

Fusion of Infrared and Visible Images for Remote Detection of Low-Altitude Slow-speed Small Targets

HaiJiang SUN*, Qiaoyuan LIU*, Jiacheng WANG, Jinchang REN, Yanfeng WU, Huimin Zhao, and Huakang Li

Abstract—Detection of the low-altitude, slow-speed, small (LSS) targets is one of the most popular research topics in remote sensing. Despite of a few existing approaches, there is still an accuracy gap for satisfying the practical needs. As the LSS targets are too small to extract useful features, deep learning based algorithms can hardly be used. To this end, we propose in this paper an effective strategy for determining the region of interest (ROI), using a multi-scale layered image fusion method to extract the most representative information for LSS target detection. In addition, an improved Self-Balanced Sensitivity Segment model (SuBSENSE) is proposed to detect the fused LSS-Target, which can further improve both the detection accuracy and the computational efficiency. We conduct extensive ablation studies to validate the efficacy of the proposed LSS target detection method on three public datasets and three self-collected datasets. The superior performance over the-state-of-the-arts (SOTA) has fully demonstrated the efficacy of the proposed approach.

Index Terms—Image fusion; LSS Target detection; saliency detection; background subtraction.

I. INTRODUCTION

THE "Low-altitude, small and slow" target (LSS-Target), such as Unmanned Aerial Vehicles (UAVs), is a general term for small aviation device/equipment with a flight altitude less than 2km and a flight speed less than 50 km/h. Detection of the LSS-targets is a key technology for precision navigation, infrared search, tracking and reconnaissance warning system, which has been widely applied into many practical applications such as air transportation [1], aerial mapping [2] and many tasks in military confrontation [3]. As a result, the LSS-Target detection indeed has important research significance in the field of ground air defense. This approach could directly determine the operating distance and detection sensitivity of the corresponding systems. Although great breakthroughs have been made in recent years [4-7,41,42], including the background subtraction algorithm based on the visible light image [6], detecting the LSS-Target with the combination of visible light image, infrared image and even hyperspectral image [8], but there are still great challenges for putting this approach into practical applications.

*The first two authors are of equal contributions.

Manuscript received: 12/09/2020, revised: 25/10/2020. This work is partially supported by the Key Laboratory of Airborne Optical Imaging & Measurement, Chinese Academy of Sciences (CAS) and International Cooperation Project of Changchun Institute of Optics, Fine Mechanics and Physics (CIOFMP) (Y9U933T190), the Dazhi Scholarship of the Guangdong Polytechnic Normal University, National Natural Science Foundation of China (62072122), and the Education Department of Guangdong Province (2019KSYS009).



Fig. 1 LSS targets in different aviation scenarios.

As the LSS-Targets usually appear like a spot, see in Fig. 1, the detection can be easily affected by the background objects. In this paper, we summarize the difficulties of the LSS-Target detection into three aspects. The first is the complex background, and how to accurately detect the LSS-targets whilst overcoming the effects of complex scenes is always a challenge in this context. The second is the weak target, as it is very hard to determine the specific trajectories for LSS-Targets under low-speeds. Due to their small sizes in the field of view, it is difficult to build accurate templates for target detection. The third is the high false alarm rate: considering that there always exist faked objects with similar features to the LSS targets, how to accurately differ the LSS-targets from complex scenes to reduce the false alarms is another major issue.

To tackle the aforementioned challenging problems, various methods have been proposed for LSS-target detection. According to the image sources used, these approaches can be divided into two categories, i.e. LSS-target detection in the visible light image and infrared images, respectively. Lou et al. [9] introduced the saliency and regional stability for feature extraction in visible-light images. A segmentation threshold was used to distinguish the target for accurate detection, but it was unsuitable for complex scenes. Xie et al. [10] proposed the Peer Group Filter (PGF) to improve the signal-to-noise ratio (SNR) of the infrared image in LSS-target detection. Although small targets could be extracted from the background accurately, it was too time-consuming for real-time tasks. Pan et al. [38] proposed a double-layer local contrast measure (DLCM) approach for small target detection in infrared images. A double-layer-diagonal gray difference contrast was used to enhance the visual saliency of the target and alleviate the impact of the background clutter and noise.

In general, the images captured by the visible cameras show clear texture, edge and high spatial resolution than the infrared

Corresponding authors: J. Ren (jinchang.ren@ieee.org) and H. Zhao (zhaohuimin@gpnu.edu.cn).

H. Sun, Q. Yuan and J. Wang are with CIOFMP, Changchun, China.

J. Ren, H. Zhao and H. Li are with School of Computer Sciences, Guangdong Polytechnic Normal University, Guangzhou, China. J Ren is also with the National Subsea Centre, Robert Gordon University, Aberdeen, UK.

Y. WU is with the 28th Research Institute of China Electronics Technology Group, Nanjing 210001, China.

ones, which are useful for target detection. However, the quality of the image produced can be affected by the weather conditions and environmental illumination et al., leading to failure of detection at night or during bad weather conditions. However, the infrared imaging can supplement the deficiency of visible images. Although the infrared image suffers from a relatively lower spatial resolution and weaker textures, it can produce particularly better images at night or in bad weathers.

In order to enhance the generality of LSS target detection, the combination of both the visible images and the infrared images have been widely used. Durga et al. proposed image fusion and secondary decomposition of the image based LSS-target detection [11]. Shi et al. [12] applied combined strategies on the visible and infrared images for accurate target extraction, in which an inter frame difference based rough entropy model was used for locating of the moving target in the infrared images, along with a Local Binary Pattern (LBP) model for the visible images and a Gaussian Mixture Model (GMM) for background modelling. Although these have improved the accuracy and generality of the algorithm, the information of visible and infrared images was not fully utilized.

To tackle the effect of the complex background in LSS-Target detection, various background modelling approaches have been proposed, either in the visible light images or the infrared images. Olivier et al. [13] proposed an unordered background modeling algorithm, which stored a sample set for all pixels and randomly selected the neighboring pixels for updating the samples. This method was robust to noise but could hardly deal with color changes such as illumination and ghost effect. The approach was further improved in Hofmann et al. [14], where a pixel level nonparametric background-modeling algorithm was proposed. Although it could adaptively detect small targets in complex background, the algorithm was quite complex and time-consuming.

In this paper, combining infrared and visible images for LSS-Target detection is also focused. Different from the existing methods, our proposed approach fuses the features and information from visible and infrared images in multi-layers. First, a background model is used from the perspective of inter frame difference from three consecutive frames rather than two, which can help to obtain the accurate motion region whilst greatly decreasing the impact of background. Second, the weighted moving average (WMA) [15] and weighted moving variance (WMV) [16] are utilized to obtain the target candidate regions from the infrared and visible images, respectively. Third, based on the target candidate regions, a self-balanced sensitivity segmentation algorithm (SuBSENSE) [17] is used to construct a local area background model for refined LSS-target detection. The Toet A.TNO (Image fusion dataset) dataset [18] and the Visible-Infrared Database (VID) dataset [19] are used for evaluating the performance of image fusion and target detection, respectively. Experimental results have validated the efficacy of the proposed model for LSS-target detection in comparison to the state-of-the-art (SOTA) approaches.

II. THE PROPOSED ALGORITHM

The diagram of the proposed LSS-Target detection approach is illustrated in Fig. 2, which has two modules for fusion based feature extraction. The ROI extraction module is designed to predict the candidate target region with multi-detector images. Given three frames of the infrared images and visible images, WMA and WMV can output a candidate target region, denoted as $rec1$ and $rec2$ respectively. The minimum bounding box of these two regions can be taken as the extracted ROI that carries target information from both infrared and visible images. In the multi-scale layered image fusion module, a fused image is generated for LSS target detection by saliency maps from the

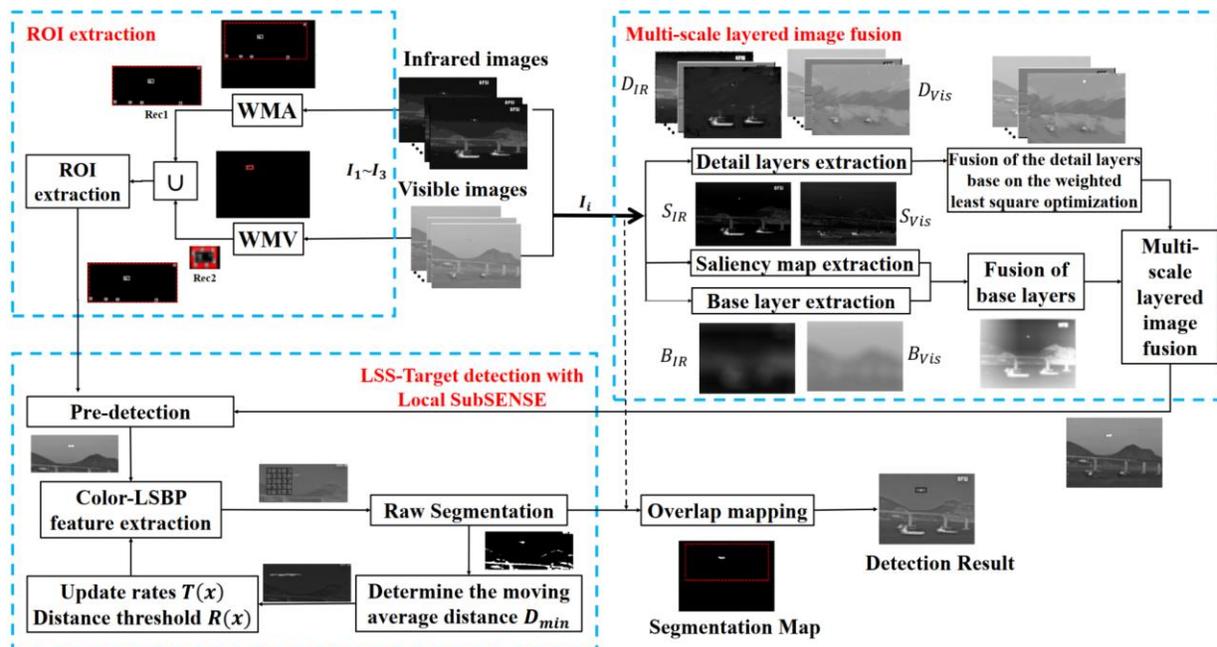


Fig. 2 The flowchart of the proposed LSS-Target detection approach.

extracted base layers and detail layers. The detected ROIs and the fused image are then inputted to an improved SubSENSE module for further refining the locations of the LSS targets.

Specifically, we firstly apply the background subtraction to the first three frames of the visible and its corresponding infrared image sequences to determine the initial candidate regions of targets. Here we denote the first three frames of the infrared image and the visible images as $I_{IR} = \{I_{IR_1}, I_{IR_2}, I_{IR_3}\}$ and $I_{Vis} = \{I_{Vis_1}, I_{Vis_2}, I_{Vis_3}\}$, respectively. By applying the WMA [15] and WMV [16] methods on I_{IR} and I_{Vis} , a rough foreground and background segmentation can be obtained, which can provide a good guidance for the subsequent target detection. Specifically, by applying the WMA on I_{IR} , several foreground masks can be obtained, where we use $rec1$ to denote the bounding box of the targets. Similarly, a bounding box $rec2$ of detected targets from the visible images can be determined from by applying WMV on the visible images. The final ROI of the target candidate region R can be obtained via a union of two sets $rec1$ and $rec2$ as follows, and the details will be presented in Section III.

$$rec1 = WMA(I_{IR}) \quad (1)$$

$$rec2 = WMV(I_{Vis}) \quad (2)$$

$$R = rec1 \cup rec2 \quad (3)$$

The target candidate region generated by the background subtraction model will be applied to each subsequent fusion image. In this way the detection algorithm SubSENSE [17] only needs be applied to the determined ROI rather than the whole image for improved efficiency.

In our approach, image fusion is the key to make the full use of infrared and visible images. Firstly, the Rolling Guide Filter (RGF) is applied to the infrared and visible images, separately, in order to decompose the images and enhance the details. This is followed by the visual saliency to fuse the decomposed image layer as the input of SubSENSE [17] detection method.

$$\{D_{IR}, B_{IR}\} = RGF(I_{IR_i}), D_{IR} = \{d_{IR}^1, d_{IR}^2, \dots, d_{IR}^N\} \quad (4)$$

$$\{D_{Vis}, B_{Vis}\} = RGF(I_{Vis_i}), D_{Vis} = \{d_{Vis}^1, d_{Vis}^2, \dots, d_{Vis}^N\} \quad (5)$$

$$\{S_{IR_i}, S_{Vis_i}\} = Saliency(I_{IR_i}, I_{Vis_i}) \quad (6)$$

where D_{IR} and B_{IR} denote respectively the extracted N detailed layers from the infrared image and base layer; D_{Vis} and B_{Vis} denote respectively the extracted N detailed layers and base layer from the visible image. As shown in Fig. 2, all the extracted detail layers D_{IR} and D_{Vis} are fused according to the weighted least square optimization, and all the extracted base layers B_{IR} and B_{Vis} combined with the corresponding saliency map S_{IR} and S_{Vis} for fused. Combine all the fused images as shown in Eq. (25), we can get the final fused image F_i .

Based on the background subtraction model R and the fused image F_i , the position of the LSS-target in the current frame can be pre-detected, which is then further refined by an improved SubSENSE [17] method as detailed in Section V.

$$l = \text{Sub}(R, F_i) \quad (7)$$

Finally, the refined mask l is mapped to the original image to produce the segmentation map and the final detection results.

III. EXTRACTION OF CANDIDATE TARGET REGIONS

For accurate detection of the moving targets, we first

introduce a three-frame image difference algorithm to extract the candidate regions of the targets, in which the adjacent three frames are taken as a group for calculating the image difference for robustness. Considering the low speed of the LSS target, only the initial three frames are used for background modelling.

A. ROI extraction in infrared images

For the infrared images, the WMA algorithm is employed for candidate target region extraction, which is more suitable for processing infrared images due to its fast calculation speed and sensitivity to illumination changes. WMA uses the weighted average of the initial three infrared frames to build a background subtraction model by replaying more weights on the observation values closer to the predicted time.

$$I_{mean} = \sum_{i=1}^3 w_i * I_i \quad (8)$$

where I_{mean} is the weighted average of the image pixels, w_i is the corresponding weight with $\sum_{i=1}^3 w_i = 1$, in our approach we set $w_1 = 0.5, w_2 = 0.3, w_3 = 0.2$ as an example.

As infrared images conform to the laws of thermodynamics, the 1-D information entropy is used as a threshold to distinguish the foreground from the background as follows:

$$H_b = - \sum_{i=0}^{\delta} \frac{p_i}{W_1(\delta)} \times \log\left(\frac{p_i}{W_1(\delta)}\right) \quad (9)$$

$$H_f = - \sum_{i=\delta+1}^{L-1} \frac{p_i}{W_2(\delta)} \times \log\left(\frac{p_i}{W_2(\delta)}\right) \quad (10)$$

where δ is the threshold to distinguish the gray values of the background before and after; W_1 and W_2 are the gray value probability of the background and the target, respectively, and L is the number of gray levels of the infrared image. p_i is the probability distribution of image gray levels, which can be computed by $p_i = \frac{n_i}{M \times N}$, where $M \times N$ is the total number of pixels, n_i indicates the number of pixels with gray level i .

The segmentation threshold can then be determined by:

$$\begin{cases} H_t(\delta) = H_b + H_f \\ \delta^* = \text{argmax} H_t(\delta) \end{cases} \quad (11)$$

where $H_t(\delta)$ is the total information entropy of the image, and δ^* is the determined segmentation threshold by maximizing the 1-D information entropy of the image [39].

For detecting LSS targets, firstly pixel-based difference between the current frame I_i and the corresponding weighted average image I_{mean} is obtained as I_{sub} . By comparing I_{sub} against the threshold δ^* , the foreground can be determined by

$$B_t = \begin{cases} 1, & I_{sub} > \delta^* \\ 0, & I_{sub} \leq \delta^* \end{cases}, I_{sub} = |I_i - I_{mean}| \quad (12)$$

The algorithm can then dynamically update the background model based on the weighted background model and the current frame to adapted the model to scene changes as follows:

$$B_{t+1} = (1 - \alpha)B_t + \alpha I_t \quad (13)$$

where α is the update learning rate, and I_t is the t^{th} frame; B_t and B_{t+1} denote the background model at time t and $t+1$.

Although the foreground can be extracted from the infrared image above, it is relatively too rough to estimate the final target location. However, it helps to provide a coarse location of the

target for efficiency. As shown in Fig. 2 the binary image outputted by WMA, the white area represents the extracted foreground, and the remaining represents the background. We first determine the bounding box for all the white areas as $rec1$, which is the candidate target region of the infrared image corresponding to the current background model.

B. ROI extraction in visible images

For the visible images, based on the original WMA background modeling, the mean I_{mean} and the standard deviation is calculated. As the visible images always contain rich texture features, this is namely weighted moving variance (WMV) background modeling.

$$\text{var}(I_i, w_i) = w_i * (I_i - I_{mean})^2 \quad (14)$$

$$\text{std} = \sqrt{\sum_{i \in [1,3]} \text{var}(I_i, w_i) / 3} \quad (15)$$

where $\text{var}(\cdot)$ and std are the weighted variance and weighted standard deviation of the image pixels, respectively. Besides, other background modeling processes are as same as the infrared images, and w_i follows the same distribution as the WMA does. Therefore, the candidate target region $rec2$ can also be generated from the visible image, where the final target candidate region could be obtained using Eq. (3).

IV. MULTI-LAYER IMAGE FUSION BASED ON ROLLING GUIDED FILTERING

Considering that the infrared images and the visible images contain different but supplement information, our approach proposes to fuse them in multi-layers to make the full use of their information. In our approach, the Rolling Guide Filter (RGF) [20] is used for multi-scale decomposition and detail enhancement. Specifically, the visible and infrared images are firstly decomposed into the corresponding base layer and detail layers, which are then fused separately. Accordingly, this could not only retain and enhance more useful image details, but also enable a favorable detection of LSS targets.

A. Extraction and fusion for the base layer

The RGF based decomposition is an iterative process [40], which can decompose an image into N layers, where the filtered detail image of the j^{th} layer is given by:

$$\mu^j = RGF(\mu^{j-1}, \sigma_s^{j-1}, \sigma_r, T) \quad (16)$$

$$d^j = \mu^{j-1} - \mu^j \quad (17)$$

where μ^j is the image after the j^{th} layer filtering, d^j is the decomposing images of the j^{th} layer, σ_s^{j-1} is the scale parameter, σ_r is the weight range parameter, and T is the number of iterations.

Nevertheless, Gaussian filter rather than RGF is applied in our approach to obtain the base-layer image. The RGF algorithm utilized in our approach lies in the advantages of its scale-aware and edge preserving properties. On the contrary, the base layer is the coarsest version of the source image, which is used to control the global contrast and appearance of the fused image. As a result, it is unnecessary to apply RGF to the base layer as there is no need to preserve edge or detail

information for this layer. For efficiency, the Gaussian filtering with a larger standard deviation is chosen in our method to get the base layer in the last filter layer with $j = N$ below:

$$\mu^N = \text{Gaussian}(\mu^{N-1}, \sigma_s^{N-1}, \sigma_r, T) \quad (18)$$

$$d^N = \mu^{N-1} - \mu^N \quad (19)$$

By setting $\sigma_s^N = 2\sigma_s^{N-1}$, the infrared base layer image B_{IR} and the visible base layer image B_{Vis} can be obtained. For fusion of the base layer, the fusion rules are given as follows:

$$\begin{cases} B_F = W_b B_{IR} + (1 - W_b) B_{Vis} \\ W_b = \frac{1}{2} (1 + S_{IR} + S_{Vis}) \end{cases} \quad (20)$$

where S_{IR} and S_{Vis} are the normalized saliency images generated from I_{IR} and I_{Vis} , using the widely used saliency map extraction method visual saliency map (VSM) [40]. With the fusion weight W_b calculated by S_{IR} and S_{Vis} , the base layer fusion image can be finally determined.

B. Extraction and fusion for the detail layers

For the fusion of the detail layers, the fusion rules are as follows. Firstly, the initial fusion detail layer M^j is obtained by:

$$M^j = W^j d_{IR}^j + (1 - W^j) d_{Vis}^j \quad (21)$$

where the fusion weight can be obtained using:

$$W^j = \begin{cases} 1 & |d_{IR}^j| > |d_{Vis}^j| \\ 0 & \text{otherwise} \end{cases} \quad (22)$$

The weighted least square method is used to obtain the final detail fusion layer D_i^j as follows:

$$\min \sum_n [(D_n^j - M_n^j)^2 + \lambda \cdot \Lambda \cdot (D_n^j - (d_{Vis}^j)_n)^2] \quad (23)$$

where $(D_n^j - M_n^j)^2$ is the Euclidean distance between the fusion detail layer D_n^j and the initial fusion detail layer M_n^j , n denotes the spatial location of a pixel, λ is a balance control parameter, $\Lambda = (\sum_n \omega_n (d_{IR}^j)_n + 0.0001)^{-1}$ is the spatial-varying weight parameter, ω_n represents the square sliding window centered at pixel n .

A larger window would not only blur the fused image but also generate a higher dimensional spatial-varying weight, resulting in increased computational cost. On the other hand, a smaller window cannot effectively remove the noise and irrelevant IR details. After testing, the size of 7×7 is used, as it produces satisfactory results, and Eq. (23) can be re-write as:

$$(1 + \lambda A^j) D^j = M^j A^j d^j \quad (24)$$

where A^j is a diagonal matrix containing the weights of all pixels. Finally, the fused image F is obtained by:

$$F = B_F + D^1 + D^2 + \dots + D^N \quad (25)$$

V. LSS-TARGET DETECTION WITH LOCAL SUBSENSE

To extract LSS-targets from the candidate target regions of the fused image, we integrate an improved SuBSENSE algorithm [17] with a proposed background modeling in local regions. The input of the background model is the fused image, and the modeling area is the comprehensive target candidate areas as determined in Sec. III. The entire detection process consists of five steps, i.e. background model initialization, recursive moving average distance calculation, noise suppression, local distance threshold control, and background

model update. Besides, a false alarm eliminating algorithm is also introduced for further improving the detection accuracy.

A. The improved SubSENSE detection algorithm

Step 1: Based on the spatio-temporal binary similarity and the Color-LBSP descriptor [22], SubSENSE could characterize the pixel representation with a non-parameter model. We initialize the background model by:

$$S_t(x) = \begin{cases} 1, & \text{if } \{ \text{dist}(F_t(x), B(x)) < R_{max} \} < m \\ 0, & \text{otherwise} \end{cases} \quad (26)$$

where $F_t(x)$ is the fused image at time t , $B(x)$ is the historical sample, the element in the background model $S_t(x)$ actually is the segmentation result, and R_{max} is the maximum distance threshold; m is the requested minimum number of matches for background classification; $\text{dist}(F_t(x), B(x))$ is the distance between the current observation and the given background. By reasonably selecting the maximum distance threshold, we can better resist unrelated changes of the model.

Step 2: Once the background model is determined, the recursive moving average distance between the current pixel and the pixel in the sample set can be calculated as follows:

$$D_{min}(x) = D_{min}(x) \cdot (1 - \alpha_2) + d_t(x) \cdot \alpha_2 \quad (27)$$

where α_2 is the learning rate and $d_t(x)$ is the minimum normalized Color-LBSP between all samples in $F_t(x)$ and $B(x)$. Since $D_{min}(x)$ is bound to $[0,1]$, an entirely static background region would have $D_{min}(x) \approx 0$, and a dynamic region to which the model cannot adapt to would have $D_{min}(x) \approx 1$. Areas with foreground objects show higher D_{min} values, because foreground detection is defined by the difference between the pixel model and local observations.

Step 3: Before applying local distance threshold control, we define a 2D map of pixel-level accumulators, v , to improve the detection accuracy. For every new segmented frame, we can obtain the binary map of all blinking pixels at time t from the background model S_t . Finally, we update v as follows:

$$v(x) = \begin{cases} v(x) + 1, & \text{if } S_t(x) = 1 \\ v(x) - 0.1, & \text{otherwise} \end{cases} \quad (28)$$

For regions with little labeling noise we will typically have $v(x) \approx 0$, while for regions with unstable labeling we have large positive $v(x)$ values. This method can guide the dynamic motion background to trigger the feedback mechanism, whilst conventional models cannot provide such feedback.

Step 4: In order to dynamically select the suitable thresholds, the local distance in the SubSENSE [17] background model is used. The local distance thresholds R can be recursively adjusted for each new frame as shown below:

$$R(x) = \begin{cases} R(x) + v(x), & \text{if } R(x) < (1 + 2 \cdot D_{min}(x))^2 \\ R(x) - \frac{1}{v(x)}, & \text{otherwise} \end{cases} \quad (29)$$

where $R(x)$ are a set of continuous values; the exponential relation between $R(x)$ and $D_{min}(x)$ is chosen over a linear relation as it favors sensitive behavior in static regions (and thus helps to generate sparse segmentation noise), but also provides robust and rapidly scaling thresholds elsewhere.

Here, the segmentation noise indicator $v(x)$ is used as a factor which, in dynamic regions, allows faster threshold increments and can even freeze $R(x)$ in place when $D_{min}(x)$ recedes to

lower values. This parameter control method is conducive to generating sparse segmentation noise and helps to provide reliable and rapidly expanding thresholds elsewhere.

Step 5: To overcome the influence of lighting, shadows, and moving targets on the detection results, the SubSENSE background model is further updated as follows, where $T(x)$ represents the model update rate,

$$T(x) = \begin{cases} T(x) + (v(x) \cdot D_{min}(x))^{-1}, & S_t(x) = 1 \\ T(x) - v(x)(D_{min}(x))^{-1}, & S_t(x) = 0 \end{cases} \quad (30)$$

If the current pixel is classified as a background, then any random pixel in the background sample has $\frac{1}{T}$ probability of being replaced by the current pixel, and its neighboring pixels have $\frac{1}{T}$ probability of being randomly replaced by any value in its neighborhood sample. Random replacement of samples and random updates of pixels can prevent static foreground objects from being quickly absorbed and ensure the authenticity of short-term and long-term background representations in the model. The spatial consistency of the background model can suppress the impact of camera shake on LSS-target detection.

B. The false alarm eliminating algorithm

In our approach, four evaluation indexes are used to reduce the false alarms, including the area ratio change, the centroid distance, the filling rate difference, and the aspect ratio difference as detailed below.

1) To automatically obtain the smallest bounding boxes of all candidate targets;

2) Calculate the target area ratio:

$$\Delta s = \left| \frac{|t|}{K(x, y)} \right| \quad (31)$$

3) Calculate the centroid distance:

$$\Delta c = \left\| \frac{c_x}{b_w} + \frac{c_y}{b_n} \right\|^2 \quad (32)$$

4) Calculate the aspect ratio:

$$\Delta a = b_w / b_n \quad (33)$$

where $|t| = \sum K(x, y)$ denotes the number of pixels of the target in the image; $K(x, y)$ represents the area of the bounding box of the candidate target; b_w and b_n denote respectively the width and height of that bounding box.

Specifically, the aspect ratio reflects the morphological state of the target to some extent as the ideal LSS target appears like a spot, i.e. the aspect ratio is close to 1. If the width or height are much disparity to each other, the suspected area becomes a long strip, which will unlikely be determined as an LSS-target even if it meets other constraints. Therefore, it will be classified as a false target to be eliminated. In our approach, the aspect ratio that can tolerate false targets is set to 0.35. Assume the set R contains m candidate targets $\{r_1, r_2, \dots, r_m\}$, and T_f is the target set after removing the false alarms. The specific process of eliminating the false targets is given below.

FALSE ALARMS ELIMINATING ALGORITHM

Input: Candidate target set R

1. $t_i = \{r_i\}, T_i = \{\emptyset\}$ Δ Initialization
2. For $i = 1, r_i \in R$

3. If $\forall S_i: \begin{cases} 0.8 \leq \Delta s < 1 \\ 0.75 \leq \Delta c < 1 \text{ Then} \\ \Delta a \geq 0.35 \end{cases}$	
4. $T_{new} = T_{exist} \cup t_i$	Δ Update target Set
5. Else $T_{exist} = T_f$	Δ Eliminate false alarms
End for	
Output: Target set T_t	

VI. EXPERIMENTAL RESULTS AND ANALYSIS

In this paper, three groups of multi-source images are selected for performance assessment. Video 1 is the ‘Take 1 sequence’ of the ‘Guanabara Bay: Outdoor’ in the Toet A.TNO (Image fusion dataset) dataset [18], which is referred to as Guanabara Bay in our paper. Video 2 is the ‘fire_sequence’ of the ‘DHV’ in the Toet A.TNO dataset, denoted as DHV in our paper. Video 3 is the ‘Duine_Image’ of ‘FEL_Image’ in the Visible-Infrared Database (VID) dataset [19], which is referred to as FEL in our paper. The specific information of the test datasets in terms of the camera scene and image resolution is shown in Table I.

TABLE I. TEST VIDEO INFORMATION

Test Video Set	Scenes	Frame size
Guanabara Bay	Forest fire	460×380
DHV	Field surveillance	360×270
FEL	Port surveillance	720×480

Due to the particularity of the LSS target detection problem, there are rare suitable datasets publicly available. In order to enhance the persuasiveness of the proposed method, three self-acquired datasets are added, which are named as V1, V2, and V3 in the paper. All the self-acquired videos are focus on the UVA surveillance, and the specific information in terms of image resolution and image resolution is shown in Table II.

For verifying the efficacy and robustness of the proposed approach, we have collected our own datasets from UAVs under various challenging scenarios, such as small-scale, fast-moving, similar background, dark illumination and complex background etc. Among them, the dataset V1 contains small-scale targets, and V2 is for targets of similar background and fast-moving. As seen, the drone in V2 has similar intensity values to the surrounding pixels, whilst the target is blurred due to the rapid rotational movement of its propeller, leading to difficulty for target detection. As for the V3 dataset, it is mainly for dark illumination and complex background, where the trees and the monitoring machines in the background may introduce false targets and extra difficulty for accurate target detection.

The parameter settings of our method are given below. The number of decomposition levels is typically set as $N=4$, which is the same as most of other fusion methods [40] and has produced satisfactory fusion results. The initial spatial weight is set as $\sigma_s^0 = 2$. Generally, the value of λ is in the range of $[0.005, 0.02]$, and we set $k=0.01$ in this paper.

The experimental results of the extracted candidate targets are shown in Fig. 3. As seen in Fig. 3(a), the target is obviously visible in the visible image but invisible in the infrared image. In Fig. 3(b), the visibility of the target is lower in the visible

image but more obvious in the infrared image. In Fig. 3(c), the target has a rough outline in the visible image and more obvious in the infrared image. However, only aircraft engines with high heat can be observed in infrared images, and image information of other parts seems missing. The corresponding results of the self-acquired result are shown in Fig. 3(d)-(f). The experimental results show that integrating the candidate target regions of infrared and visible images can reduce the missed detection from poor images with a single sensor.

In order to demonstrate the effectiveness of the proposed fusion method, seven classic fusion algorithms including WLT [23], CWT [24], CVT [25], NSCT [26], ADF [27], MSVD [28] and WLS [29] are compared. The fusion results are shown in Figs. 4-9, in which our method can extract more texture and intensity features in the fused image for more effective target detection. After applying the fusion algorithm, almost all the targets can be identified in the six test sequences, thanks to the enhanced saliency measurement, and the amount of image information can also be significantly enriched after fusion.

The detailed results are shown respectively in Tables III-VIII, in which the best results are all bolded. As seen, the proposed fusion algorithm could perform well, and obviously better than other algorithms in combination.

The idea of detecting the LSS-Target on the fused image and performing background modeling in the candidate target region was also validated as it could greatly increase the detection speed and reduce the false alarm rate. The comparison of the detection result with different SOTA algorithms are shown in Figs. 10-11, respectively, where (a) is the input image, (b) is the detection result of the LOBSTER algorithm [36], (c) is the detection result of the PBAS algorithm [37], (d) is the detection result of the (DLCM) algorithm [38], and (e) is the detection result of the proposed algorithm. As seen, the proposed algorithm can detect the complete LSS-target more accurately.

In the process of detection, multiple targets may be detected, however not every target detected could be considered as LSS-target. So, for quantitative analysis, the accuracy and F value are used to evaluate the detection performance of the approach

TABLE II. SELF-ACQUIRED VIDEO INFORMATION

Test Video Set	V1	V2	V3
Frame Size	320×240	480×320	480×640

TABLE III Fusion comparison of Guanabara Bay with the best results in each group highlighted in bold.

Video1	AVG	IE	Q ^E	SF	SD	MI	f/s
WLT	1.39	5.67	0.35	5.38	18.86	11.35	0.75
CWT	2.56	5.92	0.45	9.91	22.33	11.83	1.13
CVT	2.59	5.99	0.43	9.86	22.27	11.99	1.90
NSCT	2.58	5.94	0.40	9.925	22.67	11.88	4.39
ADF	1.73	5.77	0.39	6.180	19.09	11.54	1.02
MSVD	2.10	5.75	0.41	8.109	19.33	11.51	0.47
WLS	2.31	5.61	0.46	10.00	20.49	11.22	3.05
OUR	2.69	5.86	0.46	10.03	22.77	11.76	4.94

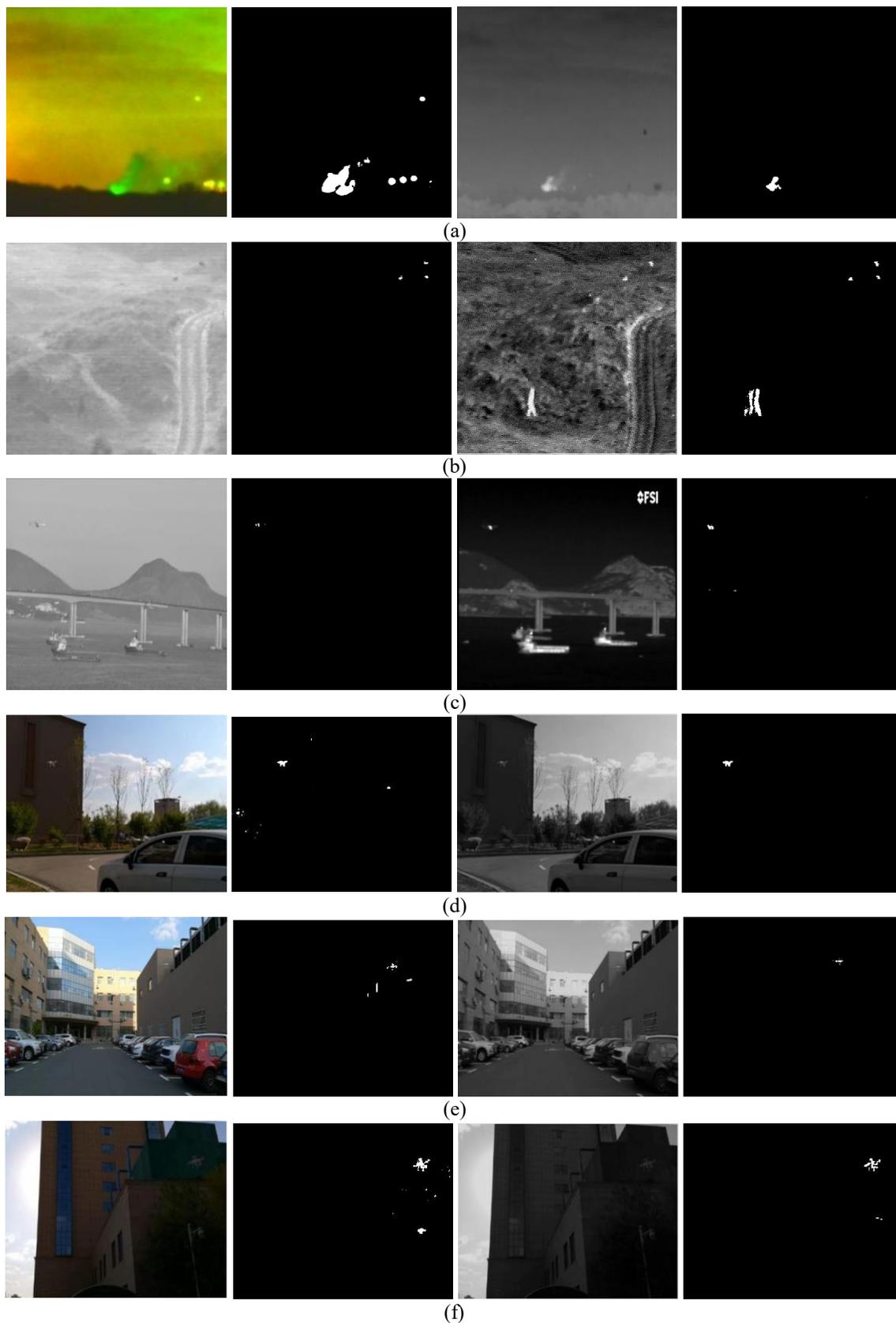


Fig. 3. Visual results of LSS-target detection from visible and infrared images with a poor resolution, separately. The first column is the original visible images, the second column is the extracted candidate regions from visible images, and the third column is the original infrared image, the fourth column is the extracted candidate regions from infrared images.

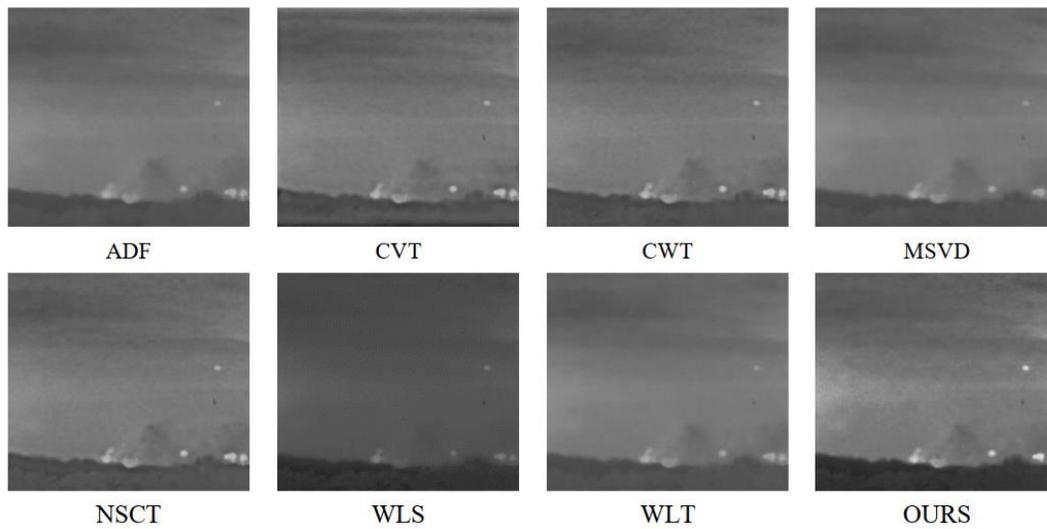


Fig. 4 Fusion comparison on the frames of the Guanabara Bay sequence

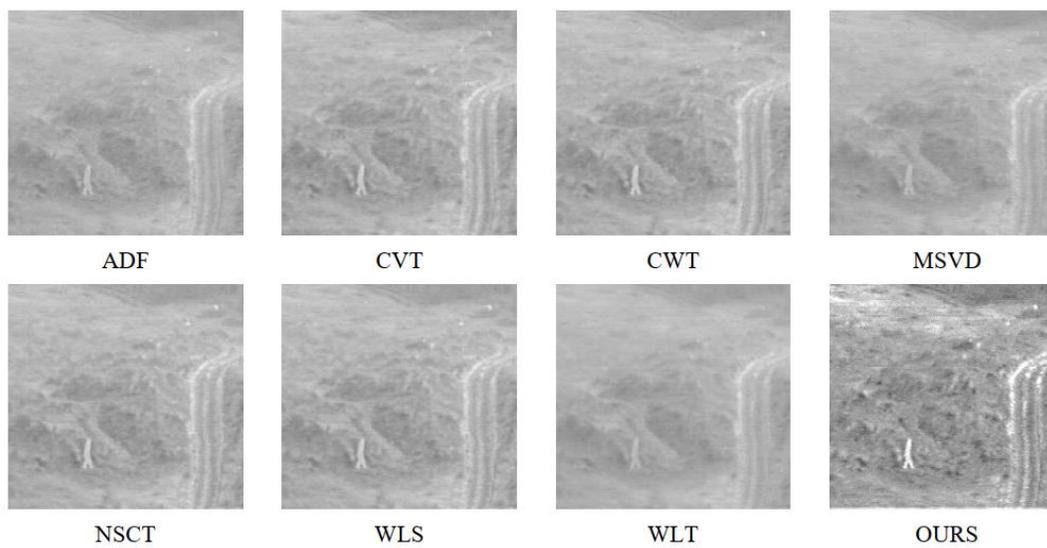


Fig. 5 Fusion comparison on the frames of the DHV sequence

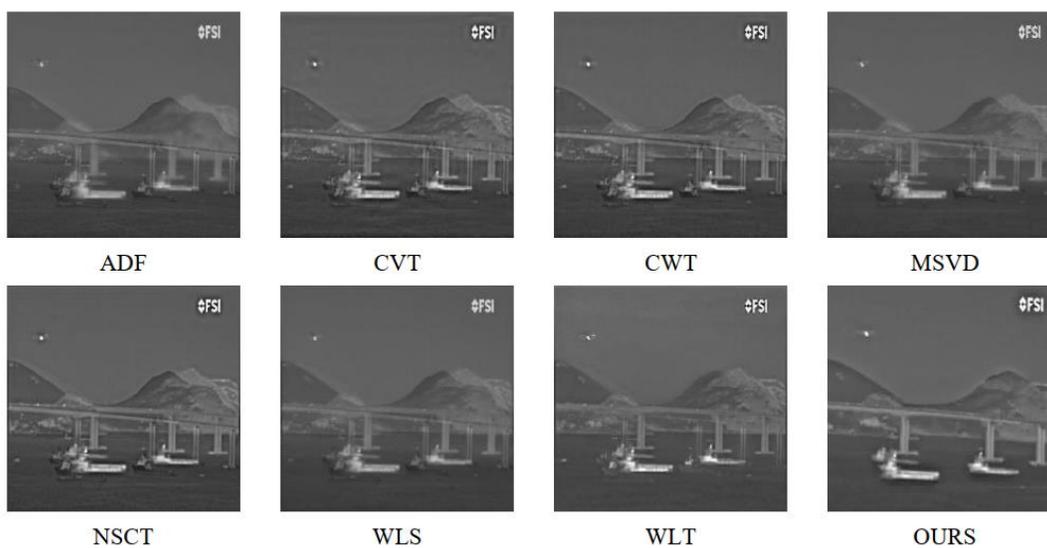


Fig. 6. Fusion comparison on the frames of the FEL sequence

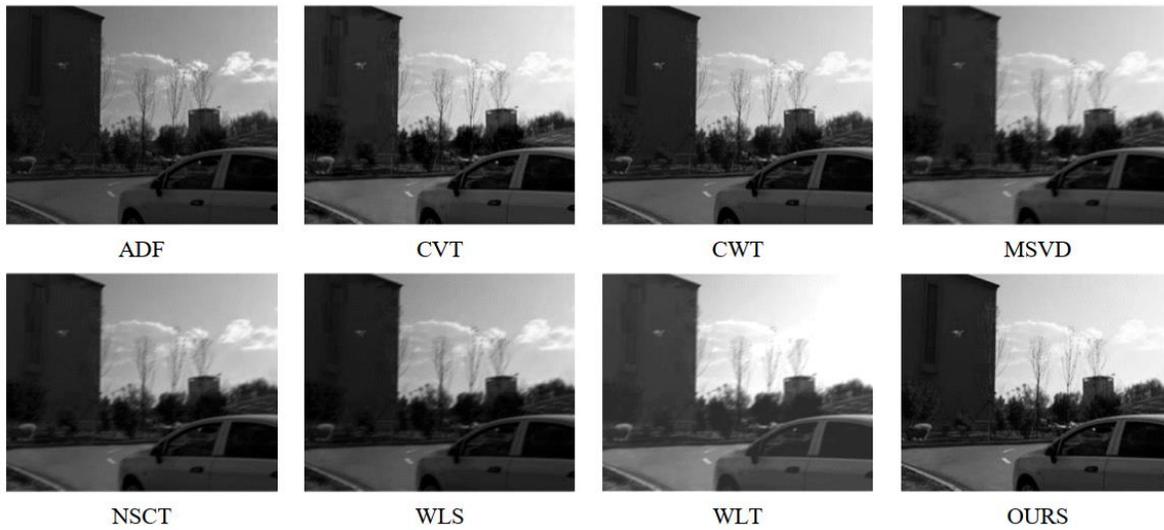


Fig. 7 Fusion comparison on the frames of the video V1

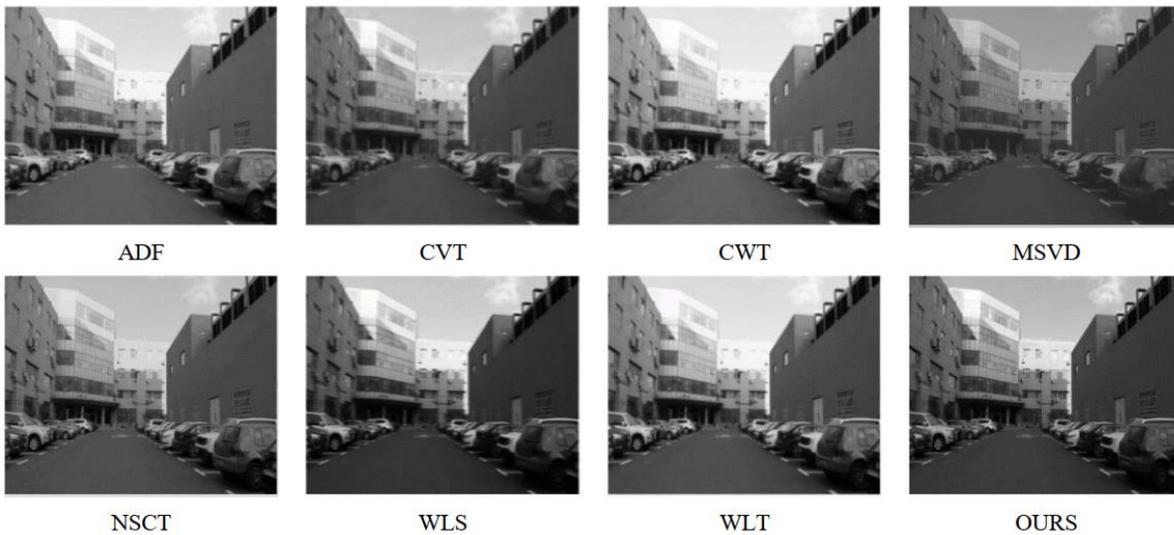


Fig. 8 Fusion comparison on the frames of the video V2

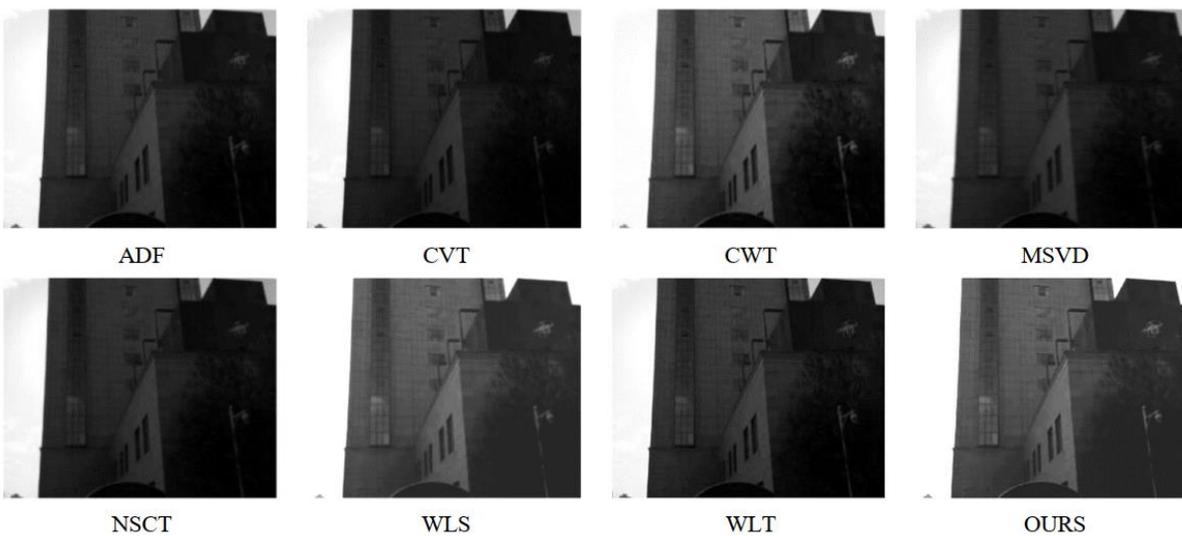


Fig. 9 Fusion comparison on the frames of the video V3

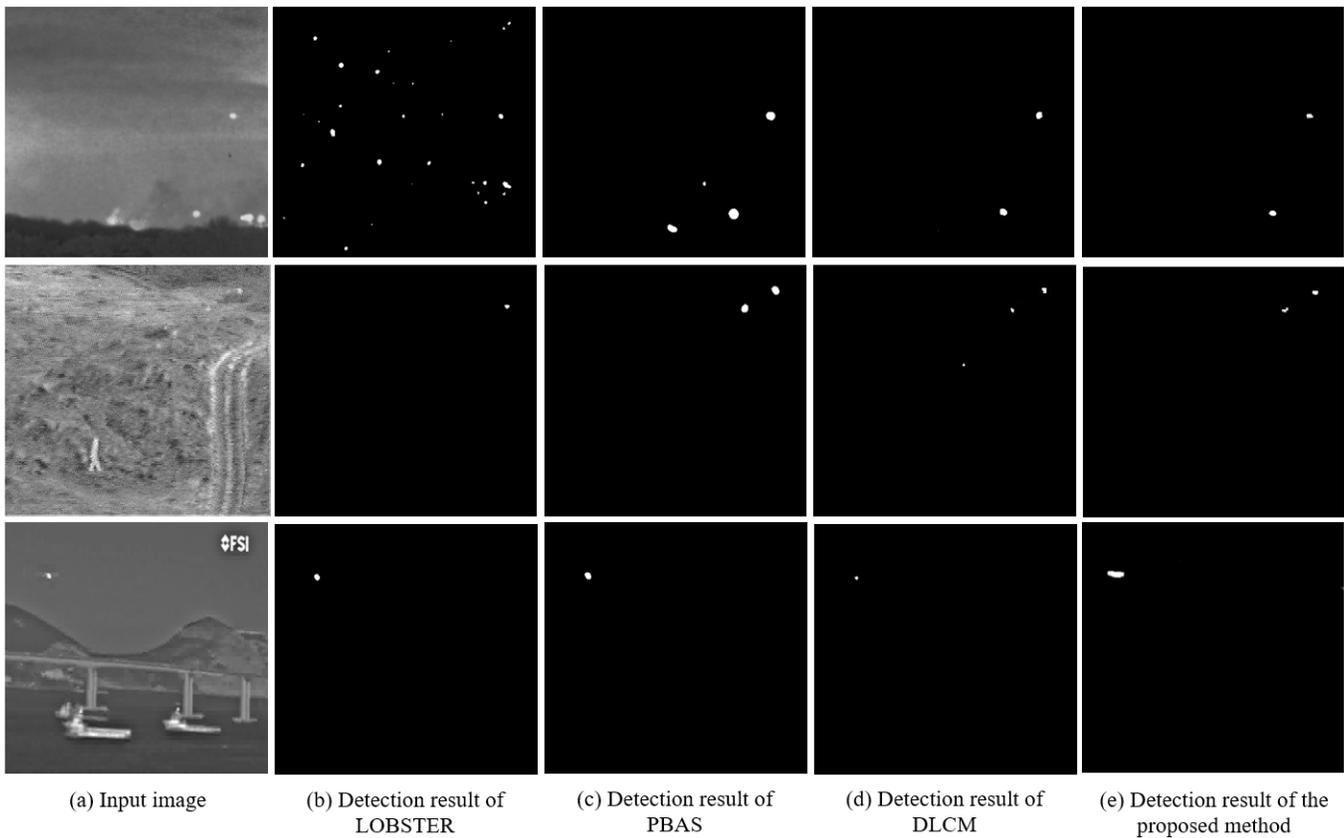


Fig. 10 Comparison of the final detection results. From top to bottom, the five rows are the input images and results from the LOBSTER, PBAS, DLCM and our proposed approach, respectively.

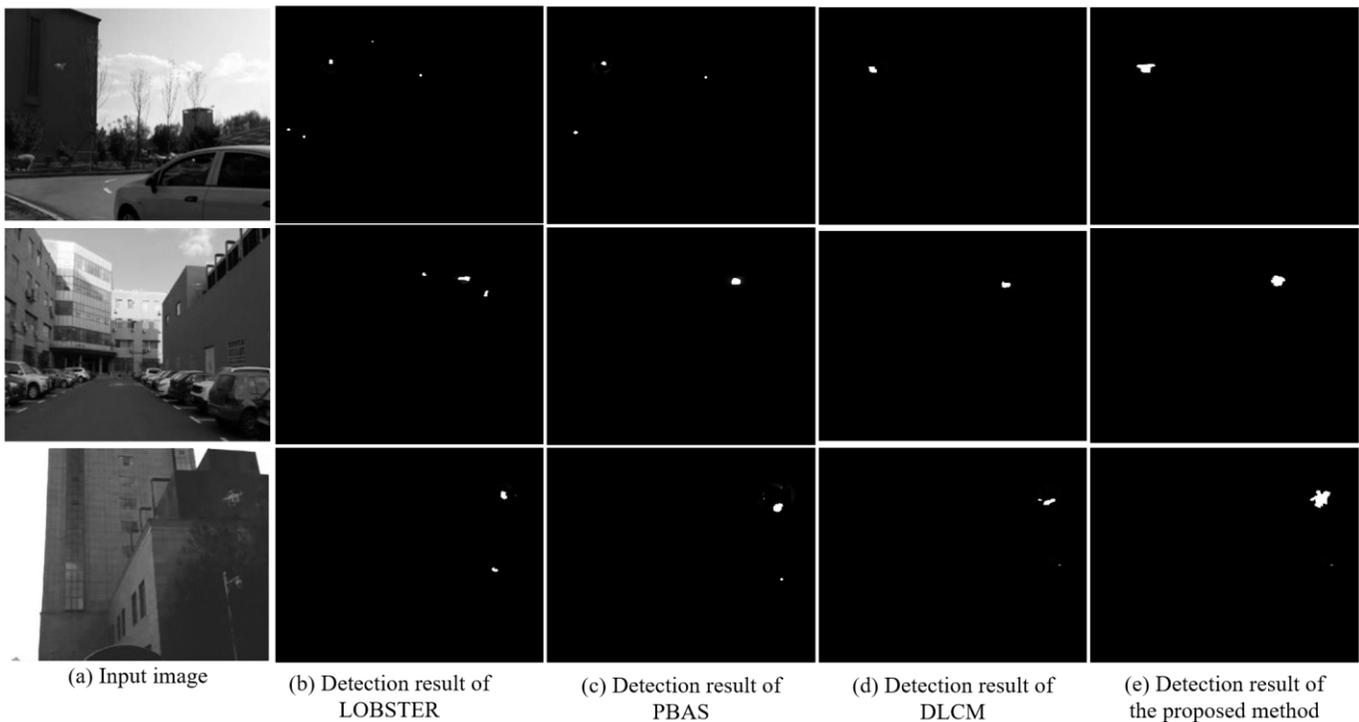


Fig. 11 Comparison of the final detection results. From top to bottom, the five rows are the input images and results from the LOBSTER, PBAS, DLCM and our proposed approach, respectively.

TABLE IV Fusion comparison of DHV with the best results in each group highlighted in bold.

Video2	AVG	IE	QE	SF	SD	MI	f/s
WLT	0.74	0.34	0.44	1.65	15.99	11.73	0.46
CWT	1.20	5.94	0.58	2.55	16.62	11.88	0.82
CVT	1.41	6.08	0.51	3.29	17.96	12.16	1.30
NSCT	1.21	5.96	0.59	2.56	16.78	11.92	2.40
ADF	0.94	5.88	0.55	1.94	16.16	11.76	0.58
MSVD	1.14	5.89	0.48	2.42	16.14	11.77	0.23
WLS	1.42	6.30	0.47	3.02	25.02	12.60	1.92
OUR	1.44	6.28	0.60	3.16	25.55	12.68	2.30

proposed. The accuracy is the proportion of the LSS-targets among all detected targets. The recall rate is the proportion of detected LSS targets in the global samples. The F value is the weighted harmonic mean of the accuracy and recall rate. As shown in Table IX, it can be seen from the results that the algorithm is still robust even when one of the sensors has poor imaging and achieves a higher detection accuracy than the traditional background-modeling algorithm.

TABLE V Fusion comparison of FEL with the best results in each group highlighted in bold.

Video3	AVG	IE	QE	SF	SD	MI	f/s
WLT	2.59	5.75	0.43	4.39	13.22	11.50	0.43
CWT	4.50	5.91	0.52	7.65	14.68	11.82	0.79
CVT	4.58	5.95	0.49	7.75	15.06	11.90	1.09
NSCT	4.53	5.95	0.54	7.70	15.05	11.89	1.35
ADF	4.18	5.82	0.57	7.14	13.87	11.65	0.71
MSVD	4.68	5.82	0.48	7.89	13.86	11.65	0.14
WLS	5.01	6.09	0.52	8.69	16.58	12.18	1.07
OUR	5.03	6.10	0.55	8.75	16.61	12.97	1.43

TABLE VI Fusion comparison of V1 with the best results in each group highlighted in bold.

Video3	AVG	IE	QE	SF	SD	MI	f/s
WLT	1.59	4.75	0.33	3.39	12.22	10.50	0.33
CWT	3.50	4.91	0.42	6.65	13.68	10.82	0.69
CVT	3.58	4.95	0.39	6.75	14.06	10.90	1.09
NSCT	3.53	4.95	0.44	6.70	14.05	10.89	1.25
ADF	3.18	4.82	0.47	6.14	12.87	10.65	0.61
MSVD	3.68	4.82	0.35	6.89	12.86	10.65	0.18
WLS	4.01	5.09	0.42	7.69	15.58	11.18	0.95
OUR	4.03	5.10	0.45	7.75	15.61	11.97	1.33

In addition, the detection speed of the proposed method is also compared with others in Table X, where the time needed for processing each frame is given, according to two publicly

available datasets and three our own datasets. Thanks to the background modelling of the candidate target region, which has greatly increased the detection speed. As can be seen, our method can achieve a high detection accuracy with a very competitive running time in comparison to SOTA approaches.

TABLE VII Fusion comparison of V2 with the best results in each group highlighted in bold.

Video3	AVG	IE	QE	SF	SD	MI	f/s
WLT	0.84	0.44	0.52	1.72	16.01	11.85	0.51
CWT	1.32	6.08	0.63	2.61	16.72	11.92	0.93
CVT	1.50	6.13	0.60	3.49	18.06	12.26	1.42
NSCT	1.36	6.12	0.59	2.63	16.88	12.05	2.33
ADF	1.25	5.95	0.61	2.02	16.26	11.86	0.69
MSVD	1.28	5.84	0.53	2.56	16.24	11.81	0.38
WLS	1.51	6.22	0.51	3.22	25.12	12.73	2.07
OUR	1.55	6.38	0.73	3.26	25.65	12.79	2.41

TABLE VIII Fusion comparison of V3 with the best results in each group highlighted in bold.

Video3	AVG	IE	QE	SF	SD	MI	f/s
WLT	0.94	0.54	0.44	1.95	15.89	11.62	0.36
CWT	2.12	5.89	0.49	2.40	16.52	11.79	0.72
CVT	1.39	6.08	0.53	2.49	17.86	12.05	1.20
NSCT	1.15	5.93	0.62	2.41	16.68	11.83	2.30
ADF	0.96	5.66	0.75	1.94	16.06	11.56	0.48
MSVD	2.04	6.79	0.48	2.35	16.04	11.68	0.13
WLS	1.72	6.27	0.73	2.53	24.92	12.40	1.82
OUR	2.34	7.18	1.50	2.92	26.45	13.51	2.20

TABLE X: Comparison of the time for LSS target detection in each frame (time in millisecond for each frame)

	Guanabara Bay	DHV	FEL	V1	V2	V3
LOBSTER	147	81	255	121	93	216
PBAS	89	58	166	73	62	149
DLCM	36	29	45	24	34	53
Ours	93	56	112	65	51	137

VII. CONCLUSIONS

In this paper, we have presented a novel algorithm for the detection of low-altitude, slow-speed small (LSS-) targets through the effective fusion of visible and infrared images. Firstly, an ROI extraction module based on 1-D information entropy and weighted average is proposed to reduce the overall calculation time and alleviate the interference from the unnecessary background information. Secondly, an efficient and effective multi-scale layered image fusion module is introduced, with a saliency map for enhancing the image details

TABLE IX: Detection accuracy (Da) and F-Measure (F-M) with the best results in each group highlighted in bold.

Test Video	Guanabara Bay		DHV		FEL		V1		V2		V3	
	Da	F-M	Da	F-M	Da	F-M	Da	F-M	Da	F-M	Da	F-M
LOBSTER	0.694	0.687	0.864	0.871	0.9535	0.953	0.9535	0.953	0.9535	0.953	0.9535	0.953
PBAS	0.847	0.855	0.843	0.834	0.966	0.968	0.966	0.968	0.966	0.968	0.966	0.968
DLCM	0.956	0.893	0.915	0.887	0.968	0.969	0.968	0.969	0.968	0.969	0.968	0.969
Ours	0.962	0.973	0.964	0.978	0.975							

of the fused image. Finally, accurate LSS-target detection is completed by local background modeling with the local SuBSENSE method. The experimental results also illustrate the efficacy of the two fusion modules, as our proposed approach outperforms the state-of-the-arts. More importantly, the proposed approach can be effectively used to detect LSS targets with a high accuracy and robustness even when a single sensor has poor images.

REFERENCES

- [1] S. A. Musa, et al., "Low-Slow-Small (LSS) target detection based on micro doppler analysis in forward scattering radar geometry," *Sensors*, 2019, vol.19, no. 15, pp.3332.
- [2] J. Li, Y. Dai, et al. "Visual detail augmented mapping for small aerial target detection," *Remote Sensing*, 2018, vol. 11, no. 1.
- [3] A. Khashman, "Neural system for military target detection," In *Proc. 8th IEEE Int. Conf. on Electronics, Circuits and Systems, Malta*, 2001, vol.3, pp. 1359-1362.
- [4] G Pedro, G Hilario, et al., "Continuous variance estimation in video surveillance sequences with high illumination changes", *Signal Processing*, 2009, vol. 89, no. 7, pp.1412-1416.
- [5] Q Shao, Z Tang, and S Han, "Hierarchical codebook for background subtraction in MRF", *Infrared Physics & Techno.*, 2013, vol. 61, no. 1, pp.259-264.
- [6] B. Wang, Y. Liu et al., "Background subtraction using spatiotemporal condition information," *Optik-Int. J. for Light and Electron Optics.*, 2014, vol. 125, no. 3, pp. 1406-1411.
- [7] Y. Lv, S. Sun, Q. Lin, and R. Liu, "Space moving target detection and tracking method in complex background," *Infrared Physics & Technology.*, 2018, vol. 91, no. 17, pp.107-118.
- [8] S. Chiang, C. Chang and I. W. Ginsberg, "Unsupervised target detection in hyperspectral images using projection pursuit," in *IEEE Trans. Geoscience and Remote Sensing*, 2001, vol. 39, no. 7, pp. 1380-1391.
- [9] J Lou, W Zhu, H Wang, et al. "Small target detection combining regional stability and saliency in a color image". *Multimedia Tools and Applications*, 2017, vol. 76, no. 13, pp. 14781-14798.
- [10] T. Xie, Z Chen, R Ma. "A novel infrared small target detection method based on PGF, BEMD and local inverse entropy". *Acta infrared and millimeter wave*, 2017, vol. 36, no. 1, pp. 92-101.
- [11] D P Bavirisetti and R Dhuli, "Two-scale image fusion of visible and infrared images using saliency detection," *Infrared Physics & Technology.*, 2016, vol. 76, no. 6, pp.52-64.
- [12] S. Qiu, J. Luo, et al, "A moving target extraction algorithm based on the fusion of infrared and visible images," *Infrared Physics & Technology*, 2019, vol. 98, no. 27, pp. 285-291.
- [13] O. Barnich and M. Van Droogenbroeck, "ViBE: A powerful random technique to estimate the background in video sequences," *IEEE Int. Conf. on Acoustics, Speech and Signal Processing*, Taipei, 2009, pp. 945-948.
- [14] M. Hofmann, P. Tiefenbacher, and G. Rigoll. "Background segmentation with feedback: The pixel-based adaptive segmenter". *Computer Vision and Pattern Recognition Workshops*, 2012.
- [15] B. P. Marcus. "The weighted moving average technique" *Wiley Encyclopedia of Operations Research and Management Science*. John Wiley & Sons, Inc. 2010.
- [16] J F Macgregor and T J Harris. "The exponentially weighted moving variance". *J. Quality Technology*, 1993, vol. 25, no. 2, pp.106-118.
- [17] P. St-Charles, G. Bilodeau and R. Bergevin, "SuBSENSE: a universal change detection method with local adaptive sensitivity," *IEEE Trans. Image Processing*, 2015, vol. 24, no. 1, pp. 359-373.
- [18] Toet, Alexander (2014): TNO Image Fusion Dataset. https://figshare.com/articles/dataset/TNO_Image_Fusion_Dataset/1008029
- [19] SMT/COPPE/Poli/UFRJ & IME-Instituto Militar de Engenharia (2016): Visible-Infrared data <http://www02.smt.ufrj.br/~fusion/>
- [20] S Liu, J Zhao, M Shi. "Medical image fusion based on rolling guidance filter and spiking cortical model". *Computational and Mathematical Methods in Medicine*, 2015, vol. 2015, pp.156043.
- [21] R. Achanta, S. Hemami, F. Estrada and S. Susstrunk, "Frequency-tuned salient region detection," *IEEE Conf. on Computer Vision and Pattern Recognition*, Miami, FL, 2009, pp. 1597-1604.
- [22] L. J. Chipman, T. M. Orr, and L. N. Graham, "Wavelets and image fusion," presented at the Proceedings of the International Conference on Image Processing., 1995, pp. 248-251
- [23] J. J. Lewis, R. J. O'Callaghan, et al, "Pixel-and region-based image fusion with complex wavelets," *Information Fusion*, 2007, vol. 8, no. 2, pp.119-130.
- [24] F. Nencini, A. Garzelli, and S. Baronti, L, "Remote sensing image fusion using the curvelet transform," *Information Fusion*, 2007, vol. 8, no. 2, pp. 143-156.
- [25] Y. Liu, S. Liu, and Z. Wang, "A general framework for image fusion based on multi-scale transform and sparse representation," *Information Fusion.*, 2015, vol. 24, no. 12, pp.147-164.
- [26] D. P. Bavirisetti and R. Dhuli, "Fusion of infrared and visible sensor images based on anisotropic diffusion and Karhunen-Loeve Transform," *IEEE Sensors Journal*, 2016, vol. 16, no. 1, pp. 203-209.
- [27] V. Naidu, "Image fusion technique using multi-resolution singular value decomposition," *Defence Science Journal.*, 2011, vol. 61, no.5, pp.479-484.
- [28] J Ma, C Chen, C Li, and J Huang, "Infrared and visible image fusion via gradient transfer and total variation minimization," *Information Fusion.*, 2016, vol. 31, no. 9, pp.100-109.
- [29] Z. Liu, Y. Song, et al, "MRI and PET image fusion using the nonparametric density model and the theory of variable-weight," *Computer Methods and Programs in Biomedicine.*, 2019, vol. 175, no. 11, pp.73-82.
- [30] X Bai, Y Zhang, F Zhou, and B Xue, "Quadtree-based multi-focus image fusion using a weighted focus-measure," *Information Fusion.*, 2015, vol. 22, no. 10, pp.105-118.
- [31] M Yin, P Duan, W Liu, and X Liang, "A novel infrared and visible image fusion algorithm based on shift-invariant dual-tree complex shearlet transform and sparse representation," *Neurocomputing.*, 2017, vol. 226, no. 19, pp.182-191.
- [32] X Kong, L Liu, Y Qian, and Y Wang, "Infrared and visible image fusion using structure-transferring fusion method," *Infrared Physics & Technology.*, 2019, vol. 98, no. 16, pp. 161-173.
- [33] M. Shahid Farid, A. Mahmood, and S. A. Al-Maadeed, "Multi-focus image fusion using content adaptive blurring," *Information Fusion.*, 2019, vol. 45, no. 9, pp. 96-112.
- [34] T Ma, J Ma, B Fang, et al, "Multi-scale decomposition based fusion of infrared and visible image via total variation and saliency analysis," *Infrared Physics and Technology.*, 2018, vol. 92, no. 23, pp.154-162.
- [35] G. Bilodeau, J. Jodoin and N. Saunier, "Change detection in feature space using local binary similarity patterns," *Int. Conf. on Computer and Robot Vision*, Regina, SK, 2013, pp. 106-112.
- [36] C Shu, D Baokuo, "Foreground detection of the adaptive LOBSTER algorithm in a dynamic background". *Journal of Image and Graphics*, 2017, no. 2, pp. 161-169.
- [37] M. Hofmann, P. Tiefenbacher and G. Rigoll, "Background segmentation with feedback: The pixel-based adaptive segmenter," *IEEE Conf. on Computer Vision and Pattern Recognition Workshops*, Providence, RI, 2012, pp. 38-43.

- [38] D Shen, S Zhang, M Zhao, "Infrared small target detection based on double-layer local contrast measure", *Acta Photonica Sinica*, 2020, vol. 49, no. 1, pp. 0110003.
- [39] Y Li, Y Wang. "Video image motion background subtraction based on one-dimensional maximum entropy", *Computer and modernization*, 2018. vol. 3, no. 4, pp. 44-53.
- [40] J Ma, Z Zhou, et.al. "Infrared and visible image fusion based on visual saliency map and weighted least square optimization" *Infrared Physics & Technology*, 2017. vol. 82, pp:8-17.
- [41] J. Zabalza et al., "Novel two-dimensional singular spectrum analysis for effective feature extraction and data classification in hyperspectral imaging," in *IEEE Trans. Geoscience and Remote Sensing*, 2015, vol. 53, no. 8, pp. 4418-4433.
- [42] Q. Liu, Y. Wang, et.al, "Decontaminate feature for tracking: adaptive tracking via evolutionary feature subset," *J. Electron. Image*. 2017, vol. 26, no. 6, pp.063025.
- [43] Y. Yan, et al, "Unsupervised image saliency detection with Gestalt-laws guided optimization and visual attention based refinement," *Pattern Recognition*, 2018, vol. 79, pp. 65-78.
- [44] Z. Wang, et al, "A deep-learning based feature hybrid framework for spatiotemporal saliency detection inside videos," *Neurocomputing*, 2018, vol. 287, pp. 768-83.
- [45] G. Sun, et al, "Deep fusion of localized spectral features and multi-scale spatial features for effective classification of hyperspectral images," *Int. J. of Applied Earth Observation and Geoinformation*, 2020, vol. 91, pp. 102157.
- [46] Q. Liu, J. Ren et al, "EACOFT: An energy-aware correlation filter for visual tracking," *Pattern Recognition*, 2021, vol. 112, pp. 107766.



and applications.

Jinchang REN received the PhD degree in Electronic Imaging and Media Communication from the University of Bradford, United Kingdom in 2009. He is currently a Professor of Computing Science with the Robert Gordon University, Aberdeen, U.K. His research interests include hyperspectral imaging, computer vision, machine learning, big data analytics



image enhancement and moving target detection.

Yanfeng WU was born in Jilin, China, in 1992. He received the bachelor's in applied physics in Jilin University, the M.S. degree in optical engineering from University of Chinese Academy of Sciences (UCAS), Beijing, China, in 2017. He is currently pursuing the PhD degree in mechanical & electronic engineering. His main research interests include



He is currently a professor and the Dean of the School of Computer Science, GPNU, Guangzhou, China. His research interests include image/video processing, and information security technology and applications.

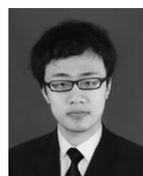
Huimin Zhao received the B.Sc. and the M.Sc. degrees in signal processing from Northwestern Polytechnical University, Xi'an, China, in 1992 and 1997, respectively, and the Ph.D. degree in electrical engineering from the Sun Yat-sen University, Guangzhou, China, in 2001.



Haijiang SUN received the PhD degree in electronic engineering from Changchun Institute of Optics, Fine Mechanics and Physics (CIOMP), Chinese Academy of Sciences in 2010. He is currently a Professor in the Perception and Display Lab., CIOMP. His research interests include target recognition & tracking, high definition video enhancement and display. He is an Associate Editor of the *Liquid Crystal and Reality* and *Optical Precision Engineering*.



Qiaoyuan LIU received the master and PhD degree in Department of Computer Science and Technology from the Northeast Normal University. Currently, she is an Assistant Professor of Changchun Institute of optics, precision mechanics and physics, Chinese Academy of Sciences. Her current research interests are visual tracking and image analysis.



Jiacheng WANG received his master degree from Xidian University, in Circuits and Systems in 2014. From 2014 to 2016, he was a Research intern in CIOMP. Since 2016, he has been a Research Assistant with the Image Processing Lab., CIOMP. His research interests include target tracking and embedded software & hardware design.

Huakang Li is a Research Assistant with the School of Computer Science, GPNU, Guangzhou, China. Her research interests include image/video processing and information security.