

HighlightNet: Highlighting Low-Light Potential Features for Real-Time UAV Tracking

Changhong Fu*, Haolin Dong, Junjie Ye, Guangze Zheng, Sihang Li, and Jilin Zhao

Abstract—Low-light environments have posed a formidable challenge for robust unmanned aerial vehicle (UAV) tracking even with state-of-the-art (SOTA) trackers since the potential image features are hard to extract under adverse light conditions. Besides, due to the low visibility, accurate online selection of the object also becomes extremely difficult for human monitors to initialize UAV tracking in ground control stations. To solve these problems, this work proposes a novel enhancer, *i.e.*, HighlightNet, to light up potential objects for both human operators and UAV trackers. By employing Transformer, HighlightNet can adjust enhancement parameters according to global features and is thus adaptive for the illumination variation. Pixel-level range mask is introduced to make HighlightNet more focused on the enhancement of the tracking object and regions without light sources. Furthermore, a soft truncation mechanism is built to prevent background noise from being mistaken for crucial features. Evaluations on image enhancement benchmarks demonstrate HighlightNet has advantages in facilitating human perception. Experiments on the public UAVDark135 benchmark show that HighlightNet is more suitable for UAV tracking tasks than other state-of-the-art (SOTA) low-light enhancers. In addition, real-world tests on a typical UAV platform verify HighlightNet’s practicability and efficiency in nighttime aerial tracking-related applications. The code and demo videos are available at <https://github.com/vision4robotics/HighlightNet>.

I. INTRODUCTION

Object tracking has been widely used in a variety of sectors, most representative in UAV applications for target following [1], autonomous landing [2], and self-localization [3]. SOTA tracking approaches [4], [5] have achieved outstanding performance under favorable illumination. However, the degradation of these trackers’ performance in low-light conditions has been discussed in recent studies [6], [7]. On the one hand, as shown in Fig. 1, human operators can hardly make an accurate initial annotation for trackers. On the other hand, due to the quick mobility of both the UAV and the tracked object, as well as numerous challenges such as occlusion, noise, illumination variation, robust and precise tracking in nighttime conditions has remained a difficult task. An effective approach to light up the target object for both human operators and trackers is in dire need.

A low-light enhancer can be a solution to this problem. However, enhancement of objects in dark areas may lead to overexposure problems in regions under favorable illumination, which may cause the destruction of potential features. Besides, artificial light sources are omnipresent in low-light conditions, which divided the image into regions with

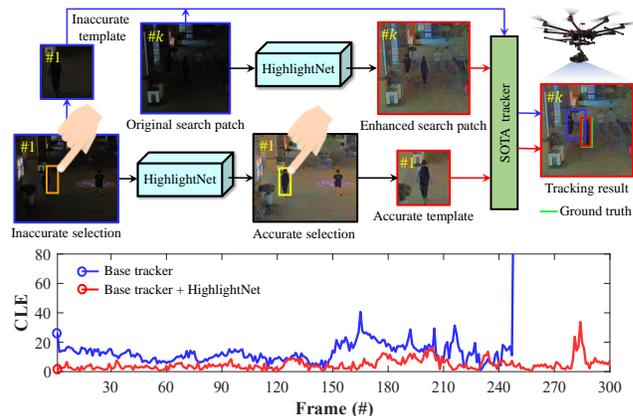


Fig. 1. Comparison of the overall tracking results from the baseline tracker with HighlightNet (in red) or not (in blue). Operators may make an inaccurate selection due to the low illumination. With the help of HighlightNet, human operators can conduct an accurate selection of targets, and the SOTA tracker can achieve promising results, as also shown in center location error (CLE) curves.

different illumination. In UAV object tracking, this problem is especially crucial because of UAVs’ wide vision field. When the target enters over-exposed or extreme dark regions, it is more difficult for the tracker to recognize the object. In this case, this work proposes a range mask to distinguish the enhancement of regions with or without artificial light sources. Moreover, experiments on typical nighttime UAV scenes show that the enhancer can focus on the enhancement of tracking objects according to the mask.

Only highlighting tracking objects and other effective local features is not enough. In nighttime UAV applications, it is inevitable that noises will damage the structural details of images and even be mistaken for tracking objects. Moreover, with excessive background information in UAVs’ wide vision field, more background noises are also introduced. However, existing low-light enhancers can hardly distinguish the noise from effective features. To avoid over-enhancement of noise, we consequently design a dark area truncation function, which produces an anti-noise mask to filter unwanted noise.

Apart from effective local features, global information is also essential for low-light enhancers. Due to the high mobility of UAVs, rapid changes of scenarios can hardly be avoided, which may cause severe global illumination variation. Nevertheless, most low-light enhancers [8]–[12] are designed for stable photography scenes. Without adjustment mechanisms, their effectiveness is unsatisfying while dealing with illumination variation. To address this problem, Transformer [13] is introduced to dynamically regulate parameters according to global illumination information.

*Corresponding author

The authors are with the School of Mechanical Engineering, Tongji University, Shanghai 201804, China. changhongfu@tongji.edu.cn

Another difficulty of UAV image enhancement is the limitation of computational resources. In consideration of the take-off weight, UAVs can hardly carry a high-power computing platform. Approaches with a high computation complexity have difficulty being processed in real-time on UAVs. To improve the computing efficiency, the main workflow of HighlightNet is based on the gray channel rather than the RGB channels. In addition, for supervised training, most low-light enhancers use paired data. In order to reduce the high expense of gathering enough paired data for UAV conditions, HighlightNet is designed to be trained with unpaired data.

The contributions of this work are summarized as follows:

- An adaptive enhancer HighlightNet is constructed to facilitate both online object selection and UAV tracking in low-light conditions.
- The range mask and the anti-noise mask are designed to highlight the object and filter unwanted noise.
- A dynamic parameter adjustment method based on an efficient Transformer is proposed to process global features for illumination retouching.
- Comprehensive evaluation on image enhancement test sets, UAV nighttime tracking benchmarks, and real-world experiments demonstrate that HighlightNet has advantages in facilitating human operator perception and improving trackers' nighttime performance.

II. RELATED WORKS

A. UAV Tracking

Object tracking methods include correlation filter-based approaches [14], [15] and convolutional neural network (CNN)-based approaches [5], [16], [17]. Among them, due to the promising tracking performance, methods based on Siamese networks have been widely used in UAV tracking. SiamAPN [18] is a Siamese-based no-prior two-stage method for adaptive anchor proposing. [4], [16], [19] employ the region proposal network (RPN) as anchor-based trackers to further improve tracking accuracy. HiFT [5] introduces a feature Transformer network to obtain hierarchical similarity maps. Ad²Attack [20] proposes a novel adaptive adversarial attack approach to help increase awareness of the potential risk. TCTrack [21] designs a comprehensive framework to fully utilize temporal contexts for UAV tracking. In spite of their high accuracy in daytime conditions, their performance is unsatisfying in low-light conditions. Damage to potential features in nighttime conditions makes their robustness drop greatly. What's worse, precise online target selection is nearly impossible in extreme low-light conditions. Therefore, image enhancement technology before tracking is urgently needed.

B. Low-light Image Enhancement

The theory of Retinex [22] has evolved into a variety of low-light enhancement approaches. LIME [8] estimates a coarse illumination map by extracting the maximum of each pixel in RGB channels, which is subsequently refined by a structure prior. Based on Retinex, convolutional neural

networks (CNNs) are introduced and remarkably improve low-light performance [9], [10]. The first CNN model LL-Net [23] employs an autoencoder to learn denoising and light enhancement simultaneously. DeepUPE [24] introduces intermediate illumination to connect the input and the anticipated enhancement result. The method in [25] takes noise into consideration and realizes satisfying performance in enhancing practical noisy low-light images. However, because of the internal locality of convolution operations, CNNs are not suitable for processing global features.

Another disadvantage of most CNN-based methods is the requirement of paired data for supervised training. In consideration of the high investment caused by collecting sufficient paired data for UAV conditions, low-light enhancement by paired data is impractical. Therefore, advanced methods turn to unsupervised training. EnlightenGAN [11] employs a generative adversarial network (GAN) so it can be trained using adversarial loss as an attention-based U-Net. Zero-DCE [12] employs image-specific curve estimation and has high efficiency, also trained with unpaired data. RUAS [26] can also be trained with unpaired data by employing a cooperative reference-free learning strategy. Nevertheless, in low-light conditions, the performance of SOTA enhancers in target object selection and target tracking can hardly meet our satisfaction. This can be partly attributed to the neglect of potential effective features.

C. Low-light UAV Tracking

Though UAV tracking tasks in nighttime scenes are crucial, only several approaches [6], [7], [27], [28] to improve UAV tracking performance in low-light conditions have been proposed. ADTrack [6] serves for correlation filter (CF) trackers, while Darklighter [7] and SCT [27] are mainly designed for trackers based on Siamese network. Darklighter [7] constructs a lightweight map estimation network to cope with poor illumination and noise. Applying task-tailored training and a Transformer structure, SCT [27] realizes stable nighttime tracking. The basic idea of these three methods is to employ a plug-and-play enhancer as a preprocessing step for UAV trackers. UDAT [28] instead considers the nighttime tracking problem as an unsupervised domain adaptation task. Despite the progress, existing approaches generally ignore the improvement of human perception, which plays a vital role in the step of online object selection. Moreover, without adjustment mechanisms and denoising modules, these methods can hardly deal with illumination variation and background noises.

III. METHODOLOGY

This work proposes a Transformer-based adaptive enhancer called HighlightNet to highlight potential features for low-light conditions. As shown in Fig. 2, to serve both human operators and UAV trackers, the template patch and search patch are processed separately on ground control stations and UAVs. The proposed HighlightNet includes three novel modules. The Transformer-based parameter adjustment module obtains the constraint α and the truncation threshold

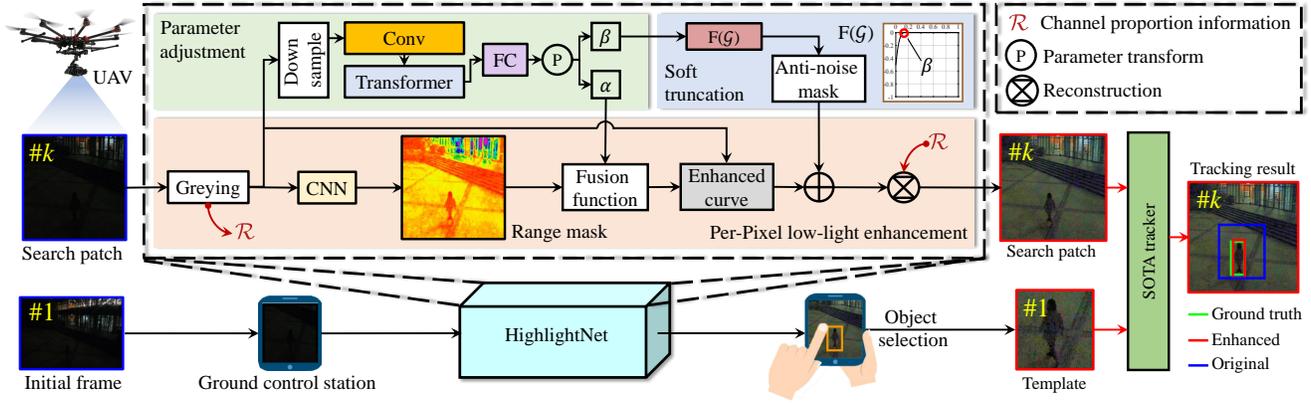


Fig. 2. Overview of our HighlightNet pipeline. Given a low-light image patch, HighlightNet includes three novel modules, *i.e.*, per-pixel low-light enhancement, Transform-based parameter adjustment, and soft truncation, to highlight potential features for both human operators and trackers with slight computational consumption. The bottom row serves for template selection by humans.

β by processing the downsampled gray-scale image. In the enhancement module, a range mask is constructed by a convolutional neural network (CNN). The mask is fused with constraint α by a fusion function to adapt the enhanced curve. As the mask corresponds to the original image pixel by pixel, each pixel has a unique enhancement range. The truncation threshold β acts on the soft truncation module. By the truncation function, the gray-scale image is translated to the anti-noise image. To improve computational efficiency, the main enhancement workflow is based on the gray-scale image. Therefore, the color information is first translated to channel proportion information and stored.

A. Per-Pixel Low-Light Enhancement

To focus on the enhancement of objects, this module is mainly designed for predicting the accurate enhancement for each pixel. Therefore, the range mask $\mathbf{M}_i \in \mathbb{R}^{1 \times H_i \times W_i}$ is with the same resolution as the input gray-scale image $\mathbf{G}_i \in \mathbb{R}^{1 \times H_i \times W_i}$. We employ a CNN with symmetrical concatenation to obtain the range mask \mathbf{M}_i . It is constructed by seven convolutional layers. In each layer, the ReLU activation function is followed by 4 convolutional kernels of size 3×3 . The sigmoid activation function is then applied to the final layer and produces the range mask \mathbf{M}_i . The range mask \mathbf{M}_i is then fused with constraint α by the fusion function to fit the enhanced curve.

As a pixel-level operation, the enhanced curve has a great influence on the computational efficiency of HighlightNet. Considering the scarce computational resources of UAV platforms, gamma transform is employed as the enhanced curve. It can realize large-scale enhancement without iteration. A standard gamma transform can be expressed as:

$$\mathbf{O}_i(r,c) = \mathbf{G}_i(r,c)^{\gamma_i(r,c)}, 0 \leq r < H_i, 0 \leq c < W_i, \quad (1)$$

where $\mathbf{G}_i(r,c)$ is the value of each pixel in gray-scale image \mathbf{G}_i , $\gamma(r,c)$ is the output of the fusion function, $\mathbf{O}_i(r,c)$ is the output of this enhancement module, H_i and W_i are the height and width of the image. Each pixel of the gray-scale image \mathbf{G}_i is normalized to $[0,1]$ before calculation. After

designating the enhanced curve, the fusion function can be designed.

The fusion function can change the codomain of each pixel's value in the range mask \mathbf{M}_i . Because of the sigmoid activation function, the value of pixels in the range mask \mathbf{M}_i is limited between 0 and 1. However, for the enhanced curve defined in Eq. (1), the input $\gamma(r,c)$ should be between the constraint α and 1. Therefore, the fusion function can be expressed as:

$$\gamma_i(r,c) = \alpha_i^{\mathbf{M}_i(r,c)}, 0 \leq r < H_i, 0 \leq c < W_i, \quad (2)$$

where $\mathbf{M}_i(r,c)$ are the value of each pixel in the range mask \mathbf{M}_i , α is the constraint obtained by the Transformer-based branch, $\gamma(r,c)$ is the input of the enhanced curve Eq. (1). The fusion function and the enhanced curve are not fixed, they can be changed to adapt to different tasks. The reason why we choose gamma transform is mainly for computational efficiency.

Three examples of the range mask are presented in Fig. 3. As shown, the range mask makes the enhancer focus on the enhancement of the potential objects. Moreover, the enhanced result exposes the features in shadow and maintains the areas with artificial light sources, which improves the

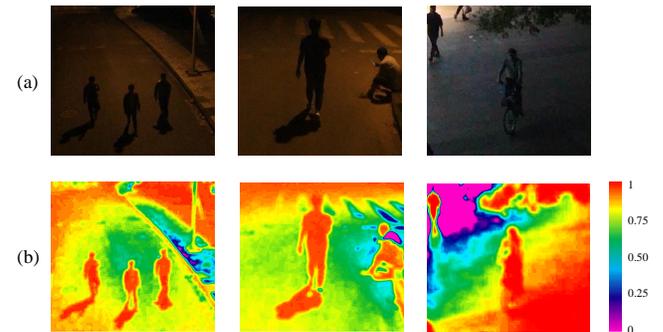


Fig. 3. Three examples of the pixel-wise range mask. (a) denotes original image. (b) denotes range mask, represented by heatmaps. For visualization, we select a particular region of the image and normalize the values to the range of $[0, 1]$. The range mask conducts HighlightNet to focus on the enhancement of tracking objects and regions without light resources.

brightness uniformity. Since the shadow is illuminated by HighlightNet, the tracker won't lose the object when it enters the dark areas without light resources.

Remark 1: The three masks in Fig. 3 are small regions cut from the high-resolution original mask. Therefore, the objects in these small masks can be considered as tiny targets. As shown in Fig. 3, with the help of the range mask, HighlightNet is able to adjust the enhancement range at pixel level and realize the targeted enhancement of small objects. Moreover, as shown in Fig. 4, HighlightNet can help the base tracker recognize small objects in low-light conditions.

B. Transformer-Based Parameter Adjustment

The main function of this module is producing constraint α and truncation threshold β by processing global features. As mentioned in the last section, constraint α is fused with the enhanced range mask to obtain the actual enhancement of each pixel. Therefore, this branch has a similar influence on all pixels. Since global information plays a more important role than local context, Transformer [13] is introduced. DETR [29] has shown the ideal efficiency of Transformer in processing global information in computer vision tasks. With the advantage of global processing, this module can dynamically regulate parameters to adapt to illumination variation caused by the high mobility of UAVs.

In this global processing module, the resolution of input does not play an important role. We find that the computation resource can be greatly saved without performance degradation by setting the resolution of input to 32×32 . Therefore, as shown in Fig. 2, the input gray-scale image $\mathbf{G}_i \in \mathbb{R}^{1 \times H_i \times W_i}$ is first identically downsampled to a low-resolution image $\mathbf{L}_i \in \mathbb{R}^{1 \times 32 \times 32}$. After processing by a convolutional layer, it is in the size of $16 \times 16 \times 16$ and then reshaped to $16 \times (16 \times 16)$ to enter the encoder of Transformer. Several identical layers make up the encoder and then two sublayers make up each layer. Each layer consists of two sub-layers. The first layer uses the concept of Multi-Head Attention [30] (MHA) to drive the model to concentrate on information from various positions and representation subspaces. MHA is defined as:

$$\text{MHA}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Cat}(\text{head}_1, \text{head}_2, \dots, \text{head}_h) \mathbf{W}^{\text{O}}, \quad (3)$$

where $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{n \times d_m}$ are input embedding matrices, n is sequence length, d_m is the embedding dimension, and h is the number of heads. Each head is defined as:

$$\text{head}_j = \text{softmax} \left[\frac{\hat{\mathbf{Q}} \hat{\mathbf{K}}^{\text{T}}}{\sqrt{d_{\mathbf{K}}}} \right] \hat{\mathbf{V}}, \quad (4)$$

where $d_{\mathbf{K}}, d_{\mathbf{V}}$ are the hidden dimensions of the projection subspaces. $\hat{\mathbf{Q}}, \hat{\mathbf{K}}, \hat{\mathbf{V}}$ are query, key, and value projected into the subspaces respectively, which can be defined as:

$$\hat{\mathbf{Q}} = \mathbf{Q} \mathbf{W}_j^{\text{Q}}, \hat{\mathbf{K}} = \mathbf{K} \mathbf{W}_j^{\text{K}}, \hat{\mathbf{V}} = \mathbf{V} \mathbf{W}_j^{\text{V}}, \quad (5)$$

where $\mathbf{W}_j^{\text{Q}}, \mathbf{W}_j^{\text{QK}} \in \mathbb{R}^{d_m \times d_{\mathbf{K}}}, \mathbf{V} \mathbf{W}_j^{\text{V}} \in \mathbb{R}^{d_m \times d_{\mathbf{V}}}$ are learned matrices.

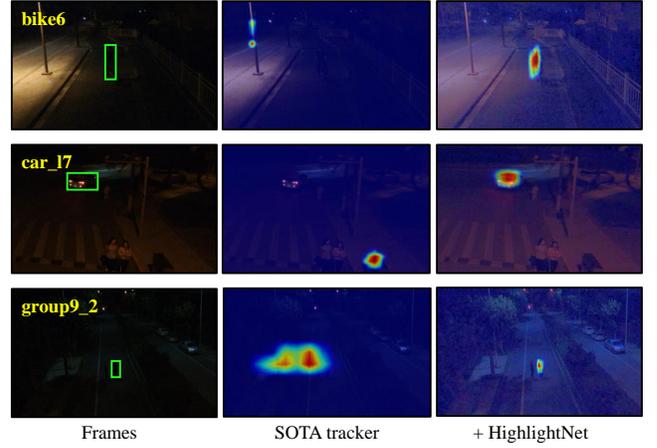


Fig. 4. Comparison of confidence maps generated by the SOTA tracker with and without HighlightNet. The tracked objects are shown by the green boxes in the first column. The sequence names are provided in the top left corner of the original frames from the public UAVDark135 benchmark. Without HighlightNet, the base tracker loses its capacity to recognize objects in low-light conditions.

The second sub-layer is a fully connected feed-forward network (FFN). The connection between each sub-layer is a residual module followed by layer normalization. Therefore, the output of each sub-layer can be expressed as:

$$\begin{aligned} \mathbf{F}' &= \text{LN}(\text{MHA}(\mathbf{F}) + \mathbf{F}) \quad , \\ \mathbf{F}'' &= \text{LN}(\text{FFN}(\mathbf{F}') + \mathbf{F}') \quad , \end{aligned} \quad (6)$$

where \mathbf{F}' and \mathbf{F}'' denote the output of MHA and FFN, respectively. LN is layer normalization. The output of Transformer is in the same size as input $16 \times (16 \times 16)$. In order to obtain two adaptive parameters, a simple fully connected layer(FC) followed by the sigmoid activation function is introduced to reshape the output to a bivector $[\alpha, \beta]$. Each dimensionality in this bivector is then translated to constraint α and truncation threshold β by the parameter transform module, which is a simple linear conversion.

C. Soft Truncation

In this module, a soft truncation function $\mathbf{T}(\mathbf{G})$ is proposed to filter useless features brought by over-enhanced noise. The intensity of enhancement weakening is negatively correlated with the gray value. Moreover, considering the limitation of computational resources on the UAV platform, this function should be easy to compute. Therefore, a cubic function determined by two parameters is employed, which can be expressed as:

$$\mathbf{T}_i(r, c) = -\tau \times (\beta_i - \mathbf{G}_i(r, c))^3, \tau \beta_i^3 \mathbf{T}_i(r, c) < 0, \quad (7)$$

$$0 \leq r < H_i, 0 \leq c < W_i,$$

where $\mathbf{G}_i(r, c)$ is the value of each pixel in gray-scale image $\mathbf{G}_i \in \mathbb{R}^{1 \times H_i \times W_i}$, β_i is the threshold of truncation, τ determines the reducing range. $\mathbf{T}_i(r, c)$ is the output of the truncation function with the codomain of $[\tau \beta_i^3, 0)$. Therefore, the value of the output anti-noise mask is subtractive, which will apply to the pixel-wise addition module to reduce the enhancement range of dark area noise.

Remark 2: The threshold β_i of truncation is an adjustable parameter produced by the parameter adjustment module. Therefore, the result of soft truncation is not only based on the input grey-scale image but also influenced by the global features acting on Transformer. Threshold β_i increases in extreme low-light conditions to filter more useless features brought by noise.

For the same purpose, a non-reference loss function of dark area noise \mathcal{L}_{dan} is introduced. Since severe noise problems always appear in dark regions of the image, this loss function limits the enhancement range of shadows and night sky. The enhanced image is first divided into local regions in the size of 16×16 . The average enhancement range of pixels in the dark area of each region is then calculated. Therefore, \mathcal{L}_{dan} can be expressed as:

$$\mathcal{L}_{\text{dan}} = \frac{1}{N} \sum_i^N \mathbf{A}_i, \quad (8)$$

where N is the number of regions, \mathbf{A}_i is the average enhancement range of pixels in the dark area of each region. The threshold to determine whether the pixel is in a dark area is set to 0.04. This threshold is half of the average of β_i . Because the impact of soft truncation becomes obvious from this threshold. Except for the \mathcal{L}_{dan} , other three non-reference losses are employed to further improve the performance of HighlightNet. The total loss can be expressed as:

$$\mathcal{L}_{\text{total}} = \lambda_1 \mathcal{L}_{\text{dan}} + \mathcal{L}_{\text{spsa}} + \lambda_2 \mathcal{L}_{\text{exp}} + \lambda_3 \mathcal{L}_{\text{tv}}, \quad (9)$$

where \mathcal{L}_{dan} is dark area noise loss, $\mathcal{L}_{\text{spsa}}$ is spatial consistency loss, \mathcal{L}_{exp} is exposure control loss, \mathcal{L}_{tv} is the illumination smoothness loss, λ_1 , λ_2 , and λ_3 are the weights of losses. $\mathcal{L}_{\text{spsa}}$, \mathcal{L}_{exp} , and \mathcal{L}_{tv} are first proposed in [12] as non-reference losses. Therefore, HighlightNet can be trained with unpaired data, which avoid the high cost of preparing a paired dataset for UAV tracking tasks.

IV. EXPERIMENT

To testify the effectiveness and robustness of HighlightNet for both online object selection and UAV object tracking, our experiments can be divided into two parts. For the template selection task, we quantitatively testify the performance of HighlightNet on the Part2 subset of SICE dataset [31] and compare HighlightNet with SOTA low-light enhancers, respectively Darklighter [7], EnlightenGAN [11], Zero-DCE [12], and LIME [8]. As a pre-processing step for UAV target tracking, its success rate and precision on UAVDark135 benchmark with baseline are compared with other SOTA low-light enhancers to testify its advantage on UAV nighttime tracking tasks. Performance on sequences with different attributes in UAVDark135 is tested to verify its ability to deal with specific UAV challenges. Moreover, to demonstrate its universality for different UAV trackers, it has been implemented on 4 SOTA trackers, *i.e.*, SiamAPN++ [32], SiamAPN [18], HiFT [5], and SiamRPN++ [16]. To show the function of each module, an ablation study is also introduced. Finally, we conduct the

real-world test by adopting HighlightNet on a typical UAV platform to verify its applicability.

A. Implementation Details

To bring the capability of the illumination-based adaptive parameter adjustment into full play, the unpaired training set includes both low-light and over-exposed images. Therefore, images from Part1 of the SCIE dataset [31] are employed to train our network. The training images are resized to 512×512 . The weights λ_1 , λ_2 , and λ_3 are set to 200, 50, and 20 respectively, to balance the scale of losses. We implement our framework with PyTorch on a PC with an Intel i9-9920X CPU, an NVIDIA TITAN RTX GPU, and 32GB RAM. The batch size is set to 8, with a total of 100 epochs. We employ the ADAM optimizer and set the learning rate to 0.001. Other parameters of the optimizer are default. Finally, to confirm the viability of HighlightNet in nighttime UAV tracking, the real-world test uses an NVIDIA Jetson AGX Xavier, which is widely applied on UAV platforms.

B. Evaluation Metrics

To verify HighlightNet's advantage on the nighttime online object selection task, we perform quantitative experiments on the Part2 subset of SCIE dataset [31], which includes 229 multi-exposure sequences. To testify its performance in low-light conditions, we choose the first low-light image in all 229 sequences and resize them to a size of $960 \times 640 \times 3$. Finally, we obtain 229 paired low/normal light images. We choose the Part2 subset of the SCIE dataset for it contains numerous outdoor sequences and is primarily designed for low-light enhancement.

We also use visual tracking evaluation measures to rate the performance because HighlightNet is also targeting nighttime UAV tracking. The studies follow the one-pass evaluation method [33], which uses two metrics: precision and success rate. The precision is calculated by measuring the CLE between the estimated position and the ground truth position. And the success rate is calculated using the intersection over union (IoU) between the estimated bounding box and the ground truth.

C. Efficacy of HighlightNet for Online Target Selection

Since target selection is mainly executed by human operators, human perception should be evaluated quantitatively to testify the efficacy of HighlightNet. Therefore, the peak signal-to-noise ratio (PSNR), and structural similarity (SSIM) are employed. TABLE I reports the comparison result of HighlightNet with other SOTA low-light enhancers. HighlightNet achieves the highest PSNR and SSIM score on the Part2 subset of the SICE dataset [31]. Darklighter is specially designed for trackers and thus its performance is unsatisfying in the test for human perception. Since initial target selection is inevitable, HighlightNet is more practical in nighttime UAV applications.

Remark 3: HighlightNet not only has a better performance on the object selection task but also be system-friendly in terms of computational complexity, for the design of HighlightNet takes the computational cost into account heavily.

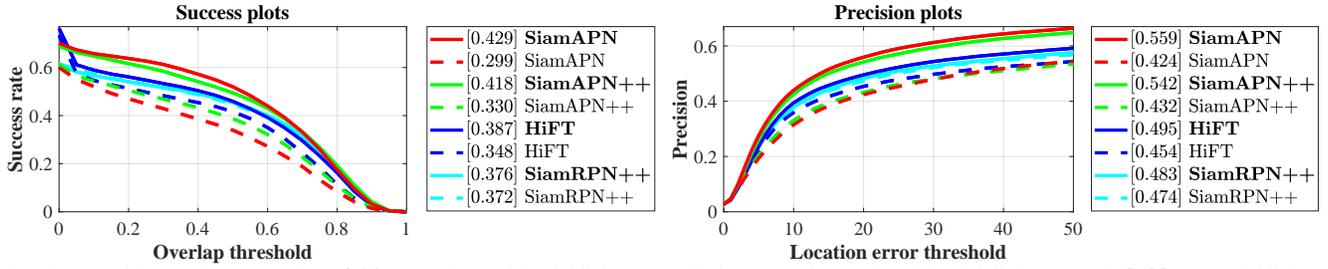


Fig. 5. Precision and success plots of SOTA trackers with HighlightNet enabled or not. The results with HighlightNet are in **bold** type. HighlightNet improves the performance of all involved trackers.

TABLE I

COMPARISON OF HIGHLIGHTNET WITH OTHER SOTA ENHANCERS IN TERMS OF FULL-REFERENCE IMAGE QUALITY ASSESSMENT METRICS. THE BEST RESULT IS IN **RED** WHEREAS THE SECOND BEST ONE IS IN **GREEN**. HIGHLIGHTNET HAS AN OBVIOUS ADVANTAGE ON FACILITATE HUMAN PERCEPTION.

Method	Darklighter [7]	LIME [8]	Zero-DCE [12]	EnlightenGAN [11]	HighlightNet (ours)
PSNR \uparrow	8.3	11.4	8.9	10.8	11.8
SSIM \uparrow	0.57	0.70	0.63	0.72	0.75

D. Efficacy of HighlightNet for UAV Dark Tracking

1) Comparison of HighlightNet and SOTA enhancers:

To demonstrate the superiority of HighlightNet for low-light UAV tracking, the performance of HighlightNet and other SOTA low-light enhancers—Darklighter [7], EnlightenGAN [11], Zero-DCE [12], and LIME [8], on UAVDark135 benchmark with tracker SiamAPN++ [32] is analyzed. SiamAPN++ is employed because it’s designed specifically for UAV tracking. TABLE II shows the tracking results with images enhanced by various enhancers. Since HighlightNet is mainly designed for UAV nighttime applications, it achieves the highest improvement in success rate and precision among the four enhancers. HighlightNet brings an increase of more than **26%** and **25%** in success rate and precision, surpassing the second-best Zero-DCE by **6.07%** and **8.8%** in the increase of success rate and precision, respectively.

2) HighlightNet on different UAV challenges:

Sequences in benchmark UAVDark135 with different attributes are employed to testify the tracking performance of specific challenges in UAV applications. The improvement of tracking performance in typical low-light aerial challenges, *i.e.*,

TABLE II

COMPARISON OF HIGHLIGHTNET WITH SOTA ENHANCERS ON UAVDARK135. THE FIRST, SECOND, AND THIRD BEST RESULTS ARE IN **RED**, **GREEN**, AND **BLUE** FONT, RESPECTIVELY. Δ_s AND Δ_p REPRESENT THE PERCENTAGES OF SUCCESS RATE AND PRECISION EXCEEDING THE BASELINE, *i.e.*, ORIGINAL TRACKER.

Method	Succ.	Prec.	Δ_s (%)	Δ_p (%)
Origin	0.336	0.428	-	-
+LIME [8]	0.415	0.518	23.5	21.0
+Darklighter [7]	0.415	0.529	23.5	23.8
+EnlightenGAN [11]	0.418	0.522	24.4	22.0
+Zero-DCE [12]	0.419	0.530	24.7	23.8
+HighlightNet (ours)	0.424	0.539	26.2	25.9

fast motion (FM), illumination variation (IV), low-resolution (LR), visual occlusion (OCC), and viewpoint change (VC) is reported in TABLE III. As shown, in the term of IV and FM, it achieves an enormous improvement of over **25%**, which is credited to the design of dynamic parameter adjustment according to global illumination. Furthermore, thanks to the pixel-level illumination enhancement brought by HighlightNet, the ability of trackers to cope with LR is also revived markedly. Though HighlightNet is not specially designed to deal with OCC and VC, it still brings an increase of over **13%**, which is due in large part to the capacity of HighlightNet to highlight potential features for the tracker.

3) HighlightNet on different trackers:

As shown in Fig. 5, HighlightNet facilitates the low-light tracking performance significantly, with obvious improvement for all trackers in both precision and success rate. Among these four trackers, HighlightNet improves the dark tracking performance of SiamAPN and SiamAPN++ obviously, with an increase of more than **42.6%** and **26.2%** in success rate, as well as an increase of **31.8%** and **25.9%** in precision.

Some tracking screenshots of the trackers with or without HighlightNet are exhibited in Figure 6. As shown, with the help of HighlightNet, trackers are able to retrieve the target missing in the shadow. We can conclude that HighlightNet boosts the accuracy and reliability of the trackers in these low-light conditions.

Remark 4: For visualization, images in Fig. 6 are enhanced by HighlightNet except for the first column. The enhanced results demonstrate that HighlightNet can markedly facilitate the perception of human operators in nighttime conditions, benefiting the initial annotation task.

E. Ablation Study

The performance of different variants of HighlightNet is investigated in this subsection. Tracking results on

TABLE III

EVALUATION OF TRACKING PERFORMANCE IN DIFFERENT UAV CHALLENGES, RESULTS ARE SHOWN AS SUCCESS/PRECISION. Δ DENOTES THE PERCENTAGES OF IMPROVEMENT BROUGHT BY HIGHLIGHTNET. IN ALL TYPICAL CHALLENGES OF DARK TRACKING, HIGHLIGHTNET BOOSTS TRACKING PERFORMANCE OBVIOUSLY.

Attributes	FM	IV	LR	OCC	VC
Original	0.315/0.404	0.312/0.408	0.323/0.464	0.332/0.416	0.353/0.426
+HighlightNet (ours)	0.399/0.507	0.398/0.515	0.402/0.571	0.392/0.497	0.408/0.485
Δ (%)	26.7/25.5	27.6/26.2	24.4/23.1	18.1/19.5	15.6/13.8



Fig. 6. Qualitative results of trackers with HighlightNet enabled (solid boxes) or not (dashed boxes). For visualization, HighlightNet is employed to enhance the images in the last three columns. The sequences, from top to bottom, are *bike3*, *bike9*, and *group1* from the public benchmark UAVDark135. HighlightNet dramatically improves tracker nighttime tracking performance.

UAVDark135 [6] are exhibited in TABLE IV. RM, TPA, and ST denote the range mask, the Transformer-based parameter adjustment, and the soft truncation, respectively. We trained all these variations in the same way. The bottom row shows the unsatisfying performance without HighlightNet, while the top row shows the completely enhanced result. Since the soft truncation and \mathcal{L}_{dan} both target filtering of useless features, their effects are discussed in the same sector.

1) *Range mask*: To ablate the range mask without impacting other modules, it is replaced with another mask. Each pixel in this mask is the same as the average value of pixels in the range mask. As shown in the sixth line of TABLE IV, without the range mask, the improvement of the tracking performance is reduced to **8.0%** in success rate and **8.5%** in precision. Moreover, with a decrease of **5.0%** and **6.0%**, PSNR and SSIM also degrade obviously. The function of the range mask to make HighlightNet focus on the tracking target and improve illumination uniformity can be verified.

2) *TPA*: Constraint α and threshold β are set to constant values, which are their average values in a complete HighlightNet. Since the parameters are invariable, the adjustment branch is consequently disabled. As shown in the fifth

TABLE IV

COMPARISON OF HIGHLIGHTNET WITH DIFFERENT MODULES ENABLED. SUCCESS RATE AND PRECISION (SUCC./PREC.) REPRESENT TRACKING PERFORMANCE WHILE PSNR AND SSIM REPRESENT HUMAN PERCEPTION. THE BEST RESULTS ARE IN RED FONT.

RM	TPA	ST	\mathcal{L}_{dan}	Succ./Prec.	PSNR/SSIM
				0.336/0.428	-/-
✓	✓			0.391/0.501	11.6/0.67
✓	✓	✓		0.419/0.530	11.3/0.73
✓	✓		✓	0.416/0.527	11.6/0.69
✓		✓	✓	0.419/0.526	10.6/0.73
	✓	✓	✓	0.418/0.524	11.2/0.71
✓	✓	✓	✓	0.424/0.539	11.8/0.75

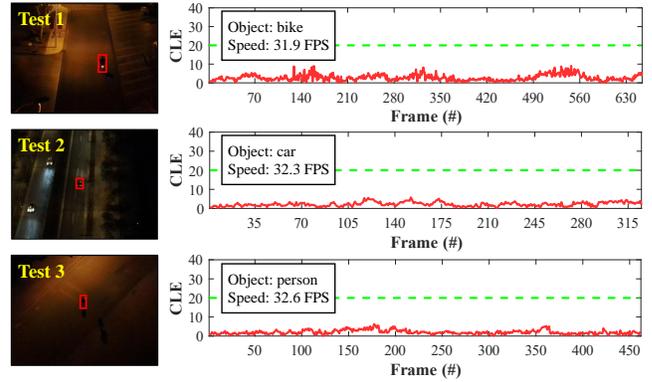


Fig. 7. Evaluations on a standard UAV platform in the real world. The estimated positions are indicated by red bounding boxes. Below are CLE curves between predictions and ground truth. The green dashed line marks a 20-pixel threshold, within which tracking mistakes are usually considered acceptable. The base tracker achieves good nighttime tracking with the help of HighlightNet.

line of TABLE IV, the gains of HighlightNet bringing to both human perception and tracking performance degraded significantly. Therefore, the Transformer-based parameter adjustment module is of clear benefit to UAV tracking tasks.

3) *ST & \mathcal{L}_{dan}* : Ablating ST and \mathcal{L}_{dan} respectively, the benefit of introducing HighlightNet decreases to some extent, validating the effectiveness of both anti-noise mask and dark area noise loss. Furthermore, as shown in the second line of TABLE IV, simultaneously disabling ST and \mathcal{L}_{dan} brings a more obvious reduction in success rate, precision, PSNR, and SSIM. Therefore, it is testified that ST should be employed along with the \mathcal{L}_{dan} to minimize the over-enhancement of noise.

E. Real-World Tests

To testify HighlightNet's practicability in real-world low-light UAV tracking applications, we employ it on a typical UAV platform with embedded devices, *i.e.*, the NVIDIA Jetson AGX Xavier. HighlightNet offers a conspicuous real-time performance of **32.2 FPS** even without TensorRT acceleration. In addition, Fig. 7 shows CLE curves and various real-world nighttime tracking experiments. The main obstacles in the testing are illumination variation, low brightness, tiny target, and partial occlusion. The prediction errors are less than 20 pixels on the CLE curves, showing that the tracking is accurate. With the help of HighlightNet, the basic tracker can provide reliable object tracking in nighttime conditions.

V. CONCLUSION

This work proposes a low-light enhancer HighlightNet for facilitating both online target selection step and tracking step in UAV nighttime tracking. Therefore, HighlightNet takes both human perception and nighttime challenges for UAVs into consideration. By three novel modules, it highlights the potential features and consequently reduces the influence of rapid illumination variation, artificial light sources, small objects, and image noise. Experiments on various tracking approaches confirm its compatibility and effectiveness in

both online target selection and tracking. Comparison with other SOTA low-light enhancers verifies the advantages of HighlightNet for facilitating human perception and tracking performance in nighttime conditions. Real-world experiments on a typical UAV platform confirm its applicability and dependability while consuming little computational resources. To summarize, we are confident that this research will aid in the expansion of UAV tracking applications to nighttime environments.

ACKNOWLEDGMENT

This work is supported by the National Natural Science Foundation of China (No. 62173249), the Natural Science Foundation of Shanghai (No. 20ZR1460100).

REFERENCES

- [1] H. Cheng, L. Lin, Z. Zheng, Y. Guan, and Z. Liu, "An Autonomous Vision-Based Target Tracking System for Rotorcraft Unmanned Aerial Vehicles," in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2017, pp. 1732–1738.
- [2] C. Meng, B. He, H. Lin, and L. Wu, "Research on UAV Point Landing Based on Visual Navigation," in *Proceedings of the IEEE CSAA Guidance, Navigation and Control Conference (CGNCC)*, 2018, pp. 1–6.
- [3] J. Ye, C. Fu, F. Lin, F. Ding, S. An, and G. Lu, "Multi-Regularized Correlation Filter for UAV Tracking and Self-Localization," *IEEE Transactions on Industrial Electronics*, vol. 69, no. 6, pp. 6004–6014, 2022.
- [4] B. Li, J. Yan, W. Wu, Z. Zhu, and X. Hu, "High Performance Visual Tracking with Siamese Region Proposal Network," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 8971–8980.
- [5] Z. Cao, C. Fu, J. Ye, B. Li, and Y. Li, "HiFT: Hierarchical Feature Transformer for Aerial Tracking," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021, pp. 15437–15446.
- [6] B. Li, C. Fu, F. Ding, J. Ye, and F. Lin, "All-Day Object Tracking for Unmanned Aerial Vehicle," *IEEE Transactions on Mobile Computing*, pp. 1–13, 2022.
- [7] J. Ye, C. Fu, G. Zheng, Z. Cao, and B. Li, "DarkLighter: Light Up the Darkness for UAV Tracking," in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2021, pp. 3079–3085.
- [8] X. Guo, Y. Li, and H. Ling, "LIME: Low-Light Image Enhancement via Illumination Map Estimation," *IEEE Transactions on Image Processing*, vol. 26, no. 2, pp. 982–993, 2017.
- [9] C. Wei, W. Wang, W. Yang, and J. Liu, "Deep Retinex Decomposition for Low-light Enhancement," in *Proceedings of the British Machine Vision Conference (BMVC)*, 2018, pp. 127–136.
- [10] Y. Zhang, J. Zhang, and X. Guo, "Kindling the Darkness: A Practical Low-Light Image Enhancer," in *Proceedings of the ACM International Conference on Multimedia (MM)*, 2019, pp. 1632–1640.
- [11] Y. Jiang, X. Gong, D. Liu, Y. Cheng, C. Fang, X. Shen, J. Yang, P. Zhou, and Z. Wang, "EnlightenGAN: Deep Light Enhancement Without Paired Supervision," *IEEE Transactions on Image Processing*, vol. 30, pp. 2340–2349, 2021.
- [12] C. Guo, C. Li, J. Guo, C. C. Loy, J. Hou, S. Kwong, and R. Cong, "Zero-Reference Deep Curve Estimation for Low-Light Image Enhancement," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 1777–1786.
- [13] A. Vaswani, N. Shazeer, N. Parmar *et al.*, "Attention Is All You Need," in *Proceedings of the Advances in Neural Information Processing Systems (NIPS)*, 2017, pp. 6000–6010.
- [14] M. Danelljan, G. Bhat, F. S. Khan, and M. Felsberg, "ECO: Efficient Convolution Operators for Tracking," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 6638–6646.
- [15] B. Li, W. Wu, Q. Wang, F. Zhang, J. Xing, and J. Yan, "SiamRPN++: Evolution of Siamese Visual Tracking With Very Deep Networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 4277–4286.
- [16] Y. Xu, Z. Wang, Z. Li, Y. Yuan, and G. Yu, "SiamFC++: Towards Robust and Accurate Visual Tracking with Target Estimation Guidelines," in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, vol. 34, no. 07, 2020, pp. 12549–12556.
- [17] C. Fu, Z. Cao, Y. Li, J. Ye, and C. Feng, "Onboard real-time aerial tracking with efficient siamese anchor proposal network," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–13, 2022.
- [18] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 6, pp. 1137–1149, 2017.
- [19] C. Fu, S. Li, X. Yuan, J. Ye, Z. Cao, and F. Ding, "Ad²Attack: Adaptive Adversarial Attack on Real-Time UAV Tracking," in *Proceedings of the International Conference on Robotics and Automation (ICRA)*, 2022, pp. 5893–5899.
- [20] Z. Cao, Z. Huang, L. Pan, S. Zhang, Z. Liu, and C. Fu, "TC-Track: Temporal Contexts for Aerial Tracking," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 14798–14808.
- [21] E. H. Land, "The Retinex Theory of Color Vision," *Scientific American*, vol. 237, no. 6, pp. 108–129, 1977.
- [22] K. G. Lore, A. Akintayo, and S. Sarkar, "LLNet: A deep autoencoder approach to natural low-light image enhancement," *Pattern Recognition*, vol. 61, pp. 650–662, 2017.
- [23] R. Wang, Q. Zhang, C.-W. Fu, X. Shen, W.-S. Zheng, and J. Jia, "Underexposed Photo Enhancement Using Deep Illumination Estimation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 6842–6850.
- [24] K. Xu, X. Yang, B. Yin, and R. W. Lau, "Learning to Restore Low-Light Images via Decomposition-and-Enhancement," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 2278–2287.
- [25] R. Liu, L. Ma, J. Zhang, X. Fan, and Z. Luo, "Retinex-inspired Unrolling with Cooperative Prior Architecture Search for Low-light Image Enhancement," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 10556–10565.
- [26] J. Ye, C. Fu, Z. Cao, S. An, G. Zheng, and B. Li, "Tracker Meets Night: A Transformer Enhancer for UAV Tracking," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 3866–3873, 2022.
- [27] J. Ye, C. Fu, G. Zheng, D. P. Paudel, and G. Chen, "Unsupervised Domain Adaptation for Nighttime Aerial Tracking," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 8896–8905.
- [28] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-End Object Detection with Transformers," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020, pp. 213–229.
- [29] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai *et al.*, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021, pp. 1–12.
- [30] J. Cai, S. Gu, and L. Zhang, "Learning a Deep Single Image Contrast Enhancer from Multi-Exposure Images," *IEEE Transactions on Image Processing*, vol. 27, no. 4, pp. 2049–2062, 2018.
- [31] Z. Cao, C. Fu, J. Ye, B. Li, and Y. Li, "SiamAPN++: Siamese Attentional Aggregation Network for Real-Time UAV Tracking," in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2021, pp. 3086–3092.
- [32] M. Mueller, N. Smith, and B. Ghanem, "A Benchmark and Simulator for UAV Tracking," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016, pp. 445–461.