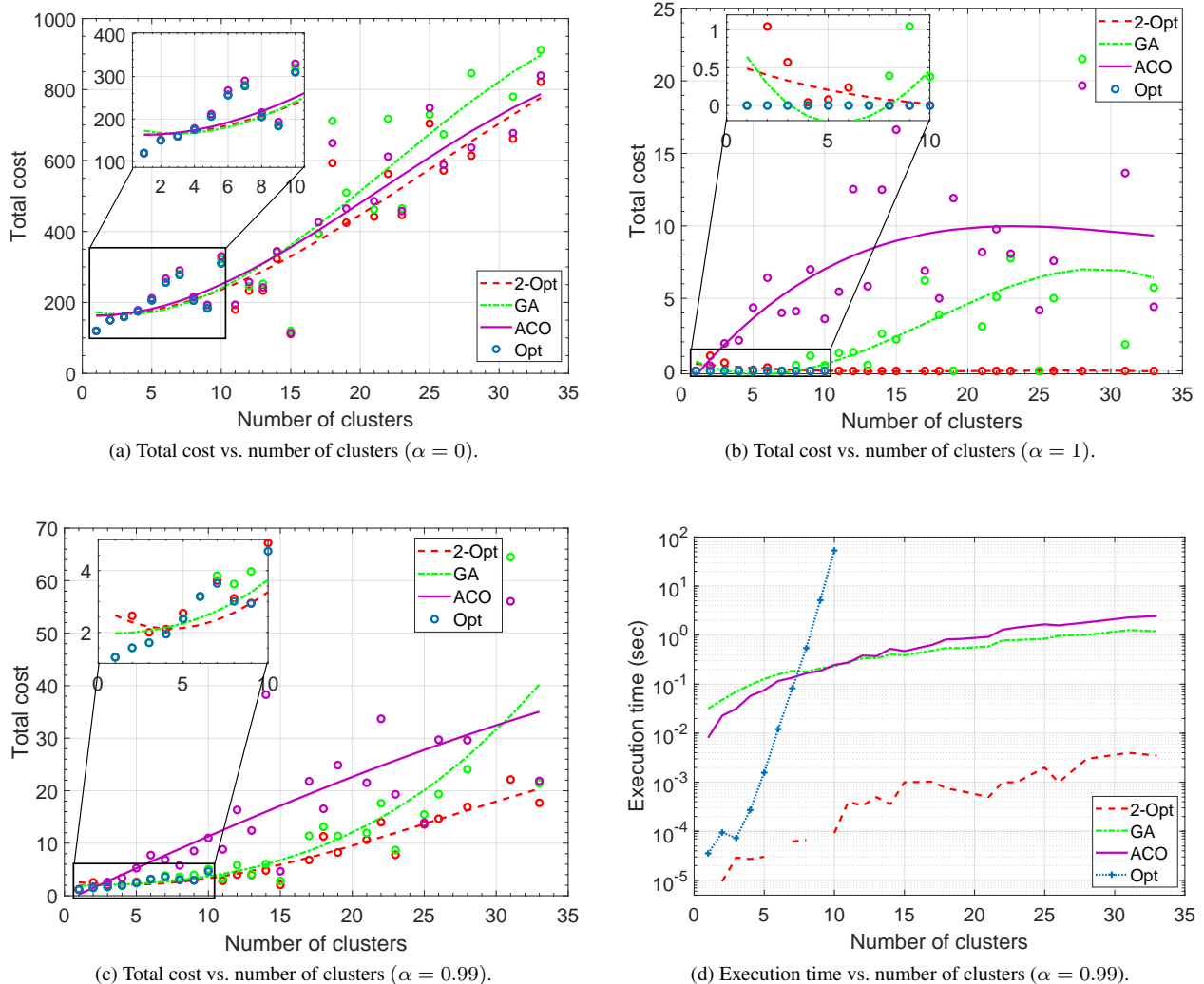(d) Execution time vs. number of clusters ($\alpha = 0.99$).

**FIGURE 11.** Performances of UAV trajectory planning algorithms (total cost and execution time) vs. number of clusters.

ages from a flying UAV, we proposed a complete surveillance framework composed of three steps: 1) Human detection and localization in captured images, 2) bounding box correction, coordinates mapping and individuals clustering in a real-world coordinates system, and 3) UAV trajectory planning for further inspection of clusters. The first step was realized using the scaled YOLOv4 approach. The second relied on optical assumptions, the pinhole camera model, and DB-SCAN clustering. Finally, the third step was realized based on a number of heuristic and meta-heuristic approaches as well as a novel proposed algorithm for cluster inspection. The obtained results draw the following insights: 1) Efficient human detection depends on the angle from which the image was captured, 2) coordinates mapping is very sensitive to the estimation error in individuals' bounding boxes drawing, and 3) 2-Opt presented the best performances, in terms of cost and execution time, compared to baseline approaches, and thus is preferred for practical real-time deployments.

## REFERENCES

[1] Dongfang Yang, Ekim Yurtsever, Vishnu Renganathan, Keith A. Redmill, and Ümit Özgüner. A vision-based social distancing and critical density detection system for COVID-19. Sensors, 21(13), Jul. 2021.

[2] Narinder Singh Punn, Sanjay Kumar Sonbhadra, Sonali Agarwal, and Gaurav Rai. Monitoring COVID-19 social distancing with person detection and tracking via fine-tuned YOLO v3 and Deepsort techniques. ArXiv:2005.01385, May 2021.

[3] Pray Somaldo, Faizal Adila Ferdiansyah, Grafika Jati, and Wisnu Jatmiko. Developing smart COVID-19 social distancing surveillance drone using YOLO implemented in robot operating system simulation environment. In Proc. IEEE R10 Humanit. Technol. Conf. (IEEE R10-HTC), pages 1–6, Dec. 2020.

[4] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In Proc. IEEE Conf. Computer Vision Pattern Recognition (CVPR), pages 580–587, 2014.

[5] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. IEEE Trans Pattern Anal. Machine Intelli., 39(6):1137–1149, Jun. 2017.

[6] Mark Everingham, Luc Van Gool, Christopher Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (VOC) challenge. Int. J. Computer Vision, 88:303–338, Jun. 2010.

[7] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You

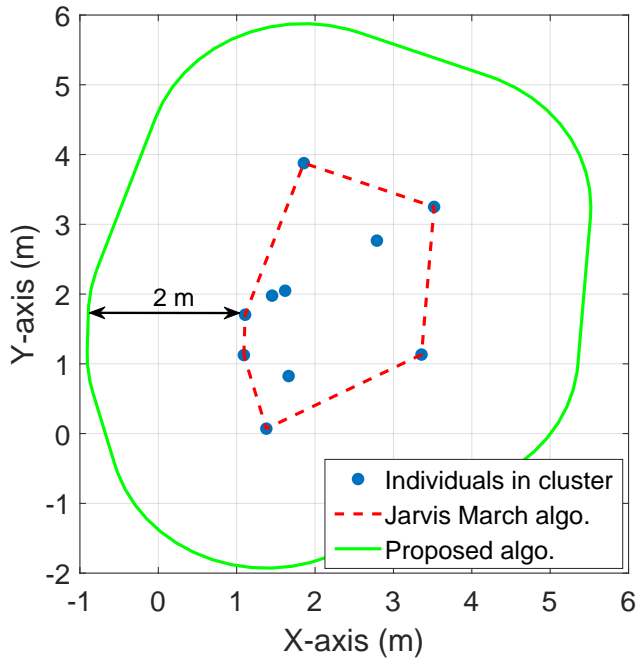**FIGURE 12.** Cluster inspection path design ($d_s = 2\ m$).
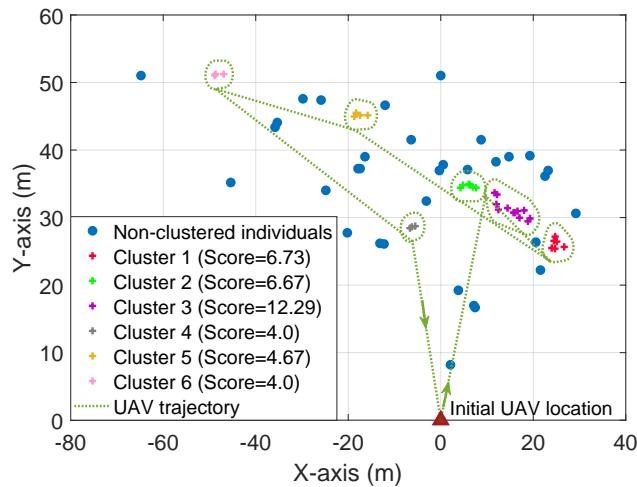


**FIGURE 13.** Output of the crowd monitoring framework (image 0000011_00234_d_0000001 of VisDrone2019 dataset).

only look once: Unified, real-time object detection. In Proc. IEEE Conf. Computer Vision Pattern Recognition (CVPR), pages 779–788, 2016.

[8] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. YOLOv4: optimal speed and accuracy of object detection. ArXiv:2004.10934, Apr. 2020.

[9] Dheeraj Reddy Pailla, Varghese Kollerathu, and Sai Saketh Chennamsetty. Object detection on aerial imagery using CenterNet. ArXiv:1908.08244, Aug. 2019.

[10] Ziyang Tang, Xiang Liu, and Baijian Yang. PENet: object detection using points estimation in high definition aerial images. In Proc. IEEE Int. Conf. Machine Learn. Appl. (ICMLA), pages 392–398, 2020.

[11] Berat Mert Albaba and Sedat Ozer. SyNet: an ensemble network for object detection in UAV images. In Proc. Int. Conf. Pattern Recogn. (ICPR), pages 10227–10234, 2021.

[12] Weiping Yu, Taojiannan Yang, and Chen Chen. Towards resolving the challenge of long-tail distribution in UAV images for object detection. In

Proc. IEEE Winter Conf. Appl. Computer Vision (WACV), pages 3257–3266, 2021.

[13] Joko Siswantoro, Anton Satria Prabuwono, and Azizi Abdullah. Real world coordinate from image coordinate using single calibrated camera based on analytic geometry. In Shahrul Azman Noah, Azizi Abdullah, Haslina Arshad, Azuraliza Abu Bakar, Zulaiha Ali Othman, Shahnorbanun Sahran, Nazlia Omar, and Zalinda Othman, editors, Soft Comput. Appl. Intelli. Syst., pages 1–11, Berlin, Heidelberg, 2013. Springer Berlin Heidelberg.

[14] Adam Babinec and Jiří Apeltauer. On accuracy of position estimation from aerial imagery captured by low-flying UAVs. Int. J. Transport. Science Technol., 5(3):152–166, Oct. 2016.

[15] Liang Zhang, Abdulkadir Celik, Shuping Dang, and Basem Shihada. Energy-efficient trajectory optimization for UAV-assisted IoT networks. IEEE Trans. Mob. Comput., pages 1–1, 2021.

[16] Marceau Coupechoux, Jérôme Darbon, Jean-Marc Kélif, and Marc Sigelle. Optimal trajectories of a UAV base station using Hamilton-Jacobi equations. ArXiv:2102.02632, Feb. 2021.

[17] Emmanouil Fountoulakis, Georgios S. Paschos, and Nikolaos Pappas. UAV trajectory optimization for time constrained applications. IEEE Network. Lett., 2(3):136–139, Jul. 2020.

[18] Oussama Ghdiri, Wael Jaafar, Safwan Alfattani, Jihene Ben Abderrazak, and Halim Yanikomeroglu. Offline and online UAV-enabled data collection in time-constrained IoT networks. IEEE Trans. Green Commun. Network., pages 1–1, Aug. 2021.

[19] Pengfei Zhu, Longyin Wen, Xiao Bian, Ling Haibin, and Qinghua Hu. Vision meets drones: A challenge. ArXiv:1804.07437, Apr. 2018.

[20] Chien-Yao Wang, Alexey Bochkovskiy, and Hong-Yuan Mark Liao. Scaled-YOLOv4: scaling cross stage partial network. In Proc. IEEE Conf. Computer Vision Pattern Recogn. (CVPR), pages 13029–13038, 2021.

[21] Alexey Bochkovskiy. Darknet. https://github.com/AlexeyAB/darknet. [Online; accessed 19-July-2020].

[22] Google. Colaboratory. https://research.google.com/colaboratory/. [Online; accessed 19-July-2020].

[23] T. Kanungo, D.M. Mount, N.S. Netanyahu, C.D. Piatko, R. Silverman, and A.Y. Wu. An efficient k-means clustering algorithm: analysis and implementation. IEEE Trans. Pattern Anal. Machine Intelli., 24(7):881–892, Aug. 2002.

[24] Zhaohui Zheng, Ping Wang, Wei Liu, Jinze Li, Rongguang Ye, and Dongwei Ren. Distance-IoU loss: Faster and better learning for bounding box regression. In Proc. AAAI Conf. Artificial Intelli., pages 12993–13000, 2020.

[25] Aleksandar Botev, Guy Lever, and David Barber. Nesterov's accelerated gradient and momentum as approximations to regularised update descent. In Proc. Int. Joint Conf. Neural Net. (IJCNN), pages 1899–1903, 2017.

[26] Peter Sturm. Pinhole Camera Model, pages 610–613. Springer US, Boston, MA, Apr. 2014.

[27] Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. In Proc. Int. Conf. Knowledge Discov. Data Mining, page 226–231, Aug. 1996.

[28] Martin Grötschel and Olaf Holland. Solution of large-scale traveling salesman problems. Math. Program., 51:141–202, Jul. 1991.

[29] G. A. Croes. A method for solving traveling-salesman problems. Operations Research, 6(6):791–812, Dec. 1958.

[30] David E. Goldberg. Genetic Algorithms in Search, Optimization and Machine Learning. Addison-Wesley Longman Publishing Co., Inc., USA, 1st edition, Oct. 1989.

[31] A. Colorni, M. Dorigo, V. Maniezzo, F. Varela, and P. Bourgine. Distributed optimization by ant colonies. In Proc. European Conf. Artificial Life, pages 134–142, Jan. 1992.

[32] R.A. Jarvis. On the identification of the convex hull of a finite set of points in the plane. Info. Process. Lett., 2(1):18–21, Mar. 1973.

[33] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. Yolov4: Optimal speed and accuracy of object detection. CoRR, abs/2004.10934, April 2020.