

Received 16 June 2023, accepted 4 July 2023, date of publication 7 July 2023, date of current version 12 July 2023.

Digital Object Identifier 10.1109/ACCESS.2023.3293124

## RESEARCH ARTICLE

# CNN-Based UAV Detection and Classification Using Sensor Fusion

HUNJE LEE<sup>1</sup>, SUJEONG HAN<sup>1</sup>, JEONG-IL BYEON<sup>1</sup>, SEOULGYU HAN<sup>1</sup>, RANGUN MYUNG<sup>1</sup>, JINGON JOUNG<sup>1,2</sup>, (Senior Member, IEEE), AND JIHOON CHOI<sup>1</sup>, (Senior Member, IEEE)

<sup>1</sup>School of Electronics and Information Engineering, Korea Aerospace University, Gyeonggi-do 10540, South Korea

<sup>2</sup>School of Electrical and Electronics Engineering, Chung-Ang University, Seoul 06974, South Korea

Corresponding author: Jihoon Choi (jihoon@kau.ac.kr)

This work was supported by the National Research Foundation of Korea (NRF) funded by the Korean Government (MSIT) under Grant 2021R1A4A2001316, Grant 2022R1F1A1073999, and Grant 2022R1A2C1003750.

**ABSTRACT** This paper proposes a detection and classification method for unmanned aerial vehicles, commonly called drones, using sensor fusion schemes. Datasets for drone detection and classification are collected by field measurements of actual drones using the optical camera, radar, and audio microphone as well as obtained from open online databases. In the first stage of the proposed method, drone detection and classification are conducted using the convolutional neural network (CNN) models separately trained by the optical images, radar range-Doppler maps, and audio spectrograms. Then, the CNN output probabilities are combined by the multinomial logistic regression to improve the drone surveillance accuracy through the fusion of the optical, radar, and audio sensors. Numerical simulations are performed with the experimental data and the open datasets. From the results, it is verified that the proposed sensor fusion method can improve the drone detection accuracy by up to 15.6% and can enhance the drone classification accuracy by up to 28.1% in terms of the F-score, compared to individual sensing schemes.

**INDEX TERMS** UAV detection, UAV classification, sensor fusion, convolutional neural network (CNN), multinomial logistic regression.

## I. INTRODUCTION

Unmanned aerial vehicles (UAVs), commonly called drones, are widely used in our lives with various applications for military missions, agriculture, entertainment, safety diagnosis, disaster relief, shipping, and wireless communications. With the rapid expansion of the drone industry, we are exposed to potential threats by drones, such as security area invasion, privacy infringement, and destructive terrors. Considering recent advances in UAV flight systems, anti-drone (or counter drone) technologies have been actively investigated to protect essential facilities and areas from accidental or intentional intrusion of drones [1], [2], [3], [4], [5], [6], [7]. Some civilian drone manufacturers have embedded geofencing software to prevent drones from flying over no-fly and flight-restricted zones, such as government buildings and airports. However,

it is formidable and challenging to apply geofencing to military-purpose drones and to enforce flight restrictions on all civilian drones. Therefore, deploying anti-drone systems for protecting security-sensitive areas is very important.

The anti-drone system requires real-time detection of drones, estimation of location, and classification of drone types (or models) to determine if the object is a threatening drone. Practically, it is challenging to detect a drone because of its small size, low flying speed, low altitude, low radar cross section (RCS), and low vibration. To overcome these difficulties, several surveillance techniques have been devised based on the video sensor, radar sensor, acoustic sensor, and radio frequency (RF) receiver. Each detection method has complementary advantages and disadvantages. Drone detection using video images is a sort of object detection problem which has been extensively studied in the field of pattern recognition and computer vision, and many research results have been reported based on image

The associate editor coordinating the review of this manuscript and approving it for publication was Aysegül Ucar<sup>1</sup>.

features such as colors, line shapes, geometric forms, and edges [8], [9], [10], [11], as well as based on motion features such as the object velocity, moving direction, and flight pattern [12], [13], [14]. Whereas optical cameras provide low-cost detection and fine-grained tracking of drones, there are shortcomings like the relatively short detection range, high sensitivity to weather conditions, and invisibility by obstacles. In an attempt to find drones under low light conditions, thermal infrared cameras detect the heat emitted from motors, batteries, and internal hardware [15], [16]. Thermal detection enables drone surveillance at night, yet the practical detection range is significantly shorter than other surveillance methods.

Although radar is commonly used for surveillance of large aircraft, it is not easy to detect drones with radar due to the limited RCS, low speed, and low altitude. With recent advancements in radar system technology, it has been possible to detect extremely small targets including drones [17]. Radar surveillance is a promising technology due to the long detection range, the high position accuracy, the weather independence, the capability for multi-target detection, and the night operability [18], [19], [20], [21], [22]. For example, the micro-Doppler signatures caused by the rotation of rotors and propellers can be used to detect and classify small drones with high accuracy [23]. Further research has been conducted to improve the detection granularity via the multi-channel passive radar [24], [25] and to provide more advanced features like high resolution and phase interferometry through the frequency-modulated continuous wave (FMCW) radar [26], [27], [28]. Despite these advantages, the use of high-power radar is strictly regulated in densely populated urban areas, and the drone detection radar necessitates high costs for installation and operation. Consequently, it is difficult to construct an anti-drone system only using radar except in the government and military areas.

Alternatively, we can exploit the sounds emitted from the rotors and propellers including inherent drone features. Acoustic drone detection can be accomplished with a single microphone [29] and multiple microphones [30] by analyzing the acoustic signatures in the time and/or frequency domains. Various techniques are jointly considered to improve the acoustic detection performance: a noise reduction technique is employed in [31]; machine learning and deep learning approaches are utilized in [29], [32], [33], and [34]; and the use of acoustic sensors equipped with drones has been investigated for target localization in [34], [35], and [36]. In [34], to mitigate the influence of noise, acoustic features are extracted by the short-time Fourier transform (STFT) in combination with convolutional neural networks (CNNs). Moreover, machine learning is applied, followed by feature extraction methods such as mel-frequency cepstral (MFC) coefficients and linear predictive cepstral coefficients in [29], and similarly, the independent vector analysis is employed for feature extraction from sounds in [37]. Acoustic sensors enable a low-cost implementation for anti-drone systems and

provide detection performance which is resilient to light and weather conditions as well as less impacted by obstacles. However, the detection range is relatively shorter than other methods (up to a few hundred meters), and the detection accuracy can be significantly degraded in the presence of background noise.

Utilizing an RF scanner is another promising technique for drone detection. RF scanning devices intercept wireless signals used to control a drone which contain various sensing data for navigation, flight commands, and so on. Commercial drones use RF signals typically in the range of 2.4 GHz to 5 GHz reserved for industrial, scientific, and medical radio bands (ISM bands), which can be detected by the RF scanner [38], [39]. As the frequency used by an illegal drone is usually unknown, an RF scanner hops among multiple frequency bands in order to find a control signal in all possible frequency ranges [40]. RF-based drone detection and identification methods can be further enhanced by using machine learning [41], [42] and deep learning [43]. The RF scanner is robust against weather conditions allowing for long-range and low-cost drone surveillance, if the frequency bands and/or the control protocols are known. However, this method has some limitations in detecting drones, if the control information is transferred without adhering to a standard communication specification or if a drone operates autonomously without communication between the drone and its controller. Also, the performance can be deteriorated by interference from other RF signals [44].

Drone detection using individual sensors reveals problems in specific scenarios due to the drawbacks of the aforementioned sensing methods. In an attempt to improve detection performance, sensor fusion technologies have been investigated [45]. The first approach for sensor fusion is to use two or more different sensors simultaneously for accurate and reliable detection. Several sensor fusion techniques are developed by combining optical and acoustic sensors in [46] and by concatenating signatures of acoustic, optical, and radar sensors in [47]. Both audio and video streams are concurrently used for drone detection by extracting features and feeding to a classifier [48]. A deep neural network (DNN) to process the RF sensing data is concatenated with a CNN to process the visual sensing data to form a combined DNN for sensor fusion [49]. The second approach is to use one sensor for acquisition and the other sensor for verification, that is, one sensor with a more extended range detects the presence of a drone, and the other sensor with higher accuracy confirms the initial detection results by adjusting the parameters such as the angle of arrival (AoA) and the zoom level of the camera. For example, a new procedure for drone detection and tracking is developed based on the fusion of daylight camera, thermal camera, and acoustic sensors in [50]. Sensor fusion can facilitate more reliable, robust, and precise drone surveillance across diverse operating scenarios, though it demands higher system complexity and deployment costs. For instance, multiple sensors need to be synchronized

in time to detect targets from the combined sensing data, and parameter optimization for joint detection is required to enhance the performance. In other words, sensor fusion methods necessitate a sophisticated design and experimental validation to accelerate the development of a practical anti-drone system.

In this paper, we consider the detection of illegal drones that utilize autonomous flight or seldom communicate with the controller. To this end, we collect experimental data for drone detection via field measurements using an optical camera, FMCW radar, and acoustic microphone. We propose a two-stage approach for drone detection and classification. In the first stage, drone detection and classification are carried out using individual sensing data and CNN models, and in the second stage, three kinds of sensing data are combined based on the multinomial logistic regression model to improve surveillance performance. The contribution of this paper is summarized as follows.

- Through field measurements at Korea Aerospace University (KAU) and Chung-Ang University (CAU), experimental sensing data are obtained for the optical image, the radar range-Doppler map, and the acoustic spectrogram. We collect experimental data on drones in flight using three types of drones with different sizes. Sensing data for non-drone objects are obtained by field measurements and also acquired from open datasets for machine learning [51], [52]. We convert the measured data into images so that drones can be detected and classified using the same type of CNN models. For this purpose, optical images are scaled to accommodate the input image size of CNN models; FMCW radar echoes are transformed to range-Doppler maps via radar signal processing; and acoustic signals collected by a microphone are used to create spectrograms through the MFC filtering and STFT.
- Two-stage deep learning models are developed corresponding to optical images, range-Doppler maps, and audio spectrograms to detect and classify drones. Using the transfer learning technique, the first-stage CNN models are adapted for drone detection with a single output node for binary classification, while the second-stage CNN models are revised for drone classification with multiple output nodes corresponding to the number of drone types.
- A new sensor fusion method is proposed based on a multinomial logistic regression model [53] to enhance the accuracy of drone detection and classification. During the training phase, the coefficients for integrating data from the three sensors are optimized using the probabilities of three CNN output nodes and the ground truth labels. During the test phase, the logistic regression models trained in this way are used for drone detection and classification based on concurrently measured data from the sensors.
- Numerical simulations are performed with the experimental data and the open datasets to evaluate the

performance of the CNN models associated with individual sensors. Moreover, the probability datasets obtained from the CNN output nodes are employed to train and test the proposed multinomial logistic regression models for sensor fusion. The results demonstrate that the proposed fusion method, incorporating data from three sensors, improves the drone detection accuracy by 2.4%~15.6% and enhances the drone classification accuracy by 8.7%~28.1% in terms of the F-score, compared to individual sensing schemes.

The organization of this paper is as follows. Section II presents the measurement setup to obtain experimental data for individual sensors using commercial drones, and Section III introduces the drone detection and classification techniques using an optical camera, FMCW radar, and acoustic microphone. In Section IV, we propose a new sensor fusion method based on the multinomial logistic regression model. Section V presents numerical results to evaluate various drone detection and classification schemes, and Section VI provides concluding remarks and future research issues.

*Notations:* Superscripts  $T$  and  $-1$  denote transposition and inversion, respectively, for any scalar  $x$ , vector  $\mathbf{x}$ , or matrix  $\mathbf{X}$ .  $\mathbf{1}$  denotes the all-ones column vector;  $\text{diag}(\mathbf{x})$  returns a diagonal matrix whose main diagonal elements are equal to  $\mathbf{x}$ ;  $\frac{\mathbf{y}}{\mathbf{x}}$  stands for elementwise division between vectors  $\mathbf{y}$  and  $\mathbf{x}$ ; and  $\frac{\partial \mathbf{y}^T}{\partial \mathbf{x}}$  means a matrix whose  $(m, n)$ th element is  $\frac{\partial y_n}{\partial x_m}$  where  $x_m$  and  $y_n$  are the  $m$ th and  $n$ th elements of  $\mathbf{x}$  and  $\mathbf{y}$ .

## II. MEASUREMENT SETUP

This section provides a description of the location and surrounding environments where the experiments were conducted, the specification of drones, and the setup of sensing devices for detection and classification of drones.

### A. LOCATIONS FOR MEASUREMENT

Experiments were mainly conducted in the Drone Airfield at Korea Aerospace University (KAU) located in Goyang-si, Republic of Korea, as shown in Fig. 1(a). The measurement data in this place is influenced by trees and grass surrounding the location. To obtain measurement data from various environments, field experiments were performed in the Futsal Field at Chung-Ang University (CAU) located in Dongjak-gu, Seoul, where the place is surrounded by buildings with more than five stories on three sides, as shown in Fig. 1(b). We collected non-UAV measurement data such as moving cars, people passing the crosswalk, and running people in the street with a 40 m width in South Seoul.

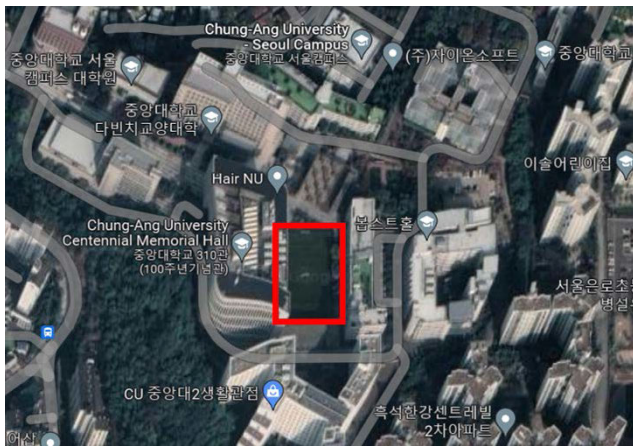
### B. DRONES FOR FIELD EXPERIMENTS

In the experiments, three types of drones are used, namely DJI Inspire2, Mavic3, and Phantom4. Each drone is different in size, shape, color, and material, as shown in Table 1. Inspire2 is the largest and heaviest in size and weight and has black and gray colors. The body is made of plastic and magnesium-aluminum alloy, and the arm is made of carbon fiber. Mavic3





(a)



(b)

FIGURE 1. Locations for field measurement: (a) Drone Airfield at KAU, (b) Futsal Field at CAU.

TABLE 1. Specifications of the drones used in experiments.

Parameter	Inspire2	Mavic3	Phantom4
Size	42.5cm(W) × 42.7cm(L) × 31.7cm(H)	34.75cm(W) × 28.3cm(L) × 10.07cm(H)	28.9cm(W) × 28.9cm(L) × 19.6cm(H)
Weight	3440g	895g	1388g
Material	Composite shell of Mg and Al, Carbon fiber	Polycarbonate, Carbon fiber reinforced nylon	Magnesium alloy
Diagonal length	65.0cm	38.1cm	35.0cm

is the thinnest and lightest and has black and gray colors. The material comprises plastic, polycarbonate and carbon fiber reinforced nylon. The body of Phantom4 has the same length and width with a medium size and weight, whose color is white. The material is made of plastic and magnesium alloy.

C. SENSING EQUIPMENT

Three types of sensors were used in the experiments: optical, radar, and acoustic sensors. For optical sensing, a camera attached to a commercial smartphone was used. As shown in



(a)



(b)



(c)

FIGURE 2. Drones used in experiments: (a) Inspire2, (b) Mavic3, (c) Phantom4.



FIGURE 3. Setup for field measurements using the smartphone camera, the FMCW radar, and the microphone.

TABLE 2. Specifications of the optical camera used in experiments.

Parameter	Value
Angle of view	77
Effective resolution	4032 × 3024
Focal length	26 mm
Size of image sensor	1/2.55 inch
Pixel size	1.4 μm
Wide dynamic range (WDR)	Smart WDR
Color filter	RGB Bayer pattern

Table 2, the focal length is 26 mm, the resolution is 12 Mpixels (the image size is 4032 × 3024), the image sensor size is 1/2.55 inch, and the size per pixel is 1.4 μm.

Table 3 presents the specifications of the radar sensor with multiple-input multiple-output (MIMO) FMCW waveforms

**TABLE 3. Specifications of the radar used in experiments.**

Parameter	Value
Waveform	FMCW
MIMO	2 TX × 4 RX antennas
RF output power	8 dBm
Antenna gain	12.6 dBi
Range Resolution	60 cm
3 dB beam width	75° in azimuth, 15° in elevation
Lower/upper frequencies	24.0/24.25 GHz
Chirp repetition interval	1 ms
Upchirp duration	512 $\mu$ s
Sampling frequency	1 MHz
Number of samples per chirp	512

**TABLE 4. Specifications of the acoustic microphone used in experiments.**

Parameter	Value
Polar pattern	Cardioid
Sensitivity	-45 dB
Sampling rate	48 kHz
Bit per sample	16 bits
Time interval per recording	5 ~ 30 sec

which are transmitted and received in the range of [24 GHz, 24.25GHz] frequency band. The radar is equipped with four receive antennas and two transmit antennas on the front. The radar can simultaneously detect the range, speed, and angle of multiple targets. The maximum distance and the velocity range can vary depending on the chirp configuration. Our experiments set the maximum distance to 80 m and the velocity range to  $[-12, 12]$  m/s. Notice that the maximum detection distance is limited due to the transmit power regulation, and a commercial anti-UAV system can increase the detection range using a high-power radar. A Cardioid microphone was used with a polar pattern as an acoustic sensor. The microphone's sensitivity is  $-45$  dB, the sampling rate is 48 kHz, and the bit depth is 16 bits, as shown in Table 4. The detection distance of the microphone depends on the background noise level caused by various sources, such as driving cars, strong wind, and talking people. In typical quiet environments, the maximum detection distance is approximately 90 – 150 m (90 m for Mavic3, 110 m for Phantom4, and 150 m for Inspire2). In a commercial anti-UAV system, microphone arrays can be used to increase the audio detection distance through acoustic beamforming [30].

As shown in Fig. 3, each sensor is attached to a tripod and placed at an identical height. In the case of the radar and microphone, the sensors are controlled by built-in softwares and the measured data are saved in Laptop1 and Laptop2 through USB cables. Since the optical sensing data is an image, no conversion procedure is required. The data obtained through the radar and the acoustic sensor are converted into images through signal processing. The echoes received from radar sensors are converted to a range-Doppler map through signal rearrangement and fast Fourier transform (FFT), and the waveform obtained by the acoustic sensor is converted into a spectrogram via the STFT. The signal processing procedure will be explained in the following section.

### III. UAV DETECTION AND CLASSIFICATION USING EACH INDIVIDUAL SENSOR

In this section, we explain the procedures for obtaining the optical images, range-Doppler maps, and spectrograms from the field measurement data obtained by the camera, FMCW radar, and microphone, respectively, and then present several example results corresponding to each sensor. Moreover, CNN models are employed for detecting and identifying drones using the optical, radar, and audio data, separately.

#### A. OPTICAL SENSING

The built-in camera of a commercial smartphone was used to obtain the optimal images of drones and non-drone objects such as helicopters, sky, surrounding buildings, background trees, and so on. The size of the image data taken through a smartphone is  $4032 \times 3024 \times 3$ . Since the drone can be operated until the sun goes down, the experiments were conducted during the daytime, and the images were taken with the sun behind. The focal length was set to 27 mm and 52 mm corresponding to the 2x zoom mode, and the continuous shooting mode of the smartphone was exploited to take as many pictures as possible. The actual drone images were obtained by taking pictures of three kinds of drones in Table 1 during flight, and the non-drone images were achieved by taking pictures of the sky, the helicopters in flight near KAU, the trees around the hill in KAU Drone Airfield, and the surrounding buildings near the CAU Futsal Field.

Moreover, additional external images were acquired for drones and non-drone objects from open datasets in [51]. The non-drone images include airplanes, warplanes, helicopters, rockets, and other objects which look like drones seen from a distance. Notice that the external datasets provide a variety of images that are difficult to obtain through measurements. We convert the size of optical images measured by the built-in camera and obtained from open datasets to  $224 \times 224 \times 3$  to fit the input image size of the pre-trained CNN models. No additional preprocessing is performed except the image size conversion because the CNN models include the convolution and pooling layers to extract features from the input image. The detailed structure of pre-trained CNN models will be described in Section III-D.

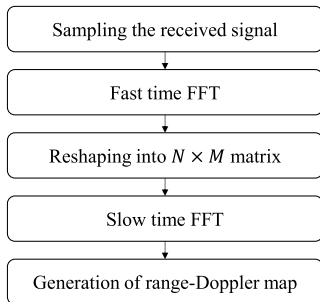
Figs. 4(a), 4(b), and 4(c) show the drone images taken by the built-in smartphone camera, and Fig. 4(d) presents an image of a military helicopter acquired from the open dataset. Whereas the Phantom4 drone is clearly recognized in Fig. 4(a), the Inspire2 drone is not well differentiated in Figs. 4(b) and 4(c) due to the background colors similar to that of the drone. These example images demonstrate the drawbacks of drone detection based on optical imaging.

#### B. RADAR SENSING

We obtain the radar sensing data by *EV-TINYRAD24G* [54]. This radar transmits the rapid chirps waveform, and the received echoes are used to construct the range-Doppler map representing the target range and velocity that can be



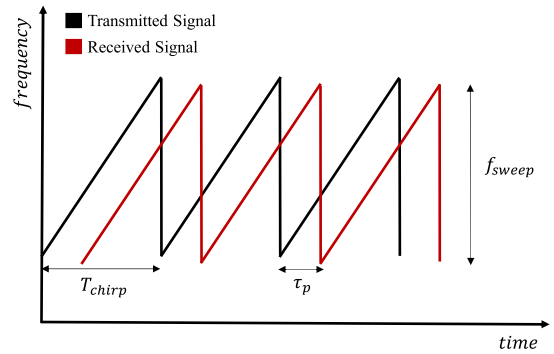
**FIGURE 4.** Images obtained by the built-in smartphone camera: (a) Phantom4 drone in flight, (b) Inspire2 drone in forest background with a similar color, (c) Inspire2 drone in building background with a similar color, (d) Military helicopter.



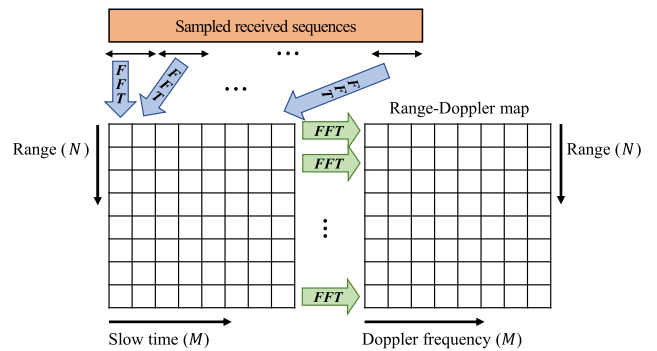
**FIGURE 5.** Overall procedure for generating the range-Doppler image from FMCW received signals.

used for drone detection and classification. Fig. 5 shows the overall procedure for generating the range-Doppler map using the received echoes of FMCW signals. Initially, the received signal is down-converted to a baseband signal and then sampled to form a matrix composed of complex samples. Subsequently, the fast time FFT is performed on the columns of the complex sample matrix, and followed by the execution of the slow time FFT on the rows in order to create a range-Doppler map corresponding to the sample matrix. The resulting range-Doppler map includes the range and velocity information of targets, which can be used for drone detection and classification.

Fig. 6 presents a conventional FMCW waveform with rapid chirps which have a very short duration  $T_{chirp}$ . By reducing this duration, the frequency components associated with the distance and velocity can be independently estimated, enabling low-complexity and high-accuracy radar signal



**FIGURE 6.** FMCW waveform with rapid chirps.



**FIGURE 7.** Generation of the range-Doppler map using two FFT operations.

processing. Specifically, the received signal is composed of reflected echoes from multiple targets as follows:

$$r(t) = \sum_{p=1}^P r_p(t), \tag{1}$$

where  $r_p(t)$  is the received FMCW echoes reflected from the  $p$ th target and  $P$  is the total number of targets. Suppose that there are no noises and clutters affecting the received signal. When the chirp is expressed as a frequency-modulated signal with instantaneous phase  $\varphi_i$ , the received signal can be expressed as [18]

$$r(t) = \sum_{p=1}^P A_p \sum_{m=0}^{M-1} \cos(\varphi_i(t - mT_{chirp} - \tau_p)) \exp(j2\pi v_p t). \tag{2}$$

where  $A_p$ ,  $\tau_p$ , and  $v_p$  denote the amplitude, time delay, and Doppler frequency shift corresponding to the  $p$ th target, respectively,  $M$  is the number of chirps, and the instantaneous phase is given by

$$\varphi_i(t) = 2\pi f_0 t + \pi k_f \alpha t^2. \tag{3}$$

Here,  $f_0$  is the lower carrier frequency,  $k_f$  is the frequency deviation, and  $\alpha$  is the modulation signal amplitude.

The frequency down-conversion is separately performed for each in-phase and quadrature component of the received



signal. After lowpass filtering, the baseband signal is accumulated to obtain the beat signal as follows:

$$b(t) = \sum_{m=0}^{M-1} \sum_{p=1}^P A_p \exp(j\varphi_{bmp}(t)), \quad (4)$$

where  $\varphi_{bmp}(t)$  is the instantaneous phase given by

$$\varphi_{bmp}(t) = \varphi_i(t - T_{chirp}) - \varphi_i(t - mT_{chirp} - \tau_p) + 2\pi v_p t. \quad (5)$$

By substituting (3) into (5), we have

$$\varphi_{bmp}(t) = \varphi_0 + 2\pi k_f \alpha \tau_p t + 2\pi v_p t, \quad (6)$$

where  $\varphi_0$  is a constant phase term independent of the time  $t$ . From (6), the instantaneous frequency of the beat signal is obtained as

$$f_{bmp} = \frac{1}{2\pi} \frac{d\varphi_{bmp}(t)}{dt} = k_f \alpha \tau_p + v_p. \quad (7)$$

Here, the first term of  $f_{bmp}$  is proportional to the delay  $\tau_p$  and the second term is equal to  $v_p$ , and used to estimate the target range and velocity, respectively. Thus, (7) can be rewritten as

$$f_{bmp} = \frac{f_{sweep}}{T_{chirp}} \tau_p + v_p = f_{Rp} + f_{Dp}, \quad (8)$$

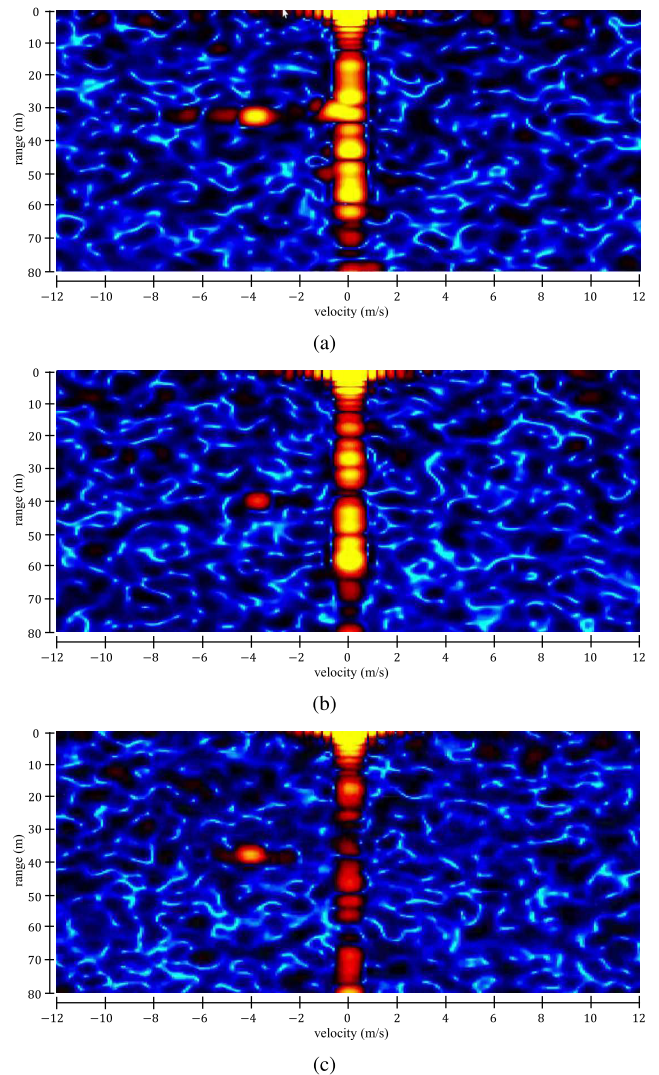
where  $f_{Rp}$  and  $f_{Dp}$  are the instantaneous frequencies related to the target range and velocity, respectively.

The received signal is sampled and the range-matched filtering is performed. When  $N$  is the number of samples during the fast time, the matched filter output is arranged as an  $N \times M$  complex matrix. The frequency  $f_{Rp}$  can be estimated by taking the fast time FFT, i.e., the  $N$ -point FFT is carried out on each column to the range direction, and the results are stored in the columns of the  $N \times M$  matrix as shown in Fig. 7. The FFT magnitudes are proportional to the amplitudes of targets (if a target exists at the considered frequency). After the FFT, the peaks of columns correspond to the target ranges, and the phases at the peaks of columns are denoted as

$$\varphi_{m,p} = \varphi_{0p} + 2\pi f_{Dp} m T_{chirp}, \quad (9)$$

where  $\varphi_{0p} = 2\pi f_0 \tau_p - \pi k_f \alpha \tau_p^2$  is a constant phase independent of the fast time and slow time indexes. Note that the phases after the fast time FFT depend on the chirp index  $m$  as shown in (9). After the slow time FFT performed to each row of the  $N \times M$  matrix, the resulting peak values are mapped to the Doppler frequencies  $f_{Dp}$ , as seen on the right side of Fig. 7. After the fast time and slow time FFTs, the final matrix represents the range-Doppler map with the range and velocity information of targets.

Figs. 8 and 9 present the range-Doppler maps obtained from the measured FMCW radar signals. The intense vertical yellow lines around the zero velocity are a kind of clutters caused by the leakage of transmit FMCW signals. In Figs. 8(a), 8(b), and 8(c), the red spots on the left plane indicate the Inspire2, Mavic3, and Phantom4 drones, respectively, moving away from the radar sensor at about 30 ~ 40 m distance with -4 m/s velocity (i.e., having

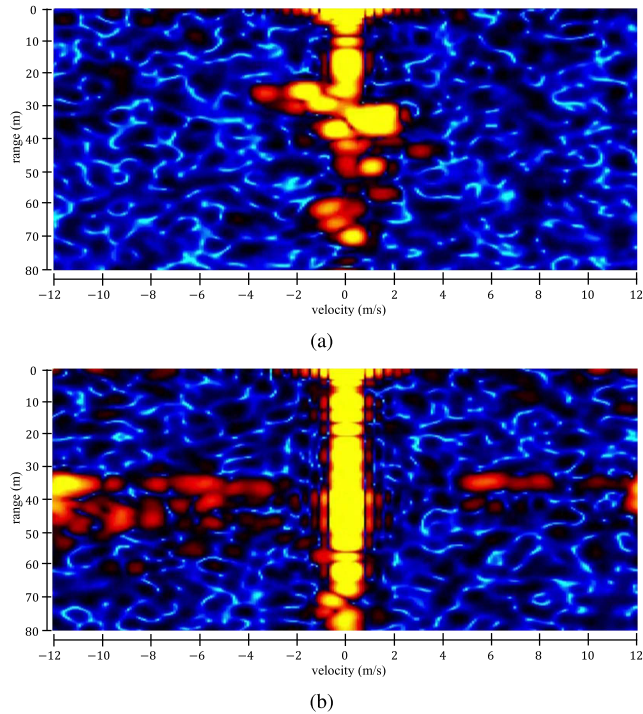


**FIGURE 8.** Range-Doppler maps obtained from the FMCW radar signals for the drones when the distance is 30 ~ 40 m, the velocity is -4 m/s, and the altitude is 8 m: (a) Inspire2, (b) Mavic3, (c) Phantom4.

a negative Doppler frequency). Moreover, 9(a) shows the range-Doppler map obtained from people playing soccer so that several spots are located near the yellow center line, and 9(b) denotes several vehicles driving on an eight-lane boulevard. These range-Doppler maps clearly demonstrate the difference between drones and non-drone objects, and also show subtle differences corresponding to three types of drones. For example, Inspire2 has the most vivid and widest yellow spot, whereas Mavic3 presents the weakest and smallest spot.

### C. ACOUSTIC SENSING

In the field test, acoustic signals are measured with the microphone shown in Fig. 3 using actual drone sounds in flight and non-drone sounds such as helicopters, vehicles, human voices, background noises, and so on. Additional non-drone sounds are obtained from the open dataset in [52] such as



**FIGURE 9.** Range-Doppler maps obtained from the FMCW radar signals for non-drone objects: (a) People playing soccer, (b) Vehicles driving on a boulevard.

those from engines, propellers, aircraft, rain and thunder, air conditioners, and background noises. The measured acoustic data is stored in the .wav file format with 48 kHz sampling rate and 16-bit quantization per sample.

As shown in Fig. 10, the audio file is converted to a mel spectrogram through audio signal processing. A spectrogram is a method for analyzing a sound waveform whose frequency characteristics change over time, which is derived by arranging the spectrum of the acoustic data in the time-frequency domain. Firstly, from the recorded audio waveform, we extract the one-second interval with the highest entropy. When the number of quantization bits is  $B$ , the entropy of the audio sequences starting at  $\ell$  is defined as

$$H(\ell) = - \sum_{m=1}^{2^B} p_{m,\ell} \log_2(p_{m,\ell}), \quad (10)$$

where  $p_{m,\ell}$  is the empirical probability corresponding to the quantization level  $q_m$ , i.e.,

$$p_{m,\ell} = E[x_n = q_m], \quad n = \ell, \ell + 1, \dots, \ell + L - 1. \quad (11)$$

Here,  $x_n$  is the  $n$ th input audio sample,  $q_m = (m - 0.5 - 2^{B-1})/2^B$ , and  $L$  is the number of samples corresponding to the one-second interval. From the definition in (10), the audio interval with the highest entropy can be selected as

$$\ell_o = \arg \max_{\ell \in \{1, 1+\Delta\ell, 1+2\Delta\ell, \dots\}} H(\ell), \quad (12)$$

where  $\Delta\ell$  is the index spacing. Note that the candidate starting index  $\ell$  is adjusted by  $\Delta\ell$  to reduce complexity.

As a next step, the Hamming window is applied to the extracted audio signals, and then 1440-point STFT is performed with 960 samples of analysis window overlap length. The number of bands for mel filtering is 32 and the number of frames is 89, when the sampling rate is 48 kHz. Finally, the size of the audio mel spectrogram is adjusted to fit the input size for the CNN model that conducts the drone detection or classification. Notice that the mel filter (or mel-scale triangle filter) describes low-frequency bands with high resolution while denoting high-frequency bands with low resolution. Therefore, the mel filter tends to emphasize the acoustic characteristics of drones in low-frequency bands.

Fig. 11 shows example spectrograms obtained by the audio signal processing in Fig. 10. Figs. 11(a), 11(b), and 11(c) present the spectrograms corresponding to Inspire2, Mavic3, and Phantom4, respectively. It is clearly seen that the spectrogram is different according to the type of drone. Fig. 11(d) is the spectrogram obtained from the sounds of people’s conversation, which is completely different from those of drones.

**D. CNN-BASED DETECTION AND CLASSIFICATION**

In this paper, we consider two-step approach composed of drone detection and classification. In the first step, we determine whether it is a drone or a non-drone object from a given image. If a drone is detected in the first step, the type of drone is identified in the second step using the same image. The input image can be either an optical image, a range-Doppler map, or an audio spectrogram. Drone detection and classification are conducted by utilizing six CNN models that individually adjust the neural network coefficients with the training data. Three CNN models are used for drone detection from the optical image, the range-Doppler map, and the audio spectrogram, respectively, and the other three models are utilized for drone classification based on the same input images obtained from three sensors. In a CNN model shown in Fig. 12, convolution layers and pooling layers that perform convolutional operations are repeatedly arranged to extract features of an input image, and the features are sent to the fully connected layer for detection and classification. In this paper, GoogLeNet [55], ResNet-101 [56], and DenseNet-201 [57] are employed among the pre-trained CNN models.

As mentioned before, measurement data obtained by the camera, radar, and microphone are used in combination with open datasets for drone detection and classification. Though, the number of images is not enough to train the CNN models, because the CNN models include a lot of parameters for feature extraction,<sup>1</sup> image processing, and metric computation for classification. To overcome this problem, we employ transfer learning that partially modifies a pre-trained CNN model for other purposes. As shown in Fig. 13, a pre-trained model is imported, and then some layers are newly configured

<sup>1</sup>GoogLeNet is composed of 22 layers with 6.8 million parameters, ResNet-101 consists of 101 layers with 1.7 million parameters, and DenseNet-201 has 201 layers with about 20 million parameters [55], [56], [57].



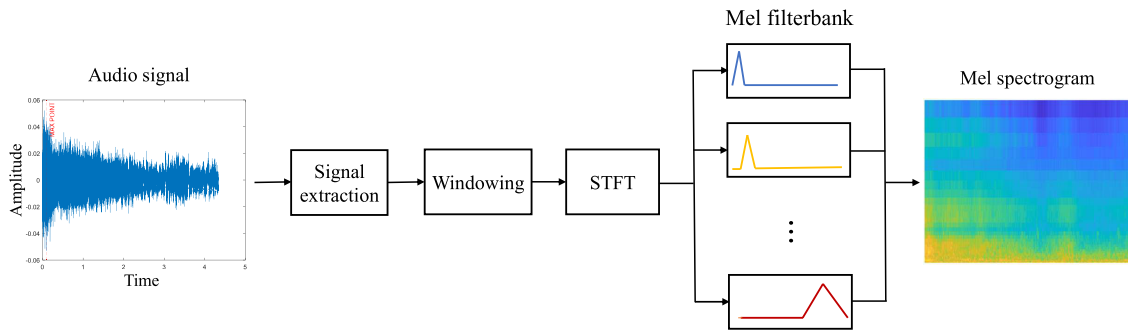


FIGURE 10. Overall procedure for generating the spectrogram from measured acoustic signals.

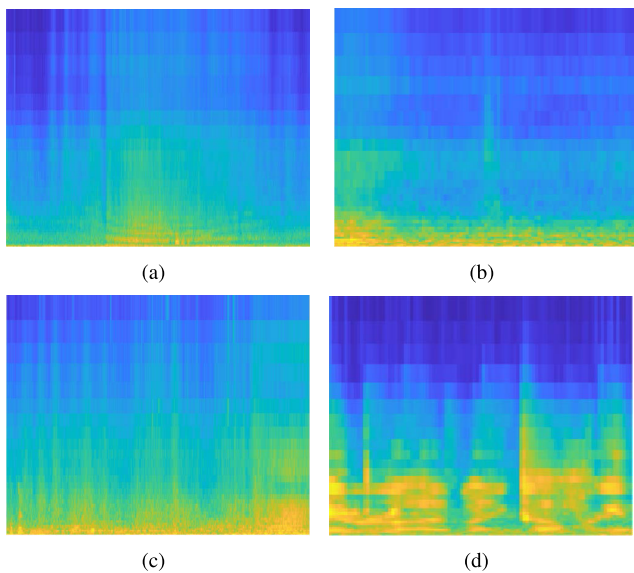


FIGURE 11. Spectrogram obtained from the audio signals: (a) Inspire2, (b) Mavic3, (c) Phantom4, (d) Sound of people conversation.

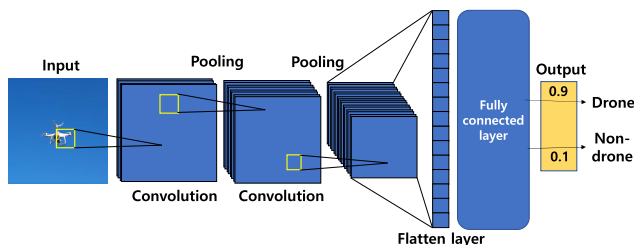


FIGURE 12. Overall processing architecture of a convolutional neural network.

and modified to produce an output suitable for new tasks. In this paper, the final layers of a pre-trained CNN model are replaced for drone surveillance, and the modified CNN model is trained using the corresponding training images. Through this procedure, six trained CNN models are developed for drone detection and classification with three kinds of sensing images (i.e., an optical image, range-Doppler map, and audio spectrogram).

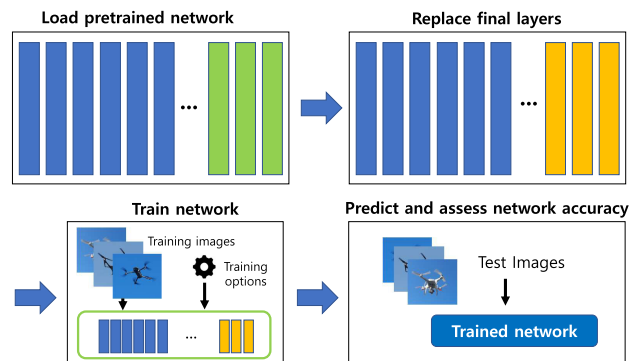


FIGURE 13. Modification of a CNN model for transfer learning.

TABLE 5. Parameters for transfer learning with optical images, radar images, and spectrogram.

Item	Optical Image	Radar Image & Spectrogram
Image augmentation	Reflection around x-axis Translation in x-axis and y-axis over [-30, 30]	No reflection, Translation in y-axis over [-30, 30]
Baseline CNN model	GoogLeNet, ResNet-101, DenseNet-201	
Input image size	224 × 224 × 3	
Minimum batch size	128	64
Number of epochs	40	40
Initial learning rate	0.0001	0.0001
Training/Validation ratio	70% / 30%	80% / 20%

As shown in Table 5, three kinds of sensing images are converted to 224 × 224 × 3 to fit the input image size of the pre-trained models. In the case of optical images, we use 70% of the data for training and 30% for verification of the trained CNN model, while we split the radar images and audio spectrograms into 80% and 20% for training and verification, respectively. Image augmentation is used to create more training examples from the measurement data and open datasets. Optical images are reflected around x-axis as well as translated in both x-axis and y-axis over [−30, 30], and range-Doppler maps and audio spectrograms are translated in y-axis over [−30, 30] with no reflection. In addition, the minimum batch size is set to 128 for optical images and 64 for range-Doppler maps and audio spectrograms considering the

number of training data, the number of epochs is set to 40, and the initial learning rate is set to 0.0001.

#### IV. PROPOSED SENSOR FUSION METHOD FOR UAV DETECTION AND CLASSIFICATION

In this section, to improve the surveillance performance, we propose a new drone detection and classification technique that combines multiple sensing schemes. When using the optical, radar, and acoustic sensing data, we can combine the optical image, range-Doppler map, and audio spectrogram to make a decision. For notational convenience, the optical image, radar sensing data, and audio signals are henceforth referred to as *image*, *radar*, and *audio* in the following of the paper. For example, the proposed sensor fusion method can combine two sensing data like *image + radar*, *image + audio*, and *radar + audio* as well as three sensing data like *image + radar + audio*. In the following of this section, we explain the proposed sensor fusion method combining three kinds of sensing data, i.e., *image + radar + audio*, because two-sensor fusion techniques are a special case of three-sensor fusion.

Fig. 14 presents the overall block diagram for drone detection and classification through sensor fusion of the image, radar, and audio data. As described in Section III-D, three kinds of sensing data are converted to the resized optical image, range-Doppler map, and audio spectrogram through pre-processing, respectively. An initial drone detection procedure is conducted by the CNN model with individual sensing data. By utilizing the logistic regression model, we combine the initial detection probabilities obtained from three CNN models corresponding to the image, radar, and audio data, and then determine whether a drone is present or not. If a drone is detected, we perform the drone classification procedure. The CNN models for drone classification separately compute the probabilities for Inspire2, Mavic3, and Phantom4 utilizing the same input data as the CNN models for drone detection. Finally, the multinomial logistic regression model is exploited to compute the combined probabilities for sensor fusion in the classification procedure.

The drone detection based on the sensor fusion is accomplished by the logistic regression model. Given training datasets, the logistic regression model is given by

$$\hat{y} = g(a_0 + a_1p_1 + a_2p_2 + a_3p_3), \quad (13)$$

where  $\hat{y}$  is an  $N \times 1$  vector predicting the probability for drone presence;  $p_1, p_2$ , and  $p_3$  represent the  $N \times 1$  probability vectors for training obtained from the CNN models with the image, radar, and audio sensing data, respectively;  $a_0$  is a bias term;  $a_1, a_2$ , and  $a_3$  are weight coefficients for the probabilities obtained from the image, radar, and audio CNN models; and  $N$  is the number of training datasets. Here,  $g(x)$  is the sigmoid function defined as

$$g(x) = \frac{1}{1 + e^{-x}}. \quad (14)$$

The logistic regression model in (13) can be rewritten in a vector-matrix form as follows:

$$\hat{y} = g(Pa), \quad (15)$$

where  $P = [1, p_1, p_2, p_3]$  and  $a = [a_0, a_1, a_2, a_3]^T$ . To find the optimal coefficients  $\{a_0, a_1, a_2, a_3\}$  for sensor fusion, we define the cost function

$$J(a) = -\frac{1}{N} \left[ y^T \log(\hat{y}) + (1 - y)^T \log(1 - \hat{y}) \right] + \frac{\lambda}{2N} a_0^T a_0, \quad (16)$$

where  $a_0 = [0, a_1, a_2, a_3]^T$  and  $\lambda$  is the regularization parameter. The gradient of  $J(a)$  is expressed as

$$\begin{aligned} \frac{\partial J(a)}{\partial a} &= \frac{1}{N} D^{-1}(\hat{y}) \frac{\partial \hat{y}^T}{\partial a} y \\ &\quad + \frac{1}{N} D^{-1}(1 - \hat{y}) \frac{\partial \hat{y}^T}{\partial a} (y - 1) + \frac{\lambda}{N} a_0, \end{aligned} \quad (17)$$

where  $D(x) = \text{diag}([x_1, x_2, \dots, x_N])$  for  $x = [x_1, x_2, \dots, x_N]^T$ . Here, from the logistic regression model in (15), we have

$$\frac{\partial \hat{y}^T}{\partial a} = P^T D^2(\hat{y}) D(\exp(-Pa)). \quad (18)$$

By substituting (18) into (17), the gradient is expressed as

$$\frac{\partial J(a)}{\partial a} = \frac{1}{N} \left[ P^T (\hat{y} - y) + \lambda a_0 \right]. \quad (19)$$

Using the gradient method, the coefficient vector  $a(j)$  at the  $j$ th iteration can be updated as

$$\begin{aligned} a(j) &= a(j - 1) - \mu \frac{\partial J(a)}{\partial a} \\ &= a(j - 1) - \mu \frac{1}{N} \left[ P^T (\hat{y} - y) + \lambda a_0 \right]. \end{aligned} \quad (20)$$

Suppose that  $a^o = [a_0^o, a_1^o, a_2^o, a_3^o]^T$  is an optimal coefficients for sensor fusion obtained by (20), the test datasets are used to evaluate the drone detection performance as follows:

$$\hat{y}^{test} = g(a_0^o + a_1^o p_1^{test} + a_2^o p_2^{test} + a_3^o p_3^{test}), \quad (21)$$

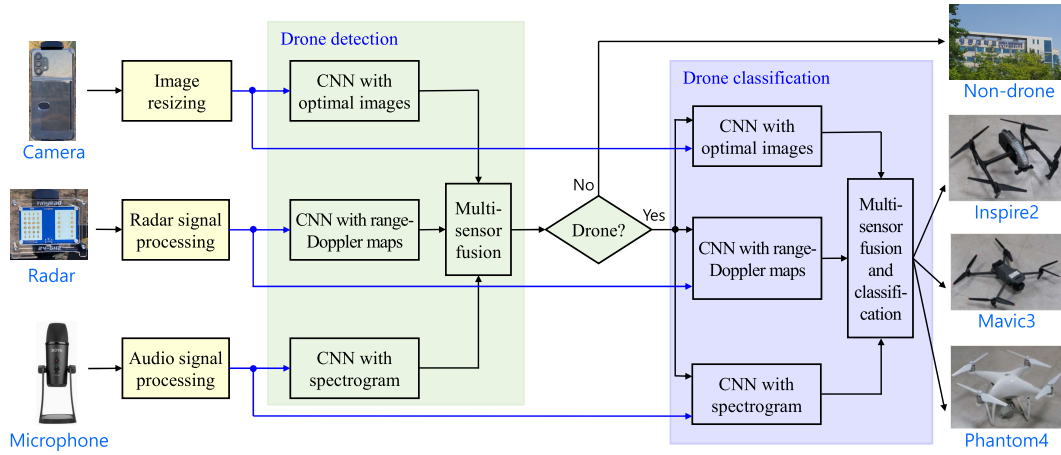
where  $\hat{y}^{test}$  is an  $M \times 1$  vector representing the probability of drone presence;  $p_1^{test}, p_2^{test}$ , and  $p_3^{test}$  are the  $M \times 1$  probability vectors of test datasets obtained from the CNN models with the image, radar, and audio sensing data, respectively; and  $M$  is the number of test datasets.

For the drone classification based on sensor fusion, we exploit the multinomial logistic regression with the logit model. Given the training datasets, the model for the relative risk is denoted as [53]

$$\begin{aligned} \log(r_{13}) &= b_0 + b_{1,M} q_{1,M} + b_{1,P} q_{1,P} + b_{2,M} q_{2,M} \\ &\quad + b_{2,P} q_{2,P} + b_{3,M} q_{3,M} + b_{3,P} q_{3,P} \end{aligned} \quad (22a)$$

$$\begin{aligned} \log(r_{23}) &= c_0 + c_{1,M} q_{1,M} + c_{1,P} q_{1,P} + c_{2,M} q_{2,M} \\ &\quad + c_{2,P} q_{2,P} + c_{3,M} q_{3,M} + c_{3,P} q_{3,P}, \end{aligned} \quad (22b)$$

where  $q_{1,*}, q_{2,*}$ , and  $q_{3,*}$  represent the  $N_c \times 1$  probability vectors obtained from the CNN models for drone classification



**FIGURE 14.** Block diagram for drone detection followed by classification through sensor fusion of the image, radar, and audio.

with the image, radar, and audio sensing data, respectively;  $\mathbf{q}_{j,M}$  and  $\mathbf{q}_{j,P}$  mean the probability vectors associated with Mavic3 and Phantom4;  $b_0$  and  $c_0$  are bias terms;  $\{b_{j,*}\}$  and  $\{c_{j,*}\}$  are weight coefficients for logistic regression; and  $N_c$  is the number of training datasets for drone classification. Here,  $\mathbf{r}_{kl}$  is given by

$$\mathbf{r}_{kl} = \left[ \frac{P(y_1 = k)}{P(y_1 = \ell)}, \frac{P(y_2 = k)}{P(y_2 = \ell)}, \dots, \frac{P(y_{N_c} = k)}{P(y_{N_c} = \ell)} \right]^T. \quad (23)$$

where  $k, \ell \in \{1, 2, 3\}$ , and  $P(y_i = 1)$ ,  $P(y_i = 2)$ , and  $P(y_i = 3)$  denote the probabilities that the  $i$ th observation is Inspire2, Mavic3, and Phantom4, respectively. The equations in (22a) can be rewritten in a vector-matrix form as below:

$$\log \frac{P(\mathbf{y} = 1)}{P(\mathbf{y} = 3)} = \mathbf{Q}\mathbf{b}, \quad (24a)$$

$$\log \frac{P(\mathbf{y} = 2)}{P(\mathbf{y} = 3)} = \mathbf{Q}\mathbf{c}, \quad (24b)$$

where  $\mathbf{Q} = [\mathbf{1} \ \mathbf{q}_{1,M} \ \mathbf{q}_{1,P} \ \mathbf{q}_{2,M} \ \mathbf{q}_{2,P} \ \mathbf{q}_{3,M} \ \mathbf{q}_{3,P}]$ ,  $\mathbf{b} = [b_0, b_{1,M}, b_{1,P}, b_{2,M}, b_{2,P}, b_{3,M}, b_{3,P}]^T$ , and  $\mathbf{c} = [c_0, c_{1,M}, c_{1,P}, c_{2,M}, c_{2,P}, c_{3,M}, c_{3,P}]^T$ . Following the approach in [53], the problem for finding the optimal coefficients is formulated as the maximum a posteriori (MAP) estimation, and can be solved by an iterative procedure such as the gradient-based optimization algorithm [53] and the coordinate descent algorithm [58].

Using the optimal coefficient vectors  $\mathbf{b}^o$  and  $\mathbf{c}^o$ , we can predict the probabilities for drone classification given test datasets. From (24a), we may write

$$P(\mathbf{y}^{test} = 1) = P(\mathbf{y}^{test} = 3) \exp(\mathbf{Q}^{test} \mathbf{b}^o) \quad (25a)$$

$$P(\mathbf{y}^{test} = 2) = P(\mathbf{y}^{test} = 3) \exp(\mathbf{Q}^{test} \mathbf{c}^o), \quad (25b)$$

where  $\mathbf{Q}^{test}$  is an  $M_c \times 7$  matrix composed of the probabilities obtained from the CNN models using the test datasets. Using (25a) and the fact that  $P(\mathbf{y} = 1) + P(\mathbf{y} = 2) + P(\mathbf{y} = 3) = 1$ , we can

predict the probabilities for drone classification as below:

$$P(\mathbf{y}^{test} = 1) = \frac{\exp(\mathbf{Q}^{test} \mathbf{b}^o)}{1 + \exp(\mathbf{Q}^{test} \mathbf{b}^o) + \exp(\mathbf{Q}^{test} \mathbf{c}^o)} \quad (26a)$$

$$P(\mathbf{y}^{test} = 2) = \frac{\exp(\mathbf{Q}^{test} \mathbf{c}^o)}{1 + \exp(\mathbf{Q}^{test} \mathbf{b}^o) + \exp(\mathbf{Q}^{test} \mathbf{c}^o)} \quad (26b)$$

$$P(\mathbf{y}^{test} = 3) = \frac{1}{1 + \exp(\mathbf{Q}^{test} \mathbf{b}^o) + \exp(\mathbf{Q}^{test} \mathbf{c}^o)}. \quad (26c)$$

## V. NUMERICAL RESULTS

In this section, we evaluate the drone surveillance performance of the proposed sensor fusion method and compare those of the schemes based on individual sensors. Specifically, we consider the following methods for drone detection and classification.

- *Image* [9]: Based on the optical images, CNN models are used for drone detection and classification as in [9]. For training and verification, the measured optical images in Section III-A are used along with the open datasets available in [51].
- *Radar* [20]: Based on the range-Doppler maps obtained from the FMCW radar, CNN models are used for drone detection and classification as in [20]. For training and verification, the measured radar signals are converted to the range-Doppler maps as explained in Section III-B.
- *Audio* [34]: Based on the audio spectrograms, CNN models are used for drone detection and classification as in [34]. For training and verification, the measured audio signals are converted to the spectrograms as in Section III-C and the open datasets in [52] are used as well.
- *Image + Radar*: The proposed sensor fusion method is designed by combining the optical images and range-Doppler maps.
- *Image + Audio*: The proposed sensor fusion method is utilized by combining the optical images and audio spectrograms.



**TABLE 6. Overall datasets for training and verification of the CNNs models corresponding to the optical images, radar range-Doppler maps, and audio spectrograms.**

Stage	Sensor	Object	Online Data	Measured Data
Detection	Image	Drone	462	8470
		Non-drone	4353	400
	Radar	Drone	-	13620
		Non-drone	-	10728
	Audio	Drone	1225	1233
		Non-drone	4478	300
Classification	Image	Inspire2	-	3576
		Mavic3	-	2665
		Phantom4	-	2229
	Radar	Inspire2	-	7085
		Mavic3	-	2777
		Phantom4	-	3758
	Audio	Inspire2	-	426
		Mavic3	-	433
		Phantom4	-	374

**TABLE 7. Open datasets obtained from [51], [52].**

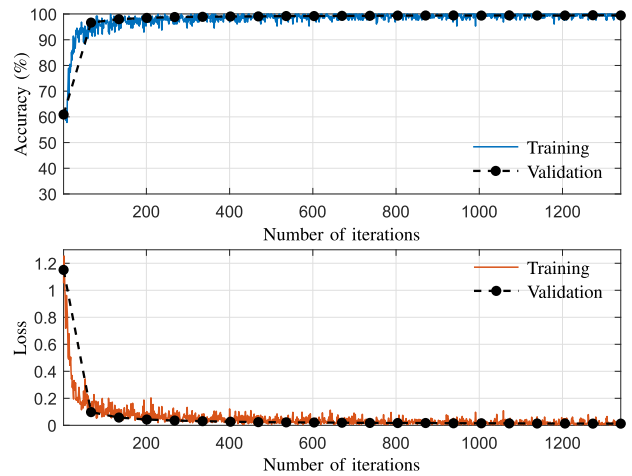
Type	Drone	Non-Drone Object	
Image	462	Airplane	1611
		Helicopter	1730
		Warplane	751
		Rocket	261
Audio	1225	Vehicle engine	837
		Aircraft propeller	402
		Rain & Thunder	320
		Air conditioner	873
		Background noise	2346

**TABLE 8. Number of datasets for coefficient training and final test in the proposed sensor fusion method based on the multinomial logistic regression.**

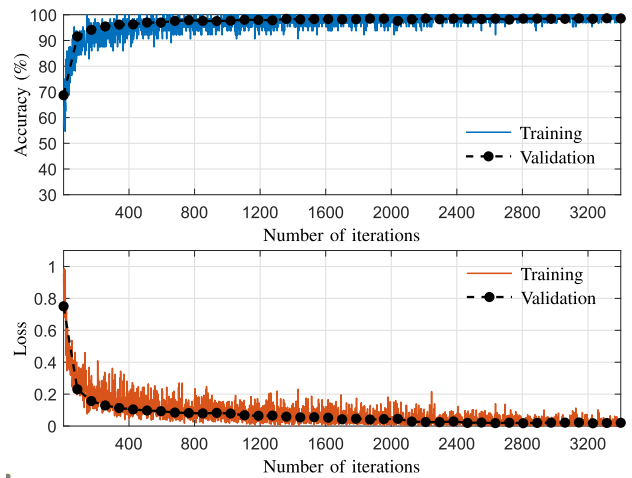
Type	Class	Coeff. Training	Final Test
Detection	drone	1200	900
	non-drone	450	300
Classification	Inspire2	450	300
	Mavic3	450	300
	Phantom4	450	300

- *Radar + Audio*: The proposed sensor fusion method is used by combining the range-Doppler maps and audio spectrograms.
- *Image + Radar + Audio*: The proposed sensor fusion method in Section IV is fully implemented by integrating the optical images, range-Doppler maps, and audio spectrograms.

We employed pre-trained CNN models provided by MATLAB deep learning toolbox, and modified the CNN models via transfer learning as described in Section III-D. Table 6 presents the overall datasets for training the CNN models with the optical images, range-Doppler maps, and audio spectrograms. In the case of the optical sensing, a total of 13685 datasets were used including 8870 field measurement images and 4815 online datasets in [51]. Through actual measurements, we obtained 400 non-drone images and 8470 drone images composed of 3576, 2665, and 2229 datasets for Inspire2, Mavic3, and phantom4,



**FIGURE 15. Learning curves of the GoogLeNet used for drone detection with optical images.**



**FIGURE 16. Learning curves of the GoogLeNet used for drone detection with radar range-Doppler maps.**

respectively. In the case of the radar sensing, a total of 24348 datasets were obtained through field measurements, i.e., 10728 range-Doppler maps for non-drone objects, 7085 datasets for Inspire2, 2777 datasets for Mavic3, and 3758 datasets for Phantom4. Note that open datasets were not employed for radar sensing because it is difficult to find range-Doppler images that match the FMCW radar specifications used in our experiments. In the case of the acoustic sensing, a total of 7236 spectrograms were used including 1533 actual measurement datasets and 5703 online datasets in [52]. In field measurements, we acquired 300 audio datasets for non-drone objects, 426 datasets for Inspire2, 433 datasets for Mavic3, and 374 datasets for Phantom4.

As shown in Table 7, the open datasets for optical images consist of images for airplanes, helicopters, warplanes, and rockets similar to drone images in flight. Also, the open datasets for acoustic signals include the sounds of vehicle engines, aircraft propellers, rain and thunder, air conditioners, and various background noises. It is noticeable that all datasets for drones and non-drone objects are utilized when

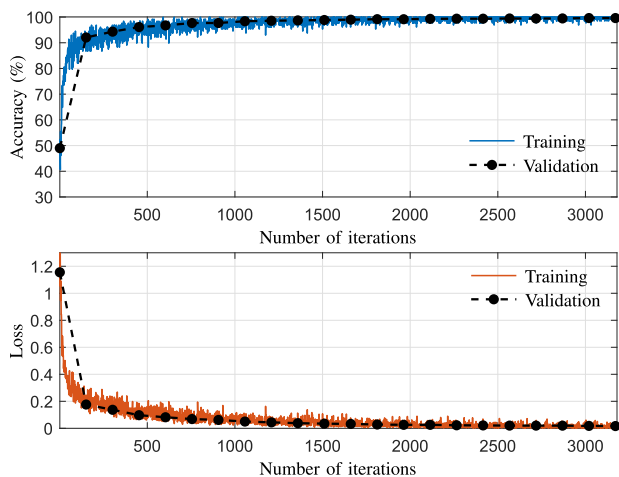


FIGURE 17. Learning curves of the GoogLeNet used for drone detection with audio spectrograms.

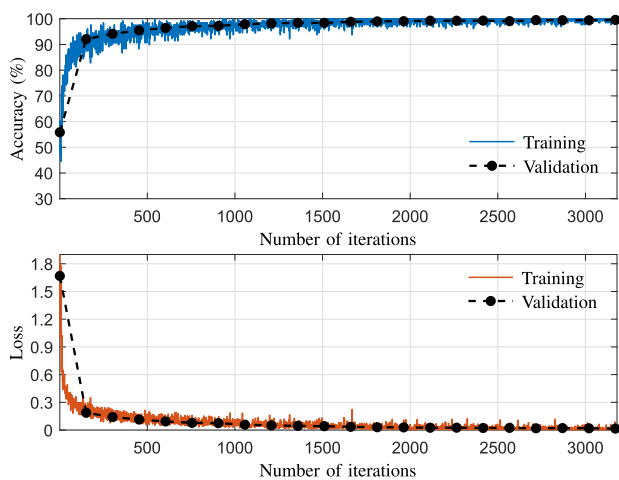


FIGURE 18. Learning curves of the GoogLeNet used for drone classification with optical images.

the CNN models are applied to drone detection while only the datasets for drones are used to the CNN models for drone classification.

The proposed sensor fusion method requires datasets concurrently measured from the camera, radar, and microphone to combine the optical image, the range-Doppler map, and the audio spectrogram obtained under the same drone flight conditions. Table 8 describes the number of datasets for coefficient training and the final test in the multinomial logistic regression model obtained by the field measurements. For drone detection, we used 900 drone datasets and 300 non-drone datasets in the training mode to find the optimal coefficients, and performed the final test for 450 drone datasets and 150 non-drone datasets. For drone classification, we used the same drone datasets as those for drone detection. So, we exploited 300 datasets for each type of drone in the training and 150 datasets for each type of drone in the final test.

Figs. 15–17 show the learning curves of the CNN models, which are applied to drone detection using optical images,

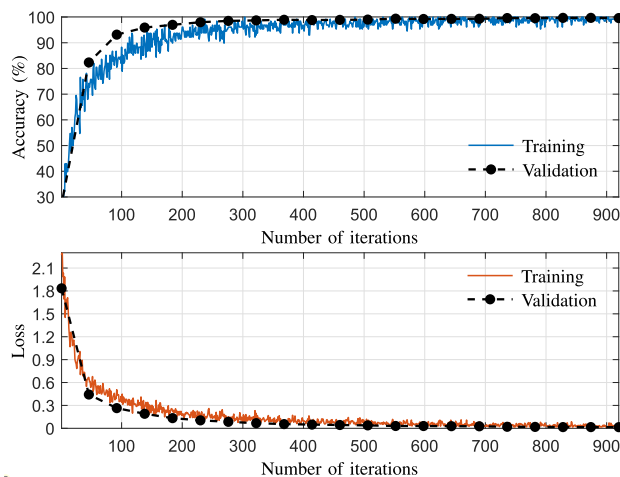


FIGURE 19. Learning curves of the GoogLeNet used for drone classification with radar range-Doppler maps.

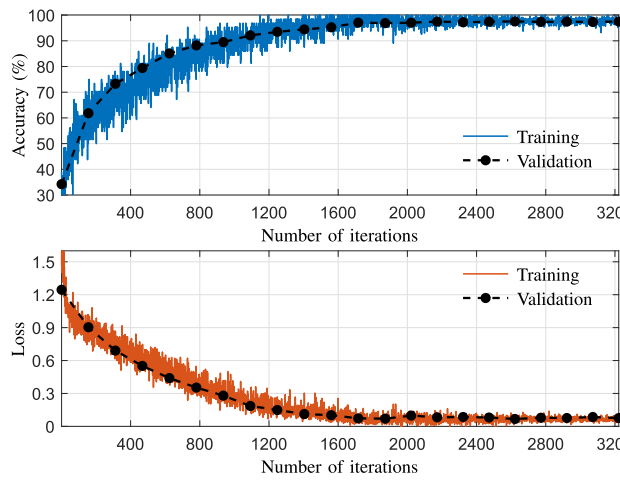


FIGURE 20. Learning curves of the GoogLeNet used for drone classification with audio spectrograms.

radar range-Doppler maps, and audio spectrograms, respectively, and Figs. 18–20 present the learning curves of the CNN models used for drone classification. In the simulations, GoogLeNet was used as a pre-trained CNN model and the parameters were set as in Table 5. Overall, with the increment of the number of iterations, the accuracy gradually increases while the loss function gradually decreases. The converging speed is somewhat different depending on the type of sensors and the sort of surveillance (detection or classification), yet the accuracy and the loss function converge to the steady-state values when the number of iterations is greater than 2000 in all cases.

### A. UAV DETECTION RESULTS

Figs. 21–23 show the confusion matrices for drone detection using GoogLeNet, ResNet-101, and DenseNet-201, respectively, with the input images obtained from the individual sensors and the two-sensor fusion techniques. Fig. 24 presents the results for drone detection using the proposed three-sensor fusion method. Moreover, Table 9 denotes the

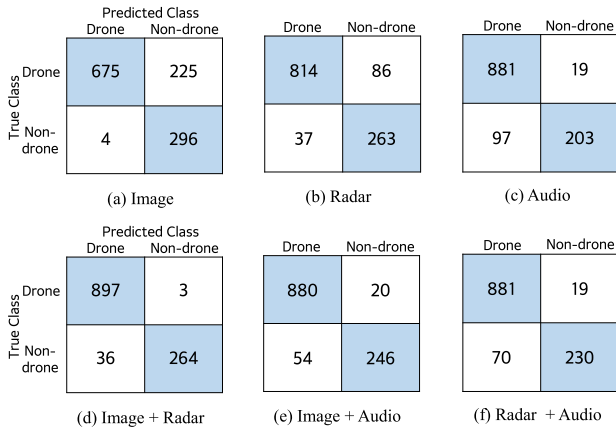


FIGURE 21. Drone detection results using the GoogLeNet with individual sensors and two-sensor fusion.

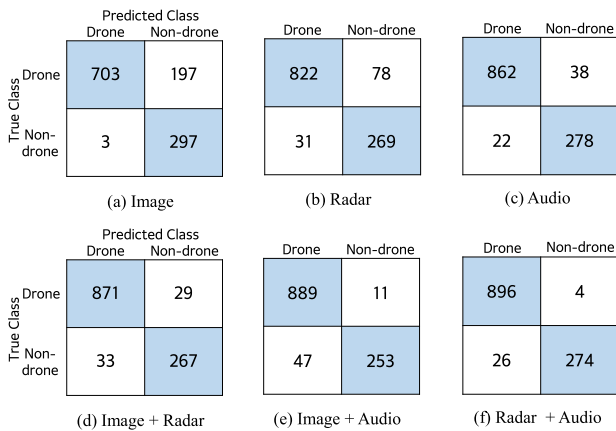


FIGURE 22. Drone detection results using the ResNet-101 with individual sensors and two-sensor fusion.

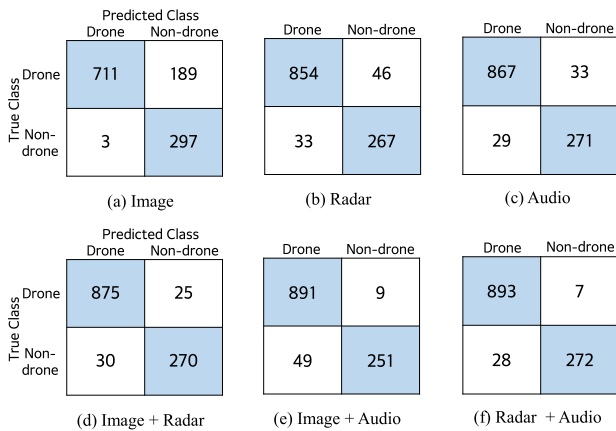


FIGURE 23. Drone detection results using the DenseNet-201 with individual sensors and two-sensor fusion.

detection accuracy of individual and combined sensing methods for drones and non-drone objects, where the positive predictive value (PPV) and the true positive rate (TPR) are also called the precision and the recall, respectively. The F-score,  $F_1$ , is defined as

$$F_1 = 2 \frac{PPV \times TPR}{PPV + TPR}. \quad (27)$$

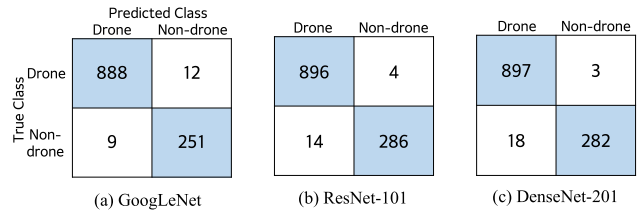


FIGURE 24. Drone detection results using various CNN models with three-sensor fusion (Image+Radar+Audio).

TABLE 9. Detection accuracy of individual and combined sensing methods for drones and non-drone objects, where the bold numbers indicate the highest value in each column. (PPV = positive predictive value, TPR = true positive rate.)

Method	Detection Accuracy (%)								
	GoogLeNet			ResNet-101			DenseNet-201		
	PPV	TPR	F-score	PPV	TPR	F-score	PPV	TPR	F-score
Image	<b>99.4</b>	75.0	85.5	<b>99.6</b>	78.1	87.5	<b>99.6</b>	79.0	88.1
Radar	95.7	90.4	93.0	96.4	91.3	93.8	96.3	94.9	95.6
Audio	90.1	97.9	93.8	97.5	95.8	96.6	96.8	96.3	96.5
Image + Radar	96.1	<b>99.7</b>	97.9	96.3	96.8	96.6	96.7	97.2	97.0
Image + Audio	94.2	97.8	96.0	95.0	98.8	96.8	94.8	99.0	96.8
Radar + Audio	92.6	97.9	95.2	97.2	99.6	98.4	97.0	99.2	98.1
Image + Radar + Audio	99.0	98.7	<b>98.8</b>	98.5	<b>99.6</b>	<b>99.0</b>	98.0	<b>99.7</b>	<b>98.8</b>

Among individual sensing methods, the optical image has the lowest detection accuracy and the audio sensor achieves the highest detection accuracy in terms of the F-score, because the verification datasets contain drone images with drone-like background colors and non-drone images that are difficult to distinguish from drones (see Fig. 4). The F-score tends to increase as the number of combined sensors increases, and thus the proposed three-sensor fusion method presents the highest F-score for all CNN models. Specifically, the F-score is improved by 2.4% ~ 15.6% by the proposed three-sensor fusion method compared to the individual sensing schemes. In the proposed three-sensor fusion method, the ResNet-101 model achieves slightly better F-score than the GoogLeNet and DenseNet-201.

### B. UAV CLASSIFICATION RESULTS

Figs. 25–27 denote the confusion matrices for drone classification using GoogLeNet, ResNet-101, and DenseNet-201, respectively, with the input images obtained from the individual sensors and the two-sensor fusion schemes. Fig. 28 shows the results for drone classification using the proposed three-sensor fusion method. Also, Table 10 presents the drone classification accuracy for individual and combined sensing techniques. Here, the average PPV, the average TPR, and the average of class-wise F-scores are defined as

$$P_{avg} = \frac{1}{3} \sum_{k=1}^3 P_k, \quad R_{avg} = \frac{1}{3} \sum_{k=1}^3 R_k \quad (28a)$$



	Inspire2	Mavic3	Phantom4
Inspire2	180	82	38
Mavic3	31	207	62
Phantom4	17	37	246

	Inspire2	Mavic3	Phantom4
Inspire2	211	56	33
Mavic3	89	194	17
Phantom4	31	75	194

	Inspire2	Mavic3	Phantom4
Inspire2	186	44	70
Mavic3	54	169	77
Phantom4	25	92	183

	Inspire2	Mavic3	Phantom4
Inspire2	240	36	24
Mavic3	54	211	35
Phantom4	8	40	252

	Inspire2	Mavic3	Phantom4
Inspire2	219	53	28
Mavic3	32	215	53
Phantom4	8	40	252

	Inspire2	Mavic3	Phantom4
Inspire2	220	53	27
Mavic3	60	206	34
Phantom4	18	74	208

FIGURE 25. Drone classification results using the GoogLeNet with individual sensors and two-sensor fusion.

	Inspire2	Mavic3	Phantom4
Inspire2	203	82	15
Mavic3	81	219	0
Phantom4	12	53	235

	Inspire2	Mavic3	Phantom4
Inspire2	210	83	7
Mavic3	58	231	11
Phantom4	79	27	194

	Inspire2	Mavic3	Phantom4
Inspire2	205	75	20
Mavic3	71	214	15
Phantom4	61	52	187

	Inspire2	Mavic3	Phantom4
Inspire2	228	64	8
Mavic3	70	224	6
Phantom4	24	19	257

	Inspire2	Mavic3	Phantom4
Inspire2	214	74	12
Mavic3	58	232	10
Phantom4	16	44	240

	Inspire2	Mavic3	Phantom4
Inspire2	201	79	20
Mavic3	50	238	12
Phantom4	38	38	224

FIGURE 26. Drone classification results using the ResNet-101 with individual sensors and two-sensor fusion.

	Inspire2	Mavic3	Phantom4
Inspire2	222	57	21
Mavic3	91	209	0
Phantom4	53	19	228

	Inspire2	Mavic3	Phantom4
Inspire2	236	64	0
Mavic3	55	208	37
Phantom4	65	34	201

	Inspire2	Mavic3	Phantom4
Inspire2	202	62	36
Mavic3	57	213	30
Phantom4	19	76	205

	Inspire2	Mavic3	Phantom4
Inspire2	237	31	32
Mavic3	54	230	16
Phantom4	58	5	237

	Inspire2	Mavic3	Phantom4
Inspire2	234	28	38
Mavic3	60	218	22
Phantom4	33	25	242

	Inspire2	Mavic3	Phantom4
Inspire2	226	19	55
Mavic3	55	213	32
Phantom4	22	40	238

FIGURE 27. Drone classification results using the DenseNet-201 with individual sensors and two-sensor fusion.

$$F_{1,avg} = \frac{1}{3} \sum_{k=1}^3 F_{1,k}, \quad (28b)$$

respectively, where  $P_k$ ,  $R_k$ , and  $F_{1,k}$  are given by

$$P_k = \frac{c_{k,k}}{c_{1,k} + c_{2,k} + c_{3,k}} \quad (29a)$$

$$R_k = \frac{c_{k,k}}{c_{k,1} + c_{k,2} + c_{k,3}} \quad (29b)$$

	Inspire2	Mavic3	Phantom4
Inspire2	219	50	31
Mavic3	40	229	31
Phantom4	23	34	243

	Inspire2	Mavic3	Phantom4
Inspire2	235	52	13
Mavic3	53	236	11
Phantom4	30	23	247

	Inspire2	Mavic3	Phantom4
Inspire2	236	28	36
Mavic3	50	232	18
Phantom4	30	15	255

FIGURE 28. Drone classification results using various CNN models with three-sensor fusion (Image+Radar+Audio).

TABLE 10. Classification accuracy of individual and combined sensing methods for entire datasets, where the bold numbers indicate the highest value in each column. (PPV = average positive predictive value, TPR = average true positive rate, F-score = average of class-wise F-scores.)

Method	Classification Accuracy (%)								
	GoogLeNet			ResNet-101			DenseNet-201		
	PPV	TPR	F-score	PPV	TPR	F-score	PPV	TPR	F-score
Image	71.2	70.3	70.2	74.8	73.0	73.5	75.2	73.2	73.7
Radar	67.6	66.6	66.8	73.3	70.6	70.9	72.9	71.7	71.8
Audio	60.4	59.8	59.9	69.3	67.3	67.6	69.7	68.9	69.0
Image + Radar	<b>78.0</b>	<b>78.1</b>	<b>78.0</b>	79.5	78.8	79.0	79.2	78.2	78.4
Image + Audio	76.7	76.2	76.2	77.4	76.2	76.5	77.4	77.1	77.1
Radar + Audio	71.0	70.4	70.6	74.7	73.7	73.8	75.4	75.2	75.2
Image + Radar + Audio	76.8	76.8	76.8	<b>80.3</b>	<b>79.8</b>	<b>79.9</b>	<b>80.5</b>	<b>80.3</b>	<b>80.4</b>

$$F_{1,k} = 2 \frac{PPV_k \times TPR_k}{PPV_k + TPR_k}. \quad (29c)$$

Here,  $c_{m,n}$  is the  $(m, n)$ th element of the confusion matrix for drone classification.

Overall, the classification accuracy in Table 10 is lower than the detection accuracy in Table 9 because distinguishing the drone type is more challenging than determining the presence of a drone. As shown in Figs. 8 and 9, the difference in range-Doppler maps between drones and non-drone objects is much more prominent than the difference among the three types of drones. Similar results are observed in the audio spectrogram in Fig. 11, and it is inferred that the optical sensor has similar trends to the radar and audio sensors. Therefore, the drone classification methods exhibit much lower accuracy than the corresponding drone detection schemes. Considering the individual sensors, the image sensor presents the highest F-score while the audio sensor obtains the lowest F-score. The proposed two-sensor fusion methods such as Image+Radar, Image+Audio, and Radar+Audio achieve higher F-scores than individual sensing schemes. The proposed three-sensor fusion method obtains the highest F-score among all drone classification techniques in the ResNet-101 and DenseNet-201 models, while the Image+Radar fusion scheme presents slightly better performance than the three-sensor fusion method in the GoogLeNet. Specifically, the F-score is improved by 8.7% ~ 28.1% by the proposed three-sensor fusion method compared to the individual sensing schemes. In the pro-

**TABLE 11. Classification accuracy of individual and combined sensing methods for Dataset-A, where the bold numbers indicate the highest value in each column.**

Method	Classification Accuracy (%)								
	GoogLeNet			ResNet-101			DenseNet-201		
	PPV	TPR	F-score	PPV	TPR	F-score	PPV	TPR	F-score
Image	90.8	87.8	87.5	85.8	85.8	85.7	88.0	87.3	87.5
Radar	75.6	72.9	73.2	81.8	81.1	81.2	84.0	81.8	82.1
Audio	69.3	67.8	67.9	81.9	78.9	79.4	80.5	80.0	80.1
Image + Radar	90.6	90.7	90.6	90.5	90.0	90.1	91.5	90.7	90.8
Image + Audio	92.8	91.6	91.3	92.0	92.0	92.0	91.9	91.6	91.6
Radar + Audio	80.6	80.2	80.2	87.7	86.9	87.0	88.1	88.0	88.0
Image + Radar + Audio	<b>95.1</b>	<b>95.1</b>	<b>95.1</b>	<b>95.3</b>	<b>95.3</b>	<b>95.3</b>	<b>95.4</b>	<b>95.3</b>	<b>95.4</b>

**TABLE 12. Classification accuracy of individual and combined sensing methods for Dataset-B, where the bold numbers indicate the highest value in each column.**

Method	Classification Accuracy (%)								
	GoogLeNet			ResNet-101			DenseNet-201		
	PPV	TPR	F-score	PPV	TPR	F-score	PPV	TPR	F-score
Image	53.3	52.9	52.4	66.6	60.2	61.8	64.1	59.1	59.8
Radar	60.5	60.2	60.3	67.4	60.0	61.2	63.0	61.6	61.6
Audio	55.9	51.8	52.4	58.0	55.8	56.1	59.2	57.8	58.1
Image + Radar	<b>65.2</b>	<b>65.6</b>	<b>65.3</b>	<b>69.7</b>	<b>67.6</b>	<b>68.3</b>	<b>67.2</b>	<b>65.8</b>	<b>66.1</b>
Image + Audio	60.7	60.9	60.6	65.5	60.4	61.6	63.8	62.7	62.5
Radar + Audio	62.0	60.7	61.0	61.5	60.4	60.5	62.9	62.4	62.2
Image + Radar + Audio	58.5	58.4	58.4	66.2	64.2	64.7	65.9	65.3	65.3

posed three-sensor fusion method, the DenseNet-201 model achieves the best F-score.

To further investigate the results of drone classification, we separate the test datasets into two groups referred to as *Dataset-A* and *Dataset-B*. Dataset-A consists of test datasets with high classification accuracy. Specifically, the distance between the sensor and the drone is less than half the maximum distance for optical, radar, and audio sensing. In optical sensing, the drone altitude is lower than 10 m, and the background is relatively simple like a blue sky. Also, the audio signal is recorded in situations with low background noises. In contrast, Dataset-B is composed of test datasets with low classification accuracy. The distance between the sensor and the drone is greater than half the maximum distance for optical, radar, and audio sensing. In optical sensing, the drone altitude is higher than 10 m, and the background is relatively complicated like many trees and buildings. Moreover, relatively high background noises are included in the audio signals. Tables 11 and 12 denote the classification accuracies for Dataset-A and Dataset-B, respectively. As expected, in Dataset-A, the individual sensing meth-

ods have the worst performance, and the performance is improved as the number of combined sensors increases. Thus, the proposed three-sensor fusion method achieves the best F-score irrespective of pre-trained CNN models. However, in Dataset-B, the proposed three-sensor fusion scheme does not guarantee the best performance. For instance, the Image+Radar method achieves higher F-scores than the Image+Radar+Audio scheme in all CNN models. These results imply that the initial classification performance of each sensor needs to exceed a certain threshold in order to enhance the classification accuracy through sensor fusion schemes.

## VI. CONCLUSION

In this paper, we proposed a sensor fusion method for drone detection and classification based on the CNN models for individual sensing and the multinomial logistic regression for combining the optical, radar, and audio sensing data. Through field experiments and numerical simulations, it was verified that the proposed sensor fusion scheme improves drone surveillance performance compared to individual sensing methods. It was also shown that the sensor fusion approach does not guarantee performance enhancement when the accuracy of individual sensing is low. Integrating multiple sensors is crucial for drone surveillance because individual sensing schemes, such as the optical camera, radar, and audio microphone, have complementary advantages and disadvantages. The results presented in this paper can be exploited to optimize the combining algorithm for sensor fusion when designing anti-UAV defense systems to protect security areas.

## REFERENCES

- [1] G. Ding, Q. Wu, L. Zhang, Y. Lin, T. A. Tsiftsis, and Y. Yao, "An amateur drone surveillance system based on the cognitive Internet of Things," *IEEE Commun. Mag.*, vol. 56, no. 1, pp. 29–35, Jan. 2018.
- [2] X. Shi, C. Yang, W. Xie, C. Liang, Z. Shi, and J. Chen, "Anti-drone system with multiple surveillance technologies: Architecture, implementation, and challenges," *IEEE Commun. Mag.*, vol. 56, no. 4, pp. 68–74, Apr. 2018.
- [3] I. Guvenc, F. Koohifar, S. Singh, M. L. Sichertiu, and D. Matolak, "Detection, tracking, and interdiction for amateur drones," *IEEE Commun. Mag.*, vol. 56, no. 4, pp. 75–81, Apr. 2018.
- [4] H. Kang, J. Joung, J. Kim, J. Kang, and Y. S. Cho, "Protect your sky: A survey of counter unmanned aerial vehicle systems," *IEEE Access*, vol. 8, pp. 168671–168710, 2020.
- [5] R. B. Netanel, B. Nassi, A. Shamir, and Y. Elovici, "Detecting spying drones," *IEEE Secur. Privacy*, vol. 19, no. 1, pp. 65–73, Jan. 2021.
- [6] S. Park, H. T. Kim, S. Lee, H. Joo, and H. Kim, "Survey on anti-drone systems: Components, designs, and challenges," *IEEE Access*, vol. 9, pp. 42635–42659, 2021.
- [7] M. A. Khan, H. Menouar, A. Eldeeb, A. Abu-Dayya, and F. D. Salim, "On the detection of unauthorized drones—Techniques and future perspectives: A review," *IEEE Sensors J.*, vol. 22, no. 12, pp. 11439–11455, Jun. 2022.
- [8] Z. Zhang, Y. Cao, M. Ding, L. Zhuang, and W. Yao, "An intruder detection algorithm for vision based sense and avoid system," in *Proc. Int. Conf. Unmanned Aircr. Syst. (ICUAS)*, Jun. 2016, pp. 550–556.
- [9] D. Lee, W. Gyu La, and H. Kim, "Drone detection and identification system using artificial intelligence," in *Proc. Int. Conf. Inf. Commun. Technol. Converg. (ICTC)*, Oct. 2018, pp. 1131–1133.
- [10] W. Zhou, S. Gao, L. Zhang, and X. Lou, "Histogram of oriented gradients feature extraction from raw Bayer pattern images," *IEEE Trans. Circuits Syst. II, Exp. Briefs*, vol. 67, no. 5, pp. 946–950, May 2020.

- [11] K. Kim, J. Kim, H. Lee, J. Choi, J. Fan, and J. Joung, "UAV chasing based on YOLOv3 and object tracker for counter UAV systems," *IEEE Access*, vol. 11, pp. 34659–34673, 2023.
- [12] F. Gökçe, G. Üçoluk, E. Şahin, and S. Kalkan, "Vision-based detection and distance estimation of micro unmanned aerial vehicles," *Sensors*, vol. 15, no. 9, pp. 23805–23846, Sep. 2015.
- [13] P. Tang, C. Wang, X. Wang, W. Liu, W. Zeng, and J. Wang, "Object detection in videos by high quality object linking," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 42, no. 5, pp. 1272–1278, May 2020.
- [14] N. J. Sie, S. Srigrarom, and S. Huang, "Field test validations of vision-based multi-camera multi-drone tracking and 3D localizing with concurrent camera pose estimation," in *Proc. IEEE Int. Conf. Control Robot. Eng. (ICCRE)*, Apr. 2021, pp. 139–144.
- [15] P. Andrašić, T. Radišić, M. Muštra, and J. Ivošević, "Night-time detection of UAVs using thermal infrared camera," *Transp. Res. Proc.*, vol. 28, pp. 183–190, Jan. 2017.
- [16] Y. Wang, Y. Chen, J. Choi, and C.-C.-J. Kuo, "Towards visible and thermal drone monitoring with convolutional neural networks," *APSIPA Trans. Signal Inf. Process.*, vol. 8, no. 1, pp. 1–13, 2019.
- [17] P. Wellig, P. Speirs, C. Schuepbach, R. Oechslein, M. Renker, U. Boeniger, and H. Pratisto, "Radar systems and challenges for C-UAV," in *Proc. 19th Int. Radar Symp. (IRS)*, Jun. 2018, pp. 1–8.
- [18] A. Macaveiu, C. Naformita, A. Isar, A. Campeanu, and I. Naformita, "A method for building the range-Doppler map for multiple automotive radar targets," in *Proc. 11th Int. Symp. Electron. Telecommun. (ISETC)*, Nov. 2014, pp. 1–6.
- [19] J. Farlik, M. Kratky, J. Casar, and V. Sary, "Radar cross section and detection of small unmanned aerial vehicles," in *Proc. IEEE Int. Conf. Mechatron. Mechatronika (ME)*, Dec. 2016, pp. 1–7.
- [20] E. Kaya and G. B. Kaplan, "Neural network based drone recognition techniques with non-coherent S-band radar," in *Proc. IEEE Radar Conf. (RadarConf)*, May 2021, pp. 1–6.
- [21] J. Liu, Q. Y. Xu, and W. S. Chen, "Classification of bird and drone targets based on motion characteristics and random forest model using surveillance radar data," *IEEE Access*, vol. 9, pp. 160135–160144, 2021.
- [22] V. Semkin, M. Yin, Y. Hu, M. Mezzavilla, and S. Rangan, "Drone detection and classification based on radar cross section signatures," in *Proc. Int. Symp. Antennas Propag. (ISAP)*, Jan. 2021, pp. 223–224.
- [23] Y. D. Zhang, X. Xiang, Y. Li, and G. Chen, "Enhanced micro-Doppler feature analysis for drone detection," in *Proc. IEEE Radar Conf. (RadarConf)*, May 2021, pp. 1–4.
- [24] H. Kuschel, D. Cristallini, and K. E. Olsen, "Tutorial: Passive radar tutorial," *IEEE Aerosp. Electron. Syst. Mag.*, vol. 34, no. 2, pp. 2–19, Feb. 2019.
- [25] N. Souli, I. Theodorou, P. Kolios, and G. Ellinas, "Detection and tracking of rogue UASs using a novel real-time passive radar system," in *Proc. Int. Conf. Unmanned Aircr. Syst. (ICUAS)*, Jun. 2022, pp. 576–582.
- [26] J. Drozdowicz, M. Wielgo, P. Samczynski, K. Kulpa, J. Krzonkalla, M. Mordzonek, M. Bryl, and Z. Pakielaszek, "35 GHz FMCW drone detection system," in *Proc. 17th Int. Radar Symp. (IRS)*, May 2016, pp. 1–4.
- [27] M. Jian, Z. Lu, and V. C. Chen, "Drone detection and tracking based on phase-interferometric Doppler radar," in *Proc. IEEE Radar Conf. (RadarConf)*, Apr. 2018, pp. 1146–1149.
- [28] P. K. Rai, H. Idsøe, R. R. Yakkati, A. Kumar, M. Z. Ali Khan, P. K. Yalavarthy, and L. R. Cenkeramaddi, "Localization and activity classification of unmanned aerial vehicle using mmWave FMCW radars," *IEEE Sensors J.*, vol. 21, no. 14, pp. 16043–16053, Jul. 2021.
- [29] Z. Shi, X. Chang, C. Yang, Z. Wu, and J. Wu, "An acoustic-based surveillance system for amateur drones detection and localization," *IEEE Trans. Veh. Technol.*, vol. 69, no. 3, pp. 2731–2739, Mar. 2020.
- [30] J. Guo, I. Ahmad, and K. Chang, "Classification, positioning, and tracking of drones by HMM using acoustic circular microphone array beamforming," *EURASIP J. Wireless Commun. Netw.*, vol. 2020, no. 1, p. 63, Jan. 2020.
- [31] B. Kang, H. Ahn, and H. Choo, "A software platform for noise reduction in sound sensor equipped drones," *IEEE Sensors J.*, vol. 19, no. 21, pp. 10121–10130, Nov. 2019.
- [32] M. Z. Anwar, Z. Kaleem, and A. Jamalipour, "Machine learning inspired sound-based amateur drone detection for public safety applications," *IEEE Trans. Veh. Technol.*, vol. 68, no. 3, pp. 2526–2534, Mar. 2019.
- [33] S. Al-Emadi, A. Al-Ali, A. Mohammad, and A. Al-Ali, "Audio based drone detection and identification using deep learning," in *Proc. 15th Int. Wireless Commun. Mobile Comput. Conf. (IWCMC)*, Jun. 2019, pp. 459–464.
- [34] Y. Seo, B. Jang, and S. Im, "Drone detection using convolutional neural networks with acoustic STFT features," in *Proc. 15th IEEE Int. Conf. Adv. Video Signal Based Surveill. (AVSS)*, Nov. 2018, pp. 1–6.
- [35] L. Wang and A. Cavallaro, "Acoustic sensing from a multi-rotor drone," *IEEE Sensors J.*, vol. 18, no. 11, pp. 4570–4582, Jun. 2018.
- [36] S. V. Sibanyoni, D. T. Ramotsoela, B. J. Silva, and G. P. Hancke, "A 2-D acoustic source localization system for drones in search and rescue missions," *IEEE Sensors J.*, vol. 19, no. 1, pp. 332–341, Jan. 2019.
- [37] Z. Uddin, J. Nebhen, M. Altaf, and F. A. Orakzai, "Independent vector analysis inspired amateur drone detection through acoustic signals," *IEEE Access*, vol. 9, pp. 63456–63462, 2021.
- [38] I. Bisio, C. Garibotto, F. Lavagetto, A. Sciarrone, and S. Zappatore, "Unauthorized amateur UAV detection based on WiFi statistical fingerprint analysis," *IEEE Commun. Mag.*, vol. 56, no. 4, pp. 106–111, Apr. 2018.
- [39] P. Flak, "Drone detection sensor with continuous 2.4 GHz ISM band coverage based on cost-effective SDR platform," *IEEE Access*, vol. 9, pp. 114574–114586, 2021.
- [40] B. Kaplan, I. Kahraman, A. R. Ekti, S. Yarkan, A. Görçin, M. K. Özdemir, and H. A. Çirpan, "Detection, identification, and direction of arrival estimation of drone FHSS signals with uniform linear antenna array," *IEEE Access*, vol. 9, pp. 152057–152069, 2021.
- [41] W. Nie, Z. Han, M. Zhou, L. Xie, and Q. Jiang, "UAV detection and identification based on WiFi signal and RF fingerprint," *IEEE Sensors J.*, vol. 21, no. 12, pp. 13540–13550, Jun. 2021.
- [42] W. Nie, Z.-C. Han, Y. Li, W. He, L. Xie, X. Yang, and M. Zhou, "UAV detection and localization based on multi-dimensional signal features," *IEEE Sensors J.*, vol. 22, no. 6, pp. 5150–5162, Mar. 2022.
- [43] S. Yang, Y. Luo, W. Miao, C. Ge, W. Sun, and C. Luo, "RF signal-based UAV detection and mode classification: A joint feature engineering generator and multi-channel deep neural network approach," *Entropy*, vol. 23, no. 12, p. 1678, Dec. 2021.
- [44] D. Noh, S. Jeong, H. Hoang, Q. Pham, T. Huynh-The, M. Hasegawa, H. Sekiya, S. Kwon, S. Chung, and W. Hwang, "Signal preprocessing technique with noise-tolerant for RF-based UAV signal classification," *IEEE Access*, vol. 10, pp. 134785–134798, 2022.
- [45] S. Samaras, E. Diamantidou, D. Ataloglou, N. Sakellariou, A. Vafeiadis, V. Magoulianitis, A. Lalas, A. Dimou, D. Zarpalas, K. Votis, P. Daras, and D. Tzovaras, "Deep learning on multi sensor data for counter UAV applications—A systematic review," *Sensors*, vol. 19, no. 22, p. 4837, Nov. 2019.
- [46] S. Jamil, Fawad, M. Rahman, A. Ullah, S. Badnava, M. Forsat, and S. S. Mirjavadi, "Malicious UAV detection using integrated audio and visual features for public safety applications," *Sensors*, vol. 20, no. 14, p. 3923, Jul. 2020.
- [47] A. Hommes, A. Shoykhetbrod, D. Noetel, S. Stanko, M. Laurenzis, S. Hengy, and F. Christnacher, "Detection of acoustic, electro-optical and RADAR signatures of small unmanned aerial vehicles," in *Proc. SPIE*, vol. 9997, Oct. 2016, Art. no. 999701.
- [48] H. Liu, Z. Wei, Y. Chen, J. Pan, L. Lin, and Y. Ren, "Drone detection based on an audio-assisted camera array," in *Proc. IEEE 3rd Int. Conf. Multimedia Big Data (BigMM)*, Apr. 2017, pp. 402–406.
- [49] M. Aledhari, R. Razzak, R. M. Parizi, and G. Srivastava, "Sensor fusion for drone detection," in *Proc. IEEE 93rd Veh. Technol. Conf. (VTC-Spring)*, Apr. 2021, pp. 1–7.
- [50] F. Svanström, C. Englund, and F. Alonso-Fernandez, "Real-time drone detection and tracking with visible, thermal and acoustic sensors," in *Proc. 25th Int. Conf. Pattern Recognit. (ICPR)*, Jan. 2021, pp. 7265–7272.
- [51] M. Ozel, "Drone dataset (UAV)," Jul. 2021. [Online]. Available: <https://www.kaggle.com/datasets/dasmehdixr/drone-dataset-uav>
- [52] S. Jamil, "Malicious UAVs detection," Aug. 2022. [Online]. Available: <https://www.kaggle.com/sonain/malicious-uavs-detection>
- [53] S. Menard, *Applied Logistic Regression Analysis*, 2 ed. London, U.K.: SAGE Publications Ltd., 2002.
- [54] Analog Devices. (2020). *EV-TINYRAD24G User Guide*. UG-1709. [Online]. Available: <https://www.analog.com>
- [55] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2015, pp. 1–9.



- [56] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016, pp. 770–778.
- [57] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017, pp. 2261–2269.
- [58] H.-F. Yu, F.-L. Huang, and C.-J. Lin, "Dual coordinate descent methods for logistic regression and maximum entropy models," *Mach. Learn.*, vol. 85, nos. 1–2, pp. 41–75, Oct. 2011.



**HUNJE LEE** received the B.S. degree from Korea Aerospace University (KAU), Goyang-si, Republic of Korea, in February 2023. He was with the Intelligent Signal Processing Laboratory (ISPL), School of Electronics and Information Engineering, KAU, as an Undergraduate Research Assistant, where he has been focusing on field measurement for air-to-ground channels and signal processing based on deep learning techniques for UAV detection. His research interests include

air-to-ground channel modeling for low altitude UAVs, UAV detection and classification using machine learning and deep learning, and signal processing for learning-based design of future wireless communication systems.



**SUJEONG HAN** received the B.S. degree from the School of Electronics and Information Engineering, KAU, Goyang-si, Republic of Korea, in August 2023. She has been focusing on mobile communication and signal processing for wireless communication during the B.S. degree. She was an Undergraduate Research Assistant with ISPL, KAU. Her research interests include air-to-ground channel modeling for low altitude UAVs, UAV detection and classification using radar and acoustic sensors, and UAV trajectory optimization for delivery and transportation.



**JEONG-IL BYEON** received the B.S. degree from KAU, Goyang-si, Republic of Korea, in February 2023, where he is currently pursuing the M.S. degree with the Department of Electronics and Information Engineering. Since 2021, he has been joining ISPL, School of Electronics and Information Engineering, KAU, as an Undergraduate Research Assistant, where he performed research on signal processing for synthetic aperture radar (SAR) imaging and compressive sensing algorithms. His research interests include imaging techniques for drone/airborne/spaceborne SAR systems, signal processing for synthetic aperture sonar, and design of wireless communication systems aided by intelligent reflecting surface (IRS).



**SEOULGYU HAN** will receive the B.S. degree from KAU, Goyang-si, Republic of Korea, in August 2023, and he will pursue the M.S. degree in the Department of Computer Science and Engineering, Seoul National University, in September 2023. He has been focusing on radar signal processing and image signal processing for computer vision during the B.S. degree. He was an Undergraduate Research Assistant with the Media Processing Laboratory, KAU. His research interests include radar signal processing, computer vision, and signal processing for machine learning and deep learning algorithms.



**RANGUN MYUNG** received the B.S. degree from the School of Electronics and Information Engineering, KAU, Goyang-si, Republic of Korea, in February 2023. She has been focusing on mobile communication and signal processing for wireless communication during the B.S. degree. Her research interests include speech and acoustic signal processing, UAV detection and classification, and signal processing based on machine learning and deep learning.



**JINGON JOUNG** (Senior Member, IEEE) received the B.S. degree in radio communication engineering from Yonsei University, Seoul, South Korea, in 2001, and the M.S. and Ph.D. degrees in electrical engineering and computer science from KAIST, Daejeon, South Korea, in 2003 and 2007, respectively.

He was a Postdoctoral Fellow with KAIST and UCLA, CA, USA, in 2007 and 2008, respectively.

He was a Scientist with the Institute for Infocomm

Research, Singapore, from 2009 to 2015, and joined Chung-Ang University (CAU), Seoul, in 2016, as a Faculty Member. He is currently a Professor with the School of Electrical and Electronics Engineering, CAU, where he is also a Principal Investigator with the Intelligent Wireless Systems Laboratory. His research interests include signal processing, numerical analysis, algorithms, and machine learning.

Dr. Joung was recognized as an Exemplary Reviewer of the IEEE COMMUNICATIONS LETTERS, in 2012, and the IEEE WIRELESS COMMUNICATIONS LETTERS, from 2012 to 2014 and in 2019. He served as the Guest Editor for the IEEE ACCESS, in 2016. He served on the editorial board for the *APSIPA Transactions on Signal and Information Processing*, from 2014 to 2019, served as a Guest Editor for the *Electronics*, in 2019, and served as an Associate Editor for the IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY, from 2018 to 2023. He is an Inventor of a *Space-Time Line Code (STLC)* that is a fully symmetric scheme to a space-time block code.



**JIHOON CHOI** (Senior Member, IEEE) received the B.S., M.S., and Ph.D. degrees from the Korea Advanced Institute of Science and Technology (KAIST), Daejeon, South Korea, in 1997, 1999, and 2003, respectively.

From 2003 to 2004, he was with the Department of Electrical and Computer Engineering, The University of Texas at Austin, where he performed research on multiple antenna systems as a Postdoctoral Fellow. From 2004 to 2008, he was with

Samsung Electronics, South Korea, where he worked on developments of commercial radio access stations for M-WiMAX and base stations for CDMA 1xEV-DO Rev.A/B. In 2008, he joined KAU, Goyang, South Korea, as a Faculty Member, where he is currently a Professor with the School of Electronics and Information Engineering. He is also a Chief Investigator with ISPL, KAU. His research interests include signal processing for wireless communications, air-to-ground channel modeling, radar signal processing, UAV trajectory optimization, machine learning, modem design for future cellular networks, wireless LANs, the IoT devices, and digital broadcasting systems.

...