

# Multi-Perspective Object Detection for Remote Criminal Analysis Using Drones

Pompílio Araújo<sup>1</sup>, Jefferson Fontinele, and Luciano Oliveira<sup>2</sup>

**Abstract**—When a crime is committed, the associated site must be preserved and reviewed by a criminal expert. Some tools are commonly used to ensure the total registration of the crime scene with a minimal human interference. As a novel tool, we propose, here, an intelligent system that remotely recognizes and localizes objects considered as important evidences at a crime scene. Starting from a general viewpoint of the scene, a drone system defines trajectories through which the aerial vehicle performs a detailed search to record evidences. A multi-perspective detection approach is introduced by analyzing several images of the same object in order to improve the reliability of the object recognition. To the best of authors' knowledge, it is the first work on remote autonomous sensing of crime scenes. Experiments demonstrated an accuracy increase of 18.2% points, when using multi-perspective detection.

**Index Terms**—Criminal scene investigation, intelligent drones, multi-perspective object detection, simultaneous localization and mapping (SLAM).

## I. INTRODUCTION

IN A scene where a crime is committed, evidences are scattered nearby and should be recorded and collected by a team of experts [1]. Evidence is not perennial, decreasing in quantity and quality over the time. The goal of collecting and recording evidences is to preserve the maximum amount of information so that experts, prosecutors, and judges can analyze the dynamics of the facts, deciding, in court, the culpability of those involved. Tools to automatize the process of collecting evidence are essential to speed up the time to solve crimes, with a minimal human interference [2].

In this letter, a novel method to increase the accuracy of object detection based on multiple perspectives is introduced as a tool for automatic detection of objects in a crime scene. To do that, we use our AirCSI system [3] in a drone equipped with stereo and monocular cameras. The stereo camera is used to provide the aircraft with a global positioning system in real time by exploiting our simultaneous localization and mapping method (Air-SSLAM) [4]. In turn, the downward-facing monocular camera is used to detect and help estimating

the coordinates of the objects in the scene. Although AirCSI can use any baseline object detector to recognize objects, in our experiments, Yolo-v3 [5] was specially trained for our purposes.

## II. OUTLINE OF OUR SYSTEM

As AirCSI initiates, a coordinate system is defined to provide the drone with a starting point. Fig. 1 summarizes our proposed system described in five steps, as follows.

- 1) *Initialization*: The drone initiates the movement in the vertical direction, stabilizing at the height  $h < h_{max}$ .
- 2) *Object Detection*: Using the monocular camera at the bottom of the drone, each detected object is classified as a type of evidence, which has a relevance coefficient  $\rho$  defined by the user.
- 3) Trajectory calculation is performed according to the coefficient  $\rho$  of the detected evidence. A coverage radius is created for each detected evidence, while the drone passes through the scene.
- 4) A control module is in charge of drone stabilization and displacement of the aircraft in the trajectories defined by the system. There are eight proportional-integral-derivative (PID) controllers: Two cascades in each direction of the quadrotor drone movements, one for velocity and other one for the position.
- 5) *Multi-Perspective Detection and Report*: From the object images collected by the detector during drone trajectory, the multi-perspective detection is performed in order to provide a more accurate report with the localized evidence (sketch, evidence list, and evidence images).

## III. SELF-LOCALIZATION AND DRONE CONTROL

To estimate the drone pose, our Air-SSLAM system [4] is used. The keypoints of the two views from the stereo camera are matched in order to calculate the transformation matrix between the two consecutive frames. Therefore, Air-SSLAM performs a periodic map maintenance around image patches, which are also used as quality indicators to improve the estimated location of the drone.

Our proposed system considers six degrees of freedom that determines the pose of the drone  $[x y z \phi \chi \psi]^T$ , where  $x$ ,  $y$ , and  $z$  are the coordinates of the drone position, and  $\phi$  is the yaw rotation. The values of the angles  $\chi$  and  $\psi$  are considered null when the drone is in equilibrium, while these values are very low during the drone movement. These

Manuscript received April 29, 2019; revised July 12, 2019 and August 14, 2019; accepted September 5, 2019. This work was supported in part by Fundação de Amparo à Pesquisa do Estado da Bahia (FAPESB) under Grant 7594/2015, in part by the Federal Police of Brazil, and in part by the Coordenação de Aperfeiçoamento de Pessoal de Nível Superior–Brasil (CAPES)–Finance Code 001. The work of L. Oliveira was supported by CNPq under Grant 307550/2018-4. (Corresponding author: Pompílio Araújo.)

The authors are with the IvisionLab, Federal University of Bahia, Salvador 40170-110, Brazil (e-mail: pjaraujojunior@gmail.com).

Color versions of one or more of the figures in this letter are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/LGRS.2019.2940546

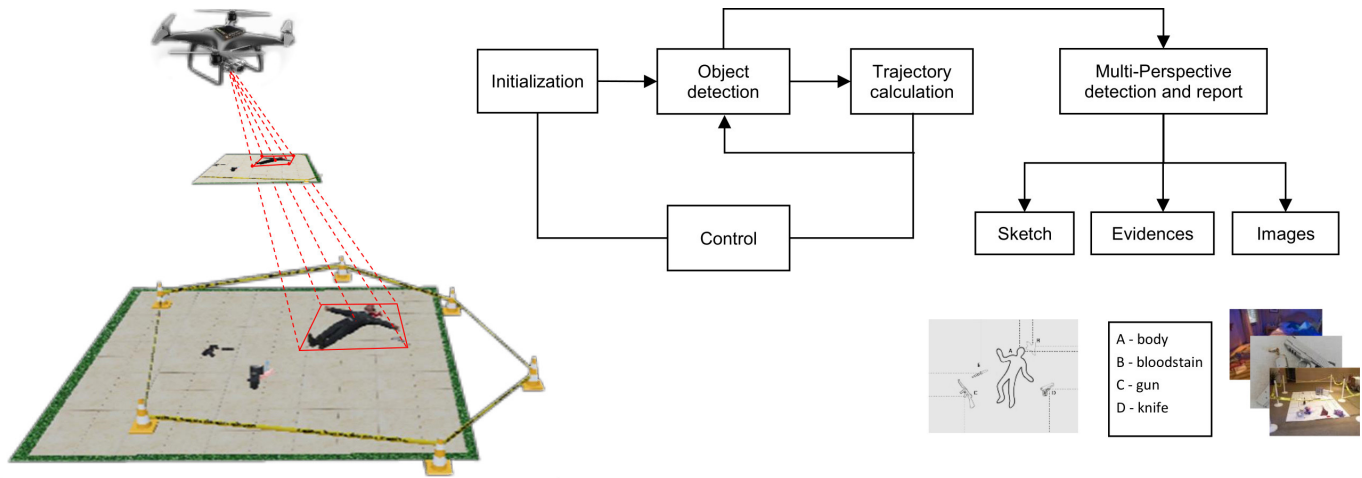


Fig. 1. Initialization—drone takes off from a position inside/near the crime scene, being positioned at a height  $h$ . Object detection—the monocular camera is used to detect suspicious objects. Trajectory calculation—with all the objects detected, a trajectory is calculated for each one of the detected objects and their locations (eventually, new objects can also be detected). Multi-perspective detection and report—the result of the scan is presented in a report. On the left, the five points of the bounding box are translated to the world coordinate system by multiplying the target vector by the inverse of the pose matrix.

constraints are completely suitable, because the drone moves at very low velocity. A double control is applied in each direction of the drone coordinates  $[x y z \psi]$ . For each variable, two controllers are used: one for velocity and another for position. The use of two controllers reduces interference from fast position variations while ensuring more efficient velocity control [6]. This improves the response of the velocity in the primary mesh. The input of the controller  $C_V$  is the velocity error  $e_{R_c}$ , given by

$$e_{\dot{x}_c} = \left( \frac{x_{c(n)} - x_{c(n-1)}}{T} \right) - P^{-1} \dot{x}_{ws} \quad (1)$$

where  $x_{c(n)}$  and  $x_{c(n-1)}$  are the positions of the drone in the camera coordinate system in the current and previous frames, respectively;  $T$  is the sampling period,  $P$  is the drone pose matrix, and  $\dot{x}_{ws}$  is the reference velocity in the global coordinate system that is received from the position controller output. The input of the controller  $C_P$  is the position error  $e_{x_w}$ , which is defined by

$$e_{x_w} = X_w - X_{ws} \quad (2)$$

where  $X_w$  is the position of the current drone and  $X_{ws}$  is the desired position. The PID controllers are used by the transfer function

$$u(t) = K_p e(t) + K_i \int e(t) dt + K_d \frac{de}{dx} \quad (3)$$

where  $K_p$ ,  $K_i$ , and  $K_d$  are the proportional, integral, and derivative constants, respectively. The 2p2z method [7] was implemented with sampling period of 160 ms. The output is given by

$$y[n] = e[n]b_0 + e[n-1]b_1 + e[n-2]b_2 + y[n-1] \quad (4)$$

where  $y[n]$  is the control signal at the output of the controller, and  $e[n]$  is the error in the controlled variable

(position or velocity). The constants  $b_0$ ,  $b_1$ , and  $b_2$  are

$$\begin{aligned} b_0 &= K_p + \frac{K_i \cdot T}{2} + \frac{K_d}{T} \\ b_1 &= K_p + \frac{K_i \cdot T}{2} - \frac{2 \cdot K_d}{T} \\ b_2 &= \frac{K_d}{T}. \end{aligned} \quad (5)$$

The controllers were tuned by the Ziegler–Nichols closed-loop method [7], adjusting the set point with a variation of 8 m in each  $x$ -,  $y$ -, and  $z$ -directions, and a variation of  $90^\circ$  in the angle yaw ( $\psi$ ). To perform experimental tests, we use the AirSim simulator [8]. AirSim is a simulator created on unreal engine that offers physically and visually realistic simulations designed to operate on high-frequency real-time looping hardware simulations. AirSim was experimentally tested with a quadrotor as a stand-alone vehicle, comparing the software components with real-world flights. In the simulator, another computer runs the AirSim program, which transmits pose information to the on-board computer (NVIDIA Jetson TX2 [9]) in the drone.

#### IV. MULTI-PERSPECTIVE DETECTION

To internally represent a detected evidence, five points (the four vertices and the geometric center of the rectangle) of the detected bounding boxes are used. Each point ( $p_i$ ) in the bounding box is defined with the coordinates  $X_{Ci} = [x_i y_i z_i]^T$  with respect to the camera coordinate system. Using this camera pose  $P$ , points are translated to the world coordinate system  $X_{wi}$

$$X_{wi} = P^{-1} \cdot X_{Ci}. \quad (6)$$

Each time an evidence is found, its position is stored. Its location is compared to that of all others stored. If the bounding box area matches an intersection over union more than 40% ( $IoU > 0.4$ ), a counter is incremented for that evidence. At the end of the scanning, each evidence will have

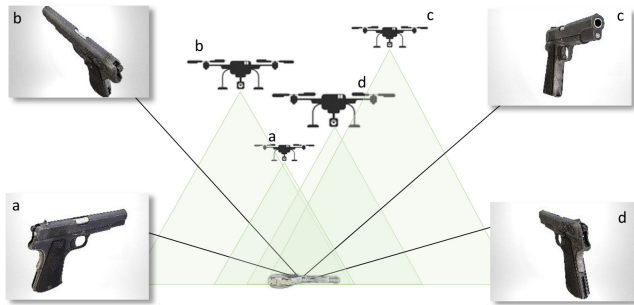


Fig. 2. Drone scans the area and objects are viewed from various perspectives. The multiple perspectives are used to improve the detection accuracy.

recorded the number of times it was detected. Then, there will be more than one perspective detection for the same evidence. During scanning, the drone camera repeatedly frames the same evidence in different perspectives. With the position of evidence recorded at the first detection, it is possible to know how many times an object was detected, as well as its detection parameters. A precision indicator (PI) is calculated, based on the number of times the object was detected, as follows:

$$PI = \frac{1}{N} \sum_{i=1}^n Cs \quad (7)$$

where  $Cs$  is the confidence score of each image provided by the baseline detector,  $N$  is the number of frames that an object should be detected, and  $n$  is the number of objects detected. Then, a PI for each evidence is given, taking into account all the available images of that evidence. In order to evaluate the proposed method, a value of Intersection over Union (IoU) is considered with respect to that evidence as the average of IoU of all images.

After detecting the objects in the scene, the object bounding box is projected onto the ground plane, providing 2-D information of the object location (as illustrated on the left of Fig. 1).

Fig. 2 shows how the drone in different positions can show the object detected in various perspectives. Position variation allows cameras to show parts of the object that could not be viewed in a single perspective.

#### A. Ablation Study

Yolo-v3 was used as a baseline detector for our multi-perspective approach. Although this detection method is not actually the most accurate nowadays, it is one of the fastest. Indeed, this issue was already shown in the work found in [5], where Yolo-v3 shows the fastest performance at the cost of presenting a lower average precision. Therefore, Yolo-v3 was considered the best choice, since precision has less relevance than the detection rate. This is so because the object will be detected from more than one perspective, and only objects that were detected more than once in all perspectives will be considered. In other words, after the first detection, object position is recorded, demanding the drone to detect the object again, in a different pose. This situation makes the object detection module to have higher precision as the drone approaches the object.

TABLE I  
COMPARATIVE EVALUATION YOLO-V3 WITH BACKBONES  
DARKNET-53 AND MOBILENET-V2. ROUNDED AVERAGE  
PRECISION (IN PERCENTAGE) FOR AP50

Backbone	Perspective								
	1	2	3	4	5	6	7	8	9
Darknet-53	35	43	45	47	48	50	51	52	53
Mobilenet-v2	31	37	38	39	40	41	45	46	50

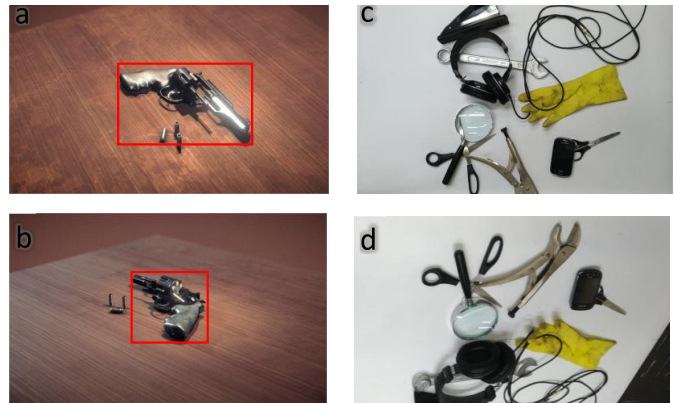


Fig. 3. (a) and (b) Two perspectives of a weapon image. (c) and (d) Two perspectives of scenes with no weapons.

The following parameters were used to train Yolo-v3: batch size = 64, momentum = 0.9, and decay = 0.0005. Images were preprocessed by changing their resolutions to  $608 \times 608$  pixels from the original images acquired. To train the detector, MS-COCO data set [10] with 3000 additional weapon images were used [11]. In order to evaluate Yolo-v3 with different backbones, we considered the original used Darknet-53 and Mobilenet-v2 [12], both made on the Jetson TX2 machine. Table I summarizes the results found.

#### V. EXPERIMENTAL ANALYSIS

MS-COCO data set was used only for training and validation, along with the 3000 additional weapon images, allowing for a model with an extensive number of categories in the future. Since there is no multi-perspective images of objects in MS-COCO data set, other 900 images containing 100 scenes in nine different perspectives were used to test the proposed system. To evaluate the proposed system, only weapon images were used housed in 50 scenes containing annotated objects and other 50 scenes considering only objects other than weapons.

Examples of scenes containing a weapon are illustrated in Fig. 3(a) and (b), while Fig. 3(c) and (d) depicts scenes without weapons. Images were submitted to the detector individually. With the object position given by Air-SSLAM, the proposed multi-perspective approach was carried out by considering the number of perspectives of each object. Considering the average precision with IoU greater than 0.5 (AP50 [13]) was possible to verify an increase of 18.2% points, going from 34.7% (with one perspective)



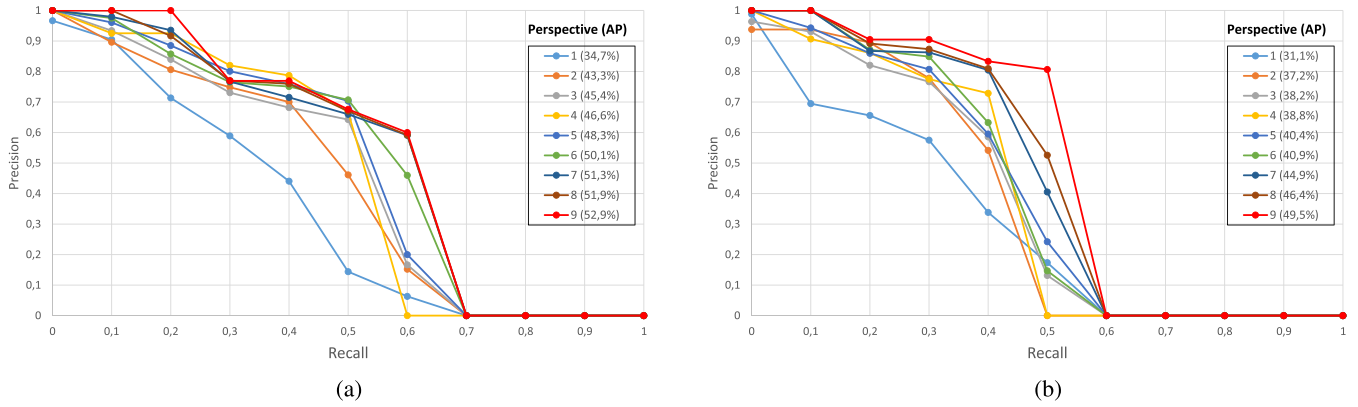


Fig. 4. Precision-recall curve of our proposed multi-detection system using Yolo-v3 with (a) Darknet-53 and (b) Mobilenet-v2 backbones. In both cases, the average precision AP50 increases with the number of perspectives.

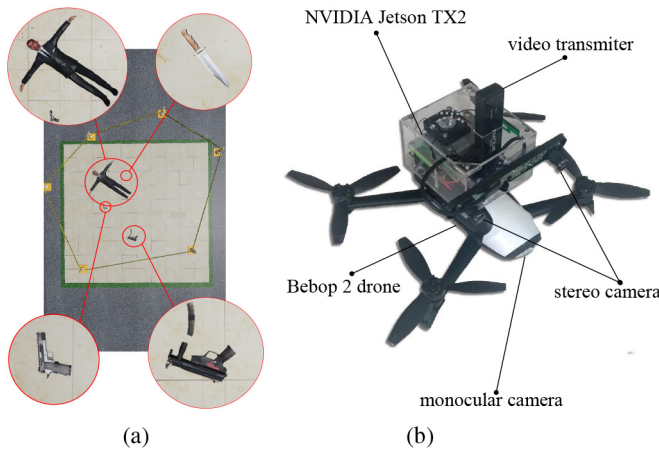


Fig. 5. (a) AirSim software simulation of a crime scene with objects detected by AirCSI. (b) AirCSI being prepared for future testing with Parrot Bebop 2 drone.

to 52.9% (with nine perspectives). Fig. 4(a) shows the AP50 plots of the proposed multi-perspective system using Darknet-53 backbone, and Fig. 4(b) illustrates the results with Mobilenet-v2 backbone. In the tests, the average detection time per image was 0.704 s for the network with the backbone Darknet-53 and 1736 s for the network with the Mobilenet-V2 backbone. An implementation of the Mobilenet-v2 backbone found in [14] and the author's implementation of Darknet-53 [5] were used.

## VI. DISCUSSION AND CONCLUSION

Although the system proposed, here, uses a drone to sweep scenes with criminal evidences (particular objects), it could also be applied to monitoring difficult areas, such as archeological parks, caves, or sites covered by dense vegetation. In searching for crime evidences, a low false-negative value is desired, since a human analysis will always be done by a specialist after the automatic search. In this sense, our proposed approach based on multiple perspective detection improved overall system accuracy successfully. In our experiments, a raise of 18.2% points in the average precision was achieved in comparison with just one perspective. The goal is to make AirCSI autonomous to detect evidences at a crime

scene [see Fig. 5(a), for example, in AirSim simulator]. In tests in the simulated environment, our drone was able to perform route calculation and detection of other objects, such as human bodies, knives, and weapons, as well as other objects present in COCO data set. We are now working on an assembly to perform testing on a Bebop 2 type drone [see Fig. 5(b)]. In the future, our challenge is an approach addressing noise and occlusion of evidence.

## REFERENCES

- [1] J. T. Fish, L. S. Miller, M. C. Braswell, and E. W. Wallace, Jr, *Crime Scene Investigation*. Abingdon, U.K.: Routledge, 2013.
- [2] M. Lega, C. Ferrara, G. Persechino, and P. Bishop, "Remote sensing in environmental police investigations: Aerial platforms and an innovative application of thermography to detect several illegal activities," *Environ. Monit. Assessment*, vol. 186, no. 12, pp. 8291–8301, 2014.
- [3] P. Araújo, M. Mendonça, and L. Oliveira, "AirCSI—remotely criminal investigator," in *Proc. Int. Conf. Adv. Signal Process. Artif. Intell.*, 2019, pp. 58–63.
- [4] P. Araújo, R. Miranda, D. Carmo, R. Alves, and L. Oliveira, "Air-SSLAM: A visual stereo indoor slam for aerial quadrotors," *IEEE Geosci. Remote Sens. Lett.*, vol. 14, no. 9, pp. 1643–1647, Sep. 2017.
- [5] J. Redmon and A. Farhadi, "Yolov3: An incremental improvement," Apr. 2018, *arXiv:1804.02767*. [Online]. Available: <https://arxiv.org/abs/1804.02767>
- [6] M. Araki and H. Taguchi, "Two-degree-of-freedom PID controllers," *Int. J. Control, Automat., Syst.*, vol. 1, no. 4, pp. 401–411, 2003.
- [7] T. E. Marlin, *Process Control: Designing Processes and Control Systems for Dynamic Performance*. New York, NY, USA: McGraw-Hill, 1995.
- [8] S. Shah, D. Dey, C. Lovett, and A. Kapoor, "AirSim: High-fidelity visual and physical simulation for autonomous vehicles," in *Field and Service Robotics* (Springer Proceedings in Advanced Robotics), vol. 5, M. Hutter and R. Siegwart, Eds. Cham, Switzerland: Springer, 2017. doi: [10.1007/978-3-319-67361-5\\_40](https://doi.org/10.1007/978-3-319-67361-5_40).
- [9] *Nvidia Autonomous Machines*. Accessed: Mar. 23, 2019. [Online]. Available: <https://www.nvidia.com/en-us/autonomous-machines/embedded-systems/jetson-tx2/>
- [10] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: Common objects in context," in *Proc. Eur. Conf. Comput. Vis.*, 2014, pp. 740–755.
- [11] F. P. Y. A. Castillo. (2018). *Weapons Detection*. [Online]. Available: <https://sci2s.ugr.es/weapons-detection>
- [12] A. G. Howard *et al.*, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," Apr. 2017, *arXiv:1704.04861*. [Online]. Available: <https://arxiv.org/abs/1704.04861>
- [13] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The Pascal visual object classes (VOC) challenge," *Int. J. Comput. Vis.*, vol. 88, no. 2, pp. 303–338, Sep. 2009.
- [14] *Mobilenet Implementation*. Accessed: 1-Jul. 1, 2019. [Online]. Available: <https://github.com/fsx950223/mobilenetv2-yolov3>