Open Access Full Text Article

Real-Time Convolutional Neural Network-Based Method for Detecting and Tracking Human Motion on Quadcopters

Quoc Duy Tran, Duc Thien Tran^{*}



Use your smartphone to scan this OR code and download this article

ABSTRACT

This paper proposes a convolutional neural network (CNN) method for human motion detection and tracking on a quadcopter. To address the challenges mentioned above, the proposed methodology is designed on computer vision techniques with an object tracking algorithm and a CNN model. The object tracking algorithm is implemented using a proportional integral differential (PID) controller to calculate the control parameters, including the pitch and yaw angles, in real time. These parameters are determined by calculating the offset between the position of the human and the camera coordinate frame. To achieve accurate object detection, a CNN model is designed based on the single shot multibox detector (SSD) architecture, which is crucial for object detection. The model above is integrated with the MobileNet base network, which is responsible for feature extraction of the object. The use of self-collected person data in model training ensures good performance for this specific application. The object detection results demonstrate that the model achieves a high level of accuracy (98%). The proposed methodology is applied to an NVIDIA Jetson NANO computer. To rigorously assess the control system, the proposed methodology was used to conduct outdoor flight tests on a campus. These tests prioritized minimal pedestrian traffic and stable weather conditions, ensuring a controlled environment for evaluation. Analysis of the flight data and signal graphs provided valuable insights into the effectiveness of the system. Key words: Human detection, SSD-MobileNet, CNN, Quadcopter, PID, Real-time processing, Embedded systems

INTRODUCTION

2 Recently, quadcopters have garnered significant at-Chi Minh City University of Technology 3 tention in various applications due to their verti-4 cal take-off and landing capabilities, as well as their ⁵ hovering functionalities¹. Additionally, quadcopters 6 can handle intricate tasks within crowded environ-7 ments and have a simpler control system than other ⁸ types of UAVs². Common applications are focused ⁹ on surveillance³, search and rescue⁴, mapping⁵, au-¹⁰ tonomous navigation⁶, obstacle avoidance⁷ and tar-¹¹ get tracking⁸. However, among the spectrum of 12 vision-based applications, object detection and track-13 ing on quadcopters present significant challenges, particularly in achieving real-time performance. Bal-14 15 ancing computational efficiency with detection accu-16 racy is crucial. Real-time operation demands fast pro-17 cessing, while high accuracy ensures reliable object 18 identification. The integration of robust vision-based 19 estimation and control algorithms is essential for ad-20 dressing these challenges and unlocking the full po-²¹ tential of quadcopters in vision-based applications. 22 In practical applications, object detection relies

on various deep learning-based algorithms, such

24 as the Faster Region-based Convolutional Neural

Network⁹, Region-based Fully Convolutional Net- 25 works¹⁰, You Only Look Once¹¹ and Single Shot De-26 tector¹². These algorithms have demonstrated remarkable capabilities in object detection tasks. However, a common challenge associated with these 29 detectors is their high computational complexity, 30 which can hinder their implementation on resource-31 constrained embedded platforms such as quad-32 copters. This limitation is particularly relevant for 33 real-time applications that demand fast processing 34 times. To address this challenge, single-shot detectors 35 have emerged as a promising approach for object detection. These detectors, such as YOLO and SSD, pro-37 cess images in a single pass, significantly reducing the computational overhead compared to two-stage de-39 tectors such as Faster R-CNN and R-FCN¹³. YOLO, 40 for instance, is renowned for its real-time processing 41 capabilities, making it suitable for applications that 42 require an immediate response¹¹. In contrast, SSD 43 strikes a balance between speed and accuracy by predicting multiple bounding boxes for each object, of-45 fering a more robust solution for tasks that demand 46 high detection accuracy^{12,14,15}. 47

Moreover, numerous control algorithms have been 48 created to address the challenges associated with 49

Cite this article : Tran Q D, Tran D T. Real-Time Convolutional Neural Network-Based Method for Detecting and Tracking Human Motion on Quadcopters. Sci. Tech. Dev. J. 2024; 27():1-15.

Department of Automation Control, Ho and Education, Vietnam

Correspondence

Duc Thien Tran, Department of Automation Control, Ho Chi Minh City University of Technology and Education, Vietnam

Email: thientd@hcmute.edu.vn

History

- Received: 2024-03-02
- Accepted: 2024-05-31
- Published Online: 2024-6-xx

DOI :

Check for updates

Copyright

© VNUHCM Press. This is an openaccess article distributed under the terms of the Creative Commons Attribution 4.0 International license.



50 tracking humans. In particular, the article presents 51 the identification and tracking of humans employ-52 ing techniques for visual data manipulation with 53 OpenCV¹⁶. In¹⁷, a fuzzy logic controller (FLC) was employed as part of a target tracking algorithm. 54 In¹⁸, they proposed a tracking algorithm grounded 55 in Euclidean space equations and image processing through cameras. While prior studies have demon-57 58 strated commendable performances, their primary focus lies within the realm of computer vision, ne-59 glecting external disturbances such as environmental 60 factors. To achieve high precision in drone control, several controllers have been applied. In¹⁹, a target-62 tracking control algorithm based on fuzzy PI was devised. This algorithm incorporates a Fuzzy-PI con-64 troller to dynamically adjust the parameters of the PI controller, utilizing positional data and changes in po-66 sition as inputs. In ²⁰, a gain-scheduled PID controller was developed to guide a UAV by continuously ad-68 justing the actuators based on real-time data from the 69 tracking unit and UAV dynamics. In²¹, a comprehen-70 sive double closed-loop proportion integral differential (PID) controller was meticulously designed, em-72 ploying estimated states to accurately track and pur-73 sue the target. Among them, PID is a promising candidate for drone control because it not only achieves 75 high accuracy but also remains robust to uncertain-76 ties from external influences²². The strengths of PID include being model-free, requiring no information about the mathematical model of the system, easy im-79 plementation on embedded boards, and high preci-80 sion²³. 81

82 This paper presents an approach for detecting and tracking target objects using an SSD object detec-83 tor on a UAV. To manage the above challenges, the 84 system is separated into two primary components: 85 (1) object motion estimation and (2) object recog-86 nition. The object motion estimation algorithm utilizes a proportional integral differential (PID) con-88 troller to compute control parameters, which include pitch and yaw angles in real time. These parameters 90 are determined based on the position of the object and are calculated by measuring the offset between 92 the position of the human and the camera coordi-93 nate frame. This module achieves robust object track-94 ing across varying relative distances. Object recog-95 nition focuses on accurately detecting "person" ob-96 97 jects using the SSD architecture. A custom-trained model differentiates between two classes: images con-98 taining objects and images without a person present. 99 Self-collected person data training enhances detec-100 tion performance. Finally, the proposed control is ap-101 102 plied to an NVIDIA Jetson NANO embedded com-¹⁰³ puter. A comprehensive outdoor flight experiment is conducted within a campus environment character-104 ized by minimal pedestrian traffic. Additionally, priority is given to selecting days with favorable weather 106 conditions and stable illumination. The analysis in-107 cludes assessing experimental flight data and signal 108 graphs to evaluate the proposed control system. 109 The remainder of this paper is structured as follows: 110 The problem statement, the object recognition algorithm and the object motion estimation algorithm are 112 described in Section II. Section III describes the ex-113 perimental analysis. Finally, Section IV offers conclusions and outlines avenues for future work. 115

116

117

147

MATERIALS AND METHODS

Preliminary

In Figure 1, the coordinate frames employed for human tracking via a quadcopter are illustrated. The system includes three coordinate frames: $O_E - x_E y_E z_E$ 120 represents the world, $O_B - x_B y_B z_B$ denotes the quad-121 copter and $O_C - x_C y_C z_C$ signifies the camera coordinates. For computational convenience, we assume 123 that the quadcopter and camera share the same co-124 ordinate frame. To address the challenge of the human motion estimation problem, the key challenge is 126 keeping the transformation matrix between the quadcopter and the human being tracked unchanged. To 128 achieve this transformation matrix, which involves 129 both orientation and position, a proposed camera sys-130 tem aims to determine both the orientation and position through the entirety of the captured image. The relative location of the human concerning the quadcopter is calculated using the camera model, which 134 is expressed as (P_R) in the camera coordinates. The 135 target coordinates (P_0) are then determined in quadcopter coordinates. The relationship between these 137 two coordinates is mathematically expressed as fol-138 lows: 139

$$\begin{bmatrix} P_B\\1 \end{bmatrix} = T_{B0}P_0 = \begin{bmatrix} R_{B0} & t_{B0}\\0 & 1 \end{bmatrix} \begin{bmatrix} P_0\\1 \end{bmatrix}$$
(1)

where R_{B0} and T_{B0} represent the matrix for rotation 140 and the matrix for transformation between the camera framework and quadcopter framework, respectively. t_{B0} denotes the position of camera²⁴. 143 Additionally, to identify human subjects from the camera output, a CNN (convolutional neural network) system is utilized for object detection. 146

Hardware Specifications

To address the challenges mentioned above, the quadcopter system comprises an executive structure and ground station control, as illustrated in Figure 2. The 150





151 ground station control is responsible for gathering
152 data from the quadcopter, while the executive struc153 ture runs the tracking and detection algorithms.

154 Control System Overview

The proposed control aims to maintain the transformation matrix between the quadcopter and the
tracked human by ensuring consistent output responses. As analyzed in Section II, this transformation matrix involves both the orientation and posi-

tion of the transformation matrix T_{B0} . This control 160 consists of two main components: vision-based estimation and object tracking control. These parts handle the detection of the targeted human and subsequent human tracking, respectively. The design of this 164 suggested control system overview is outlined in Figure 3. 166



167 Vision-Based Estimation

As illustrated in Figure 4, once an image of an objectis received, a CNN algorithm is implemented.

In this study, the SSD method relies on a feed forward convolutional network that generates a bounding box. 171 A subsequent nonmaximum suppression step is ap-172 blied to produce the final detection results 12. Fig-173 ure 4 illustrates the MobileNetSSD system, which is 174 an extension of MobileNet²⁵. However, it eliminates 175 the fully connected layers and softmax components. 176 MobileNet employs depthwise separable convolution 177 for constructing streamlined deep neural networks, 178 leading to enhancements in computational speed and model size^{26,27}. Additionally, MobileNet exhibits 180 strong performance in high-quality image classifica-181 tion tasks, contributing to its popularity in scenarios 182 where transfer learning aids in performance improve-183 ment. 184

The aim of the project is to scale the image to a size 185 of 300x300x3 and feed it into the model through 13 186 depthwise-separable convolution layers to extract the 187 feature maps, as shown in Figure 5 25. A feature layer 188 with dimensions of 10x10x1024 is selected to detect 189 objects of various sizes. The initial layers (1-5) in this 190 project are utilized for identifying typical character-191 istics present in the object image. The following lay-192 ers (from 6 onward) contain more specific informa-193 tion about the object. Next, the output of Conv_13 in the MobileNet base network is sequentially con-195 volved with a 3x3 kernel, Stride = 2, and a 1x1 kernel, 196 Stride = 1, to generate subsequent downsized feature 197 maps. The project requires a total of 6 feature maps 198 to serve as object detection layers. For every cell in the detection feature map, 4 default boxes are set up, 200 201 each having 5 distinct aspect ratios to encompass size variations. To obtain a single bounding box for a recognized object (person), the prediction box with the greatest level of confidence is selected. Any bounding boxes with an intersection over union (IoU) threshold greater than the set threshold are removed. This process is repeated until only one bounding box remains to be output. 208

Following the application of SSD-MobileNet, Figure 6²⁰⁹ depicts the presentation of a bounding box around the identified person. The positional data of the detected²¹¹ target are then extracted and employed as an input for²¹² initiating the object motion estimation algorithm to²¹³ commence the estimation process.²¹⁴

215

Control of Object Motion Estimation

Figure 7 indicates the human's position in the camera 216 coordinate system. To track a human using the entire 217 captured image, it is essential to determine the hu- 218 man's position in the coordinate framework fixed to 219 the camera. The $O_C - x_C y_C z_C$ coordinate framework 220 represents the camera coordinates. $P_0(x_0, y_0, z_0)$ rep- 221 resents the human's position at the center of the camera coordinates, where signifies the width [in pixels] 223 and represents the height [in pixels] of the entire im- 224 age. Figure 8 illustrates the connection between the 225 camera coordinates and the global coordinates. Cal- 226 culating the coordinates (y_C, z_C) is feasible because 227 the whole image is two-dimensional. However, it is 228 difficult to calculate the distance in x_C . Consequently, 229 x_C is computed as follows: 230

$$\theta_1 = \theta_0 \frac{2|y_c| + b_w}{2W} \tag{2}$$

$$x_C = \frac{2|y_C| + b_w}{2\tan\left(\theta_1\right)} \tag{3}$$

where θ_0 is the angle of view of the camera and θ_1 is ²³¹ the angle between the straight line and the z_C -axis. ²³²









233 Figure 1 illustrates the coordination frames utilized ²³⁴ for human tracking. $O_B - x_B y_B z_B$ represents the quadcopter coordinate system. Within this system, 235 v_x [m/s] denotes the translational velocities of the 236 quadcopter along the x_B -axis in $O_B - x_B y_B z_B$. Ad-237 ditionally, ψ_{z} [rad/s] signifies the angular velocity of 238 239 the quadcopter around the z_B -axis in $O_B - x_B y_B z_B$. The desired human position is designated $\bar{P}_0(\bar{x}_0)$ 240 *const*), $\bar{y}_0 (= 0)$, $\bar{z}_0 (= 0)$). In Figure 3, a block di-241 agram of position conversion (PID) concerning the 242 quadcopter velocity for human tracking is depicted. It 243 is necessary to give velocities such that $P_0(x_0, y_0, z_0)$ 244 245 comes to the center ($y_C = z_C = 0$) of images captured by the camera of the quadcopter while maintaining 246 the distance ($x_C = const$) between the quadcopter and 247 248 the human. Subsequently, the translational velocities ²⁴⁹ v_x [m/s] and the angular velocity ψ_z [rad/s], which 250 enable the quadcopter to track the human, are deter-251 mined as follows:

$$v_x = k_{px}e_x + k_{ix}\int_0^t e_x(\tau)\,d\tau + k_{dx}\frac{de_x}{dt} \tag{4}$$

$$\psi_z = k_{pz} e_z + k_{iz} \int_0^t e_z(\tau) d\tau + k_{dz} \frac{de_z}{dt}$$
(5)

where $e_x = \bar{x}_0 - x_0$ and $e_z = \bar{y}_0 - y_0$ are the errors between the position of the human in the center of the camera coordinate frame and the desired human position. When the human is undetected in the captured images, the values of v_x [m/s] and ψ_z [rad/s] are both set to zero. The quadcopter continues human tracking until a terminal command signal is received. The proposed method effectively enables quadcopters to track humans. 260

The object tracking algorithm is shown in Algorithm 261 1. The algorithm takes as input from the image of a 262 person. Following initialization, the quadcopter un- 263 dergoes a series of checks to ensure safe and reliable 264 operation. This initialization phase might involve cal- 265 ibrating sensors, verifying battery levels, and confirm- 266 ing proper motor function. Once it is given the all- 267 clear, the quadcopter autonomously ascends to a pre- 268 determined altitude. This chosen altitude offers a suitable vantage point for the search mission, allowing the 270 camera to capture a wider field of view and potentially 271 increasing the chance of human detection. The quad- 272 copter then starts on a 36-second search mission for 273 a human target. It continuously scans the environ- 274 ment using the SSD-MobileNet model. Upon success- 275 ful detection, the center offset method is used to track 276 the target by calculating the offset between the person 277 and the center of the image captured by the camera. If 278 the offset exceeds zero and the image center lies out- 279 side the bounding box, the quadcopter rotates accord- 280 ingly; otherwise, it moves forward and backward. In 281 the absence of human detection within a designated 282 timeframe, the system assumes that the target is no 283 longer present. To optimize the search efficiency, the 284 quadcopter performs a preprogrammed 10-degree ro- 285 tation, expanding the search area and increasing the 286





²⁸⁷ probability of detection. This iterative process of scanning, tracking (if detected), and rotating continues for
²⁸⁸ a total of 36 seconds. If no human is detected through²⁹⁰ out this period, prioritizing safety, the system auto²⁹¹ matically initiates a landing sequence, returning the
²⁹² quadcopter to the ground.

293 RESULTS AND DISCUSSION

²⁹⁴ To further assess the benefits of the suggested control,
²⁹⁵ a series of experiments and evaluations on an actual
²⁹⁶ system are carried out.

297 Experiment description

²⁹⁸ Figure 9 illustrates the basic movements of the quad²⁹⁹ copter during object detection and tracking. We con³⁰⁰ ducted a series of experiments to quantitatively evalu³⁰¹ ate the algorithm's performance on real hardware. We

utilized an NVIDIA Jetson NANO embedded com- 302 puter for this purpose. The algorithm was imple- 303 mented in Python within the Ubuntu Linux environ- 304 ment. The experiments were carried out outdoors on 305 the HCMUTE campus. To minimize the presence of 306 multiple objects in the scene, we chose a location with 307 minimal pedestrian traffic. Additionally, favorable 308 weather conditions were ensured to obtain accurate 309 evaluation results. The experiments and results were 310 divided into three parts. First, we evaluated the post- 311 training data to assess the algorithm's ability to detect 312 humans accurately using metrics such as precision, re- 313 call, and F1-score. Second, the flight data evaluation 314 focused on system stability and tracking performance. 315 This involved assessing the quadcopter's stability dur- 316 ing takeoff, hovering, landing, and directional move- 317 ments (forward, backward, and rotational). Finally, 318

Table 1: Algorithm 1: Object Tracking Algorithm 1.

. 1

-

Algorithm 1: Object Tracking Algorithm
input: Image person
outputs : v_x and ψ_z
begin
/* Initialize */
Sensor calibration, battery level verification, motor confirmation
Take off quadcopter
while (within 36 seconds)
Detect human using SSD-MobileNet
if (objects) then
Calculate the center of the frame, the person P_0
Calculate the offset between the person and the frame (e_x, e_y)
if $(offset > 0) & (not centered)$ then
Calculate PID control for rotation ψ_z
Send the rotation control command
end
else
Calculate PID control for forward, backward v_x
Send the forward, backward control command
end
end
Rotation by an angle of 10 degrees
end while
Landing
end

319 the data are evaluated when combining object detec-320 tion and object tracking.

Experimental results 321

CNN Training 322

Figure 10 illustrates the process of collecting and 323 preparing data for model training. 324

This study employed a single shot detector (SSD) im-325 plemented on a powerful processing unit for human 326 detection on a quadcopter. The SSD model was specif-327 ically trained to recognize a single class: individu-328 als (persons). To train and evaluate this model effec-329 tively, we constructed a comprehensive image dataset 330 containing two distinct categories: images with ob-331 jects (primarily featuring individuals) and images de-332 void of objects. The images were carefully curated 333 to ensure their suitability for real-world applications 334 involving human detection in a quadcopter environ-335 ment. The image acquisition process involved cap-336 turing video footage from the quadcopter's camera. 337 The footpad showcased a diverse range of human sub-338 jects, including group members and other individuals 339 within the research laboratory. This footpad was then painstakingly segmented into individual frames, re-341 ³⁴² sulting in a raw dataset of approximately 1000 images. To augment the dataset and enhance its learning po- 343 tential, we employed data augmentation techniques. 344 Redundant images were removed, and a subset of im- 345 ages was transformed using basic manipulations (ro- 346 tation, scaling, flipping, and brightness) to introduce 347 variations, enrich the dataset and promote model gen- 348 eralizability. Figure 11 shows the process of labeling 349 the data from the dataset. 350

The dataset comprised a total of 1000 images, main- 351 taining a 3:1 ratio between images with and without 352 objects. Each image featuring a person was meticu- 353 lously labeled for accurate object identification during training. Subsequently, these images were di- 355 vided into three distinct sets-training (70%), valida- 356 tion (20%), and testing (10%)-for network training 357 and evaluation. The network configuration included 358 a dropout ratio of 0.7, a kernel size of 3x3, a box code 359 size of 4, and a learning rate of 0.001. The training 360 process was conducted through 200 iterations using 361 Google Colab. Figure 12 illustrates the model's out-362 comes after completion of the training process. To assess how well the proposed object detection 364

method performs on an embedded computer, exper- 365 iments were conducted using the confusion matrix 366 method. The experiments were conducted 50 times 367 and included both positive and negative person in- 368







Figure 11: Images from the dataset are labeled.



³⁶⁹ stances. These experiments yielded the following met-³⁷⁰ rics: precision = 0.96078, recall = 0.98, and F1 = ³⁷¹ 0.9703. Figure 13 illustrates the results of the train-³⁷² ing model.

Additionally, the object detection process analyzed a
frame and generated an output for the detected object
within a time span of 5 ms. During this 5 ms interval,
frames captured by the quadcopter's camera underwent processing, and the CNN provided the output
in the form of an image featuring the bounding box
around the identified object, achieving a frame rate of
27 frames per second (FPS), as shown in Figure 14.

381 Altitude Control

To test the performance of the proposed controller, an 382 383 altitude experiment is first carried out. The objective of this altitude test is to control the quadcopter to take 384 off vertically to a desired altitude of approximately 2.5 385 meters and maintain that altitude for approximately 386 100 seconds before landing. In Figure 15, the initial 387 altitude (Z) is set to zero because the altitude of the quadcopter remains fixed at 2.5 meters upon takeoff. 389 After removal, the quadcopter hovers at this fixed al-390 titude (Z). 391

392 Forward, reverse motion control

³⁹³ Following the altitude experiment, rotation and ³⁹⁴ forward-backward experiments are conducted. The ³⁹⁵ aim of this experiment is to control the aircraft to ³⁹⁶ move at predetermined speeds and angles. The pre-³⁹⁷ set speed is 1 m/s, and the rotation angle is set to 90 degrees; this process is repeated three times within a39860-second flight time. The resulting data are represented as squares in Figure 16.399

401

Combined Control

After conducting two flight experiments involving 402 tracking in the forward, backward, and object rota- 403 tion directions, the goal is for the quadcopter to de- 404 tect objects within the frame and simultaneously per- 405 form forward-backward movement and object track- 406 ing. Figure 17 shows the real-world object tracking 407 experiment. The detected object will move freely to 408 verify the accuracy of the system. The validation flight 409 process took place over approximately 300 seconds. Based on the signals from the graphs, we can observe 411 the aircraft's status during the tracking process in the 412 forward, backward, and rotation directions. The roll 413 angle is approximately equal to 0. In the yaw angle re- 414 sponse graph, the aircraft rotates from approximately 415 180 degrees to 0 degrees within 80 seconds, from the 416 50th to the 130th second, after which it moves north- 417 ward. During this time, the yaw angle experiences 418 only slight rotation in the north direction. This in- 419 dicates that the quadcopter tracks the detected object 420 relatively well. From the velocity response graph in 421 the x-direction, it is evident that the aircraft's velocity 422 in the x-direction is very low, indicating slow forward 423 movement. However, it still responds effectively to 424 track the object. Furthermore, the engine pulse out- 425 put graph shows that the engines pulse continuously 426 when the aircraft is in a combined state. Pulse gener- 427



Figure 13: Output of the CNN model



Figure 14: Evaluate the processing speed of the model

428 ation during takeoff and landing is very fast, demon-429 strating stable takeoff and landing. The resulting data430 are represented as squares in Figure 18.

While our object detection and tracking algorithm
demonstrated promising results, its real-world implementation presented unforeseen hurdles. A significant challenge arose from the delayed response
data received from the quadcopter. This latency, attributed to the limitations of Bluetooth data transmission, created a disadvantage in the real-time process-

ing pipeline. Furthermore, the hardware of onboard438cameras occasionally hinders the ability of SSD ob-439ject detectors to consistently identify target objects.440This limitation was particularly evident under varying441lighting conditions, where real-time object detection442proved challenging.443

CONCLUSION

This study presents a novel approach for human mo- 445 tion detection and tracking on a quadcopter, lever- 446

11

444



aging the power of convolutional neural networks 447 (CNNs). The proposed system, implemented on an embedded computer, comprises two key compo-449 nents: object recognition and object motion estima-450 tion. The object recognition module employs a CNN-451 based SSD model to identify moving objects within 452 the camera's field of view. This model effectively gen-453 erates bounding boxes around detected objects, ex-454 tracting their center positions for precise tracking. 455 Simultaneously, the object motion estimation mod-456 ule, powered by a PID controller, dynamically ad-457 justs the quadcopter's flight path to pursue the tar-458 get object even under varying speeds. The experi-459 mental results demonstrate the impressive capabili-460 ties of the system. The object recognition algorithm 461 boasts high accuracy in object detection and catego-462 rization while maintaining low power consumption 463 and achieving a high frame rate (fps). However, real-464 time implementation has revealed limitations asso-465 ciated with communication latency due to Bluetooth 466 data transmission and onboard camera hardware con-467 straints. These limitations manifested as occasional 469 delays in receiving data and hindered object detection ⁴⁷⁰ accuracy under varying lighting conditions.

In the future, this work paves the way for further 471 advancements. Integrating vision-based techniques 472 with a stereo camera to estimate the distance between 473 the quadcopter and the target object has emerged as 474 a crucial area for future research and development. 475 This advancement would enable more precise object 476 tracking and navigation, particularly in complex en- 477 vironments. Additionally, the focus will shift toward 478 developing more sophisticated algorithms for han- 479 dling multiple objects. By incorporating techniques 480 for multiobject tracking, the system could effectively 481 track and differentiate between multiple people in 482 high-density environments. This advancement would 483 be invaluable for applications such as search and 484 rescue operations in crowded areas or autonomous 485 surveillance tasks involving multiple targets. 486

COMPETING INTERESTS	487
None	488
NONFINANCIAL COMPETING INTERESTS	489 490
None	401



Figure 16: Output responses of the forward and reverse motion control experiments



Figure 17: Object tracking experiment in the real world



492 ACKNOWLEDGMENTS

- ⁴⁹³ This research was implemented at the Robotics and
- 494 Intelligent Control Laboratory (RIC Lab), Faculty of
- ⁴⁹⁵ Electrical and Electronics Engineering, Ho Chi Minh
- 496 City University of Technology and Education, Viet-
- 497 nam.

498 AUTHOR CONTRIBUTION

⁴⁹⁹ The authors acknowledge the support of time and fa-⁵⁰⁰ cilities from the Robotics and Intelligent Control Lab-⁵⁰¹ oratory for this study.

502 **REFERENCES**

- Shakhatreh H, Sawalmeh AH, Al-Fuqaha A, Dou Z, Almaita
 E, Khalil I, et al. Unmanned aerial vehicles (UAVs): A sur-
- vey on civil applications and key research challenges. leee

Access. 2019;7:48572-634;Available from: https://doi.org/10.
 1109/ACCESS.2019.2909530.

- Deepak B, Singh P. A survey on design and development of an unmanned aerial vehicle (quadcopter). International Journal of Intelligent Unmanned Systems. 2016;4(2):70-106;Available from: https://doi.org/10.1108/IJIUS-10-2015-0012.
- Zaheer Z, Usmani A, Khan E, Qadeer MA, editors. Aerial surveillance system using UAV. 2016 thirteenth international con-
- ference on wireless and optical communications networks
 (WOCN); 2016: IEEE;Available from: https://doi.org/10.1109/
 WOCN.2016.7759885.
- Bevacqua G, Cacace J, Finzi A, Lippiello V, editors. Mixedinitiative planning and execution for multiple drones in search and rescue missions. Proceedings of the International Conference on Automated Planning and Scheduling; 2015;Available from: https://doi.org/10.1609/icaps.v25i1.13700.
- 522 5. Qin H, Meng Z, Meng W, Chen X, Sun H, Lin F, et al. Autonomous exploration and mapping system using het
 - erogeneous UAVs and UGVs in GPS-denied environments.
- 525 IEEE Transactions on Vehicular Technology. 2019;68(2):1339-
- 526 50;Available from: https://doi.org/10.1109/TVT.2018.2890416.

- Padhy RP, Xia F, Choudhury SK, Sa PK, Bakshi S. Monocular vision aided autonomous UAV navigation in indoor corridor environments. IEEE Transactions on Sustainable Computing. 2018;4(1):96-108;Available from: https://doi.org/10.1109/ TSUSC.2018.2810952.
- McGuire K, De Croon G, De Wagter C, Tuyls K, Kappen H. Efficient optical flow and stereo vision for velocity estimation and obstacle avoidance on an autonomous pocket drone. IEEE Robotics and Automation Letters. 2017;2(2):1070-6;Available 536 from: https://doi.org/10.1109/LRA.2017.2658940. 536
- Le T-L, Quynh NV, Long NK, Hong SK. Multilayer 537 interval type-2 fuzzy controller design for quadcopter unmanned aerial vehicles using Jaya algorithm. IEEE Access. 2020;8:181246-57;Available from: 540 https://doi.org/10.1109/ACCESS.2020.3028617. 541
- Faster R. Toward real-time object detection with region proposal networks. Advances in neural information processing systems. 2015;9199(10.5555):2969239-50;.
 544
- Dai J, Li Y, He K, Sun J. R-fcn: Object detection via region-based fully convolutional networks. Advances in neural information processing systems. 2016;29;.
- Redmon J, Divvala S, Girshick R, Farhadi A, editors. You only look once: Unified, real-time object detection. Proceedings of the IEEE conference on computer vision and pattern recognition; 2016;Available from: https://doi.org/10.1109/CVPR.2016.551 91.552
- Liu W, Anguelov D, Erhan D, Szegedy C, Reed S, Fu C-Y, et al., editors. Ssd: Single shot multibox detector. Computer Vision-ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part I 14; 2016: Springer;Available from: https://doi.org/10.1007/978-3-319-46448-0_2.
- Zaidi SSA, Ansari MS, Aslam A, Kanwal N, Asghar M, Lee B. 559 A survey of modern deep learning based object detection models. Digital Signal Processing. 2022;126:103514;Available from: https://doi.org/10.1016/j.dsp.2022.103514.
- Phadtare M, Choudhari V, Pedram R, Vartak S. Comparison between yolo and ssd mobile net for object detection in a surveillance drone. Int J Sci Res Eng Man. 2021;5:1-5;.
- Alkentar SM, Alsahwa B, Assalem A, Karakolla D. Practical comparation of the accuracy and speed of YOLO, SSD and 567

524

- Faster RCNN for drone detection. Journal of Engineering.
 2021;27(8):19-31;Available from: https://doi.org/10.31026/j.
 eng.2021.08.02.
- 571 16. George RP, Prakash V, editors. Real-time human detection
- and tracking using quadcopter. Intelligent Embedded Systems: Select Proceedings of ICNETS2, Volume II; 2018:
- 574 Springer;Available from: https://doi.org/10.1007/978-981-10-575 8575-8 29.
- 576 17. Pei C, Zhang J, Wang X, Zhang Q. Research of a nonlinearity
- control algorithm for UAV target tracking based on fuzzylogic systems. Microsystem Technologies. 2018;24:2237-
- 579 52;Available from: https://doi.org/10.1007/s00542-017-3641-
- 580

0

- 581 18. Chen P, Dang Y, Liang R, Zhu W, He X. Real-time object track-
- ing on a drone with multi-inertial sensing data. IEEE Transac tions on Intelligent Transportation Systems. 2017;19(1):131-
- ⁵⁸⁴ 9;Available from: https://doi.org/10.1109/TITS.2017.2750091.
- 585 19. Rabah M, Rohan A, Mohamed SA, Kim S-H. Autonomous mov-
- ing target-tracking for a UAV quadcopter based on fuzzy-PI.
 IEEE access. 2019;7:38407-19;Available from: https://doi.org/
 10.1109/ACCESS.2019.2906345.
- Rabah M, Rohan A, Haghbayan M-H, Plosila J, Kim S-H. Het erogeneous parallelization for object detection and tracking
 in UAVs. IEEE access. 2020;8:42784-93;Available from: https:
- 592 //doi.org/10.1109/ACCESS.2020.2977120.
- 593 21. Wu S, Li R, Shi Y, Liu Q. Vision-based target detection and
- tracking system for a quadcopter. IEEE Access. 2021;9:62043 54;Available from: https://doi.org/10.1109/ACCESS.2021.
 3074413.
- Praveen V, Pillai S. Modeling and simulation of quadcopter us gn PID controller. International Journal of Control Theory and
 Applications. 2016;9(15):7151-8;.
- 23. Joseph SB, Dada EG, Abidemi A, Oyewola DO, Khammas
 BM. Metaheuristic algorithms for PID controller parameters
- tuning: Review, approaches and open problems. Heliyon.
 2022;8(5);PMID: 35600459. Available from: https://doi.org/10.
- 604 1016/j.heliyon.2022.e09399.
- Luukkonen T. Modeling and control of quadcopter. Inde pendent research project in applied mathematics, Espoo.
 2011;22(22);.
- 608 25. Sandler M, Howard A, Zhu M, Zhmoginov A, Chen L-C, edi-
- 609 tors. Mobilenetv2: Inverted residuals and linear bottlenecks.
- 610 Proceedings of the IEEE conference on computer vision and
- pattern recognition; 2018;Available from: https://doi.org/10.1109/CVPR.2018.00474.
- 613 26. Howard AG, Zhu M, Chen B, Kalenichenko D, Wang W,
- 614 Weyand T, et al. Mobilenets: Efficient convolutional neu-615 ral networks for mobile vision applications. arXiv preprint
- 616 arXiv:170404861.2017;.
- Szegedy C, Reed S, Erhan D, Anguelov D, loffe S. Scalable, high quality object detection. arXiv preprint arXiv:14121441.2014;.