

Extremely Low-Light Image Enhancement with Scene Text Restoration

Po-Hao Hsu^{*§}, Che-Tsung Lin^{†§}, Chun Chet Ng^{‡§}, Jie Long Kew[‡], Mei Yih Tan[‡],
Shang-Hong Lai^{*¶}, Chee Seng Chan[‡] and Christopher Zach[†]

^{*}National Tsing Hua University, Hsinchu, Taiwan

[†]Chalmers University of Technology, Gothenburg, Sweden

[‡]Universiti Malaya, Kuala Lumpur, Malaysia

[¶]Microsoft AI R & D Center, Taipei, Taiwan

Abstract—Deep learning-based methods have made impressive progress in enhancing extremely low-light images - the image quality of the reconstructed images has generally improved. However, we found out that most of these methods could not sufficiently recover the image details, for instance, the texts in the scene. In this paper, a novel image enhancement framework is proposed to precisely restore the scene texts, as well as the overall quality of the image simultaneously under extremely low-light conditions. Mainly, we employed a self-regularised attention map, an edge map, and a novel text detection loss. In addition, leveraging the synthetic low-light images is beneficial for image enhancement on the genuine ones in terms of text detection. The quantitative and qualitative experimental results have shown that the proposed model outperforms state-of-the-art methods in image restoration, text detection, and text spotting on See In the Dark and ICDAR15 datasets.

I. INTRODUCTION

Image enhancement under extremely low-light conditions is very challenging as under such an extreme scenario, a large amount of information is usually not available, as shown in Figure 1(a). At the same time, these images are also highly susceptible to noise due to demosaicing in the image sensing pipeline. Although it is possible to capture a better image with either a larger aperture or a longer exposure time, the images might then suffer from being overexposed or blurred due to improper camera settings or/and object motions.

Recently, many low-light enhancement algorithms have been proposed. Traditional approaches [1]–[5] usually aim at restoring the statistics of the original images to that of the normal or natural ones; while deep learning-based methods [6]–[15] mainly aim to learn the mapping between low-light images and brighter ones via regression.

Overall, the quality of these low-light images has been improved to a certain extent. However, we found out that the degree of visibility or readability of contextual information, for instance, the scene texts in these enhanced low-light images have never been explicitly discussed or stressed. As an example, empirically, we notice that current low-light image enhancement models generally suffer from losing finer details, especially in the text regions. That is to say, although the existing methods can achieve moderate image quality scores,

[§]These authors contributed equally to this work. Corresponding author: Chee Seng Chan (*cs.chan@um.edu.my*)

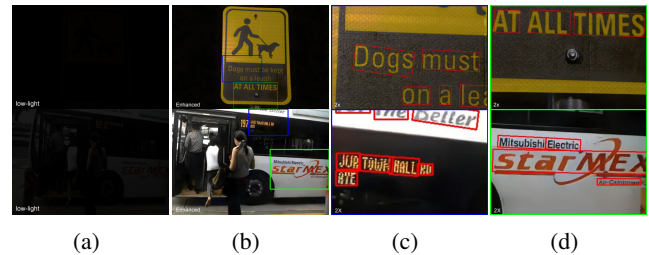


Fig. 1: From left to right: (a) Original low-light images; (b) Enhanced results with our proposed method; (c-d) Zoomed-in (2x) of the blue and green bounding box. It is obvious that the texts are clearly visible with sharp edges.

e.g., peak signal-to-noise ratio (PSNR) and structural similarity (SSIM) in overall, they seem to fail in the downstream tasks, such as text detection and text spotting in our case as shown in Figure 3-4 (b)-(h). On the other hand, with the success of deep neural networks, text detection [16], [17] and text recognition [18], [19] models have made significant progress. For example, they can work on extremely low-light images to a small extent (see Figure 4 (a)), but we believe that the results could be significantly improved (see Figure 4 (i)) if the images are better restored in terms of both the overall image quality and the local delicate text features.

In this paper, we are inspired to address this problem by focusing on restoring texts in scene images captured under extremely low-light conditions and the overall image quality simultaneously. Diving into more details, below is the summary of contributions of this paper:

- Firstly, we proposed an image enhancement framework that is capable to improve both the low-light image quality (in overall) and the scene text simultaneously through a novel text detection loss. It enables the enhanced images to preserve finer low-level details such as character edges and shapes without compromising the overall quality of the image, leading to a successful text detection (see Figure 1, Tables I-II).
- Secondly, in terms of datasets, we (i) annotated the texts in the real low-light dataset - See In the Dark (SID) Sony [20], and (ii) created a synthetic low-light dataset

based on the commonly used ICDAR15 [21] dataset. Code and datasets are available at https://github.com/SheepHow/ELIE_STR.

- Thirdly, extensive experiments have shown that our proposed model achieve the best scores in terms of text detection and spotting tasks on the enhanced low-light images in both See In the Dark and ICDAR15 datasets. We also showed that our image enhancement model trained with additional synthetic low-light images can achieve better results as compared to training on the real SID Sony set only (Table V).

II. RELATED WORKS

Generally, low-light image enhancement can be categorized into two main approaches: traditional methods and deep learning-based ones. Over the years, traditional Histogram Equalization (HE) and its variations [1]–[3] have been widely studied. To overcome their limitations, Retinex-based algorithms, such as SSR [4], MSR [5], SIRE [22], LIME [23], and BIMEF [24] were proposed and assume that an image can be decomposed into illumination and reflectance. Recently, deep learning-based image denoising and enhancement tasks have achieved significant improvements. DCNN [6] proposes a joint framework with a CNN-based denoising module. LLCNN [7] designs a special module to utilize multi-scale feature maps. LLNet [8] uses a stacked-sparse denoising autoencoder to identify signal features and adaptively enhances and denoises the image. Inspired by bilateral grid processing and local affine color transforms, HDRNet [9] predicts image transformation by learning to make local, global, and content-dependent decisions. Retinex-Net [10] combines deep learning and retinex theory, and adjusts illumination for enhancement after image decomposition.

Furthermore, the recent successes of GANs [25] have also attracted attention from the low-light image enhancement community because GANs have proven successful in image synthesis and translation. Specifically, Pix2pix [13] can provide visually plausible images in the target domain given paired training data. However, it is challenging to obtain paired images in many practical applications. As such, cycle-consistency introduced by CycleGAN [11] had opened up the possibility of performing unpaired image-to-image translation. To overcome the complexity of CycleGAN, EnlightenGAN [12] proposed an unsupervised one-path GAN structure which includes a global discriminator and a local one, a self-regularized perceptual loss fusion, and an attention mechanism. EEMEFN [14] proposed an edge enhancement module to enhance the initial image generated by the multi-exposure fusion (MEF) module; however, the model focuses on raw images only. RetinexGAN [15] presented a method combining the Retinex theory and GAN.

One of the nearest works to us is possibly Xue et al. [26] that proposed a low light text detector based on spatial and frequency feature fusion to enhance the fine details of low light image. However, in contrast to their work, we propose an

image enhancement framework to jointly improve the image and text quality as a whole, as illustrated in Figure 2.

III. PROPOSED METHOD

This paper introduces a novel framework that can simultaneously enhance the scene text and overall image quality under extreme low-light conditions. Our proposed framework consists of two modules, as illustrated in Figure 2. The first module is the image-translation module that consists of a U-net which inputs a combination of a low-light image, an edge map and an attention map, while the second module is the loss calculation module which is only involved in the training phase and consists of an ℓ_1 loss, MS-SSIM loss, and our novel text detection loss.

A. Image-Translation Module

Herein, we adopt the U-net generator with refinements. i.e., edge map and self-regularized attention map which will be described next. As such, our U-net generator concatenates the low-light image I and the edge map E as a 4-channel input with self-regularized attention map S :

$$I' = \mathcal{F}(I, E, S). \quad (1)$$

1) **Edge map:** RCF edge detector [27] can be directly employed to predict edges from the low-light images. The idea is that ℓ_1 or mean squared error loss tends to blur sharp edges and other fine image details. Therefore, leveraging edge information can recover abundant textures and sharp edges and is helpful in the following downstream tasks, such as text detection. Technically, the RCF edge detector is pretrained on the BSDS500 dataset [28] based on the VGG16 architecture. The network consists of five stages that will extract multiple levels of edge features containing each layer’s essential fine details. Then, these features from all different stages are combined by a fusion layer. In this paper, instead of directly using RCF to predict edges, which might lead to blurred edges and artifacts in the dark regions, we trained another U-Net to produce the edge maps.

2) **Self-regularized attention map:** Low-light image enhancement can easily cause the output image to be over or under-exposure. Recently, EnlightenGAN [12] proposed an easy-to-use attention mechanism called self-regularized attention to handle this issue by focusing on the dark areas in the images. Technically, the self-regularized attention map S is attained by taking the illumination channel Y of the input RGB image, normalizing it to $[0,1]$, followed by $1 - Y$ (element-wise difference). Herein, we also found this trick helpful in extremely low-light image enhancement. In our work, the attention map is consistently re-scaled through max pooling layers and element-wisely multiplied with the feature maps in the encoding part. Then, these re-weighted features are passed to the corresponding feature maps in the decoding part at different scales as skip-connections.

As a whole, our U-Net consists of 9 convolutional blocks and a 1×1 convolutional layer on top. Each block consists of two 3×3 convolutional layers followed by LeakyReLU.

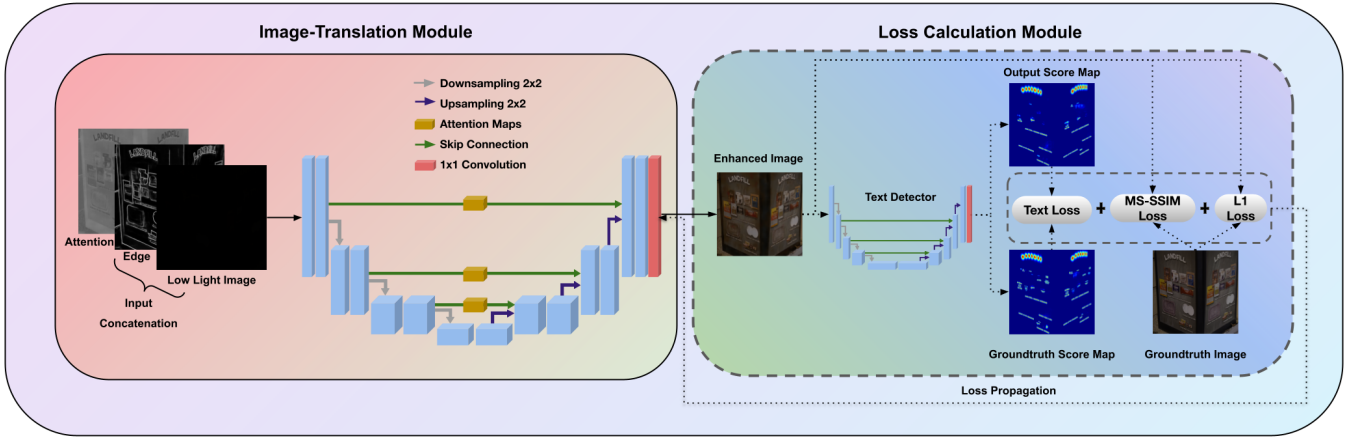


Fig. 2: Illustration of the overall architecture of our proposed low-light scene text restoration framework. Starting with the input concatenation of low-light image and edge map, the core enhancement network is guided by (1) an attention module through the multiplication of feature map with attention map, and followed by (2) the calculation of novel text detection loss where the model focuses on scene text regions. The dashed lines show how the text detection loss is formulated and propagated back to the main enhancement network so that the image enhancement and text restoration tasks can be trained simultaneously.

Each block is connected with a 2×2 max pooling layer at the downsampling step. We also apply max pooling on the attention map simultaneously. At the upsampling step, we extract the feature map with 2×2 transposed convolution and concatenate it with the corresponding feature map from the downsampling step, which passes through the attention maps.

B. Loss Calculation Module

As aforementioned, previous image enhancement methods suffered from not being able to restore scene text regions sufficiently in the low-light images. Therefore, in this module, we introduced a text detection loss \mathcal{L}_{text} , as well as two other image enhancement loss ($\mathcal{L}_{SSIM_{MS}}$ and \mathcal{L}_{ℓ_1}) as a joint loss function to simultaneously improve the overall image quality, as well as the scene text quality.

1) **Text detection loss:** Intuitively, a well-restored scene text implies that we could obtain a very similar text detection results on the enhanced image and the ground truth. In this work, we employ CRAFT [16] as our text detector to effectively localize the text characters and then estimate the affinity between them for text detection. Technically, our CRAFT network architecture is based on VGG16 with batch normalization as the backbone. There are skip connections in the decoding part, which is similar to U-Net in aggregating low-level features. CRAFT model will predict two Gaussian heatmaps: (i) *region score* and (ii) *affinity score*. In the former, the *region score* is the probability of the characters, and in the latter, the *affinity score* represents the probability that adjacent characters are in the same word. As such, we propose the text detection loss \mathcal{L}_{text} as:

$$\mathcal{L}_{text} = \|R(I') - R(I^{GT})\|_1, \quad (2)$$

where $R(I')$ and $R(I^{GT})$ denote the region score maps of the enhanced image and the ground-truth image, respectively; while w and h denote the width and height of the region score map.

2) **Image enhancement loss:** The multi-scale SSIM method was proposed in [29] for reference-based image quality assessment, focusing on the image structure consistency. An M -scale SSIM between I' and I^{GT} is given by

$$SSIM_{MS}(I', I^{GT}) = [l_M(I', I^{GT})]^\alpha \prod_{j=1}^M [c_j(I', I^{GT})]^\beta [s_j(I', I^{GT})]^\gamma, \quad (3)$$

where l_M is the luminance at M -scale; c_j and s_j represent the contrast and the structure similarity measures at the j -th scale, respectively; while α , β , and γ are parameters to adjust the importance of the three components.

Inspired by [29], we adopted the multi-scale SSIM loss function in our work to enforce the image structure of the enhanced image I' to be close to that of the ground-truth image I^{GT} :

$$\mathcal{L}_{SSIM_{MS}} = 1 - SSIM_{MS}(I', I^{GT}). \quad (4)$$

Additionally, in order to better enforce correctness at the low frequencies [13], we also employ ℓ_1 loss between I' and I^{GT} as:

$$\mathcal{L}_{\ell_1} = \|I' - I^{GT}\|_1. \quad (5)$$

As summary, the overall joint loss function to train our image enhancement model is given by:

$$\mathcal{L}_{Total} = \omega_1 \mathcal{L}_{\ell_1} + \omega_2 \mathcal{L}_{SSIM_{MS}} + \omega_3 \mathcal{L}_{text}, \quad (6)$$

where ω is the hyperparameter.

	SRIE [22]	LIME [23]	BIMEF [24]	RetinexNet [10]	CycleGAN [11]	EnlightenGAN [12]	Pix2pix [13]	Ours
Sony	12.59/0.104	13.87/0.135	12.87/0.110	15.49/0.368	15.34/0.453	14.59/0.426	21.07/0.662	25.51/0.716
Syn. ICDAR15	10.42/0.409	12.04/0.522	16.00/0.581	15.55/0.637	23.42/0.720	21.03/0.661	24.87/0.727	28.41/0.840

TABLE I: Quantitative evaluation of low-light image enhancement algorithms in terms of PSNR/SSIM.

Model	Text Detection H-Mean				Two-Stage Text Spotting Case Insensitive Accuracy			
	Sony		Syn. ICDAR 15		Syn. ICDAR 15			
	CRAFT	PAN	CRAFT	PAN	CRAFT + TRBA	CRAFT + ASTER	PAN + TRBA	PAN + ASTER
Input	0.057	0.026	0.355	0.192	0.114	0.118	0.062	0.067
SRIE [22]	0.133	0.076	0.465	0.423	0.127	0.144	0.117	0.122
LIME [23]	0.127	0.057	0.454	0.428	0.129	0.146	0.120	0.128
BIMEF [24]	0.136	0.079	0.450	0.411	0.124	0.138	0.123	0.122
RetinexNet [10]	0.115	0.040	0.374	0.325	0.090	0.096	0.069	0.076
CycleGAN [11]	0.090	0.053	0.428	0.458	0.119	0.144	0.122	0.138
EnlightenGAN [12]	0.146	0.075	0.458	0.461	0.140	0.157	0.123	0.139
Pix2pix [13]	0.266	0.190	0.559	0.542	0.183	0.207	0.182	0.195
Ours	0.324	0.266	0.623	0.631	0.193	0.219	0.197	0.221
GT	0.842	0.661	0.800	0.830	0.526	0.584	0.555	0.591

TABLE II: Quantitative evaluation of enhanced images of all image enhancement algorithms in terms of H-Mean for text detection, and we use case insensitive word accuracy as the main text spotting score. Scores in bold are the best of all.

IV. EXPERIMENTAL RESULTS

A. Experiment Setup

Datasets. Two public datasets, **SID** [20] Sony and **ICDAR15** [21], are employed in this work. **SID** Sony was captured by Sony α 7S II contains 2697 short-exposure images and 231 long-exposure images at the resolution of 4240×2832. The exposure time of the short-exposure images was set to 1/30, 1/24, and 1/10 seconds. The corresponding long-exposure images were captured with 10 and 30 seconds. In our experiments, we convert the SID Sony images to 24-bit RGB format and manually label the scene text bounding boxes. As a result, there are 8210 and 611 labels for short-exposure and long-exposure images, respectively. **ICDAR15** dataset was introduced in the ICDAR 2015 Robust Reading Competition for incidental scene text detection and recognition. It contains 1500 scene text images at the resolution of 1280×720. In this work, the brightness of each image in ICDAR15 is reduced to make it visually similar to SID Sony dataset, so that it forms paired images in the low-light image enhancement application. **Metrics.** We compare our model with state-of-the-art methods in terms of PSNR and SSIM to measure the image quality. Also, we follow the common standard to use H-Mean as the evaluation metric for text detection. For text spotting, case insensitive word accuracy is used but lexicon is not employed so that we can study the raw impact of different sets of the enhanced images. Finally, the SID Sony dataset is only analyzed in terms of text detection because text recognition labels are not available.

Implementation Details. We train our network for 4000 epochs using Adam optimizer [30]. The initial learning rate is set to 1e-4 which is decreased to 1e-5 after 2000 epochs. In each training iteration, we randomly crop a 512×512 patch, and apply random flipping together with rotation as data augmentation. The parameters in the loss function ($\omega_1, \omega_2, \omega_3$) are set to (0.85, 0.15, 0.425) empirically.

B. Quantitative comparison

We compare our model with traditional approaches, as well as deep learning-based. Table I reports PSNR and SSIM to indicate how much the image quality has improved. It is noticed that our method is able to obtain the highest scores on both SID Sony and ICDAR15 datasets, which signifies that our method can enhance the low-light images to be as close as possible to the original images. Table II shows the scene text detection scores and text spotting accuracies of our model and other related works.

Scene Text Detection. We run two state-of-the-art scene text detectors (CRAFT [16] and PAN [17]) on the enhanced images. As shown in Table II, our model achieves the highest H-Mean on SID Sony and synthetic ICDAR15 datasets.

Scene Text Spotting. In order to show that both detection and recognition tasks are highly correlated and better enhanced images will lead to a quantitatively higher text detection and text recognition accuracy at the same time, we carry out a two-stage text spotting experiment using the aforementioned detectors and two robust scene text recognizers (TRBA [18] and ASTER [19]). In detail, we first select detection results that have Intersection over Union (IoU) greater than 0.5 when compared to the ground-truth bounding boxes. Then, we crop out the detected text regions before passing them to the text recognizers. As shown in Table II, when both recognizers are fed with the text detection results of CRAFT [16] and PAN [17] respectively, it is evident that the recognition accuracy of the images enhanced by our model consistently outperforms the ones of other competing methods. This shows that our proposed method is able to enhance the overall images, as well as restoring the text regions where our results excel in terms two-stage text spotting. In contrary, current existing methods might introduce a lot of noises and artifacts due to the inferior image enhancement capability.

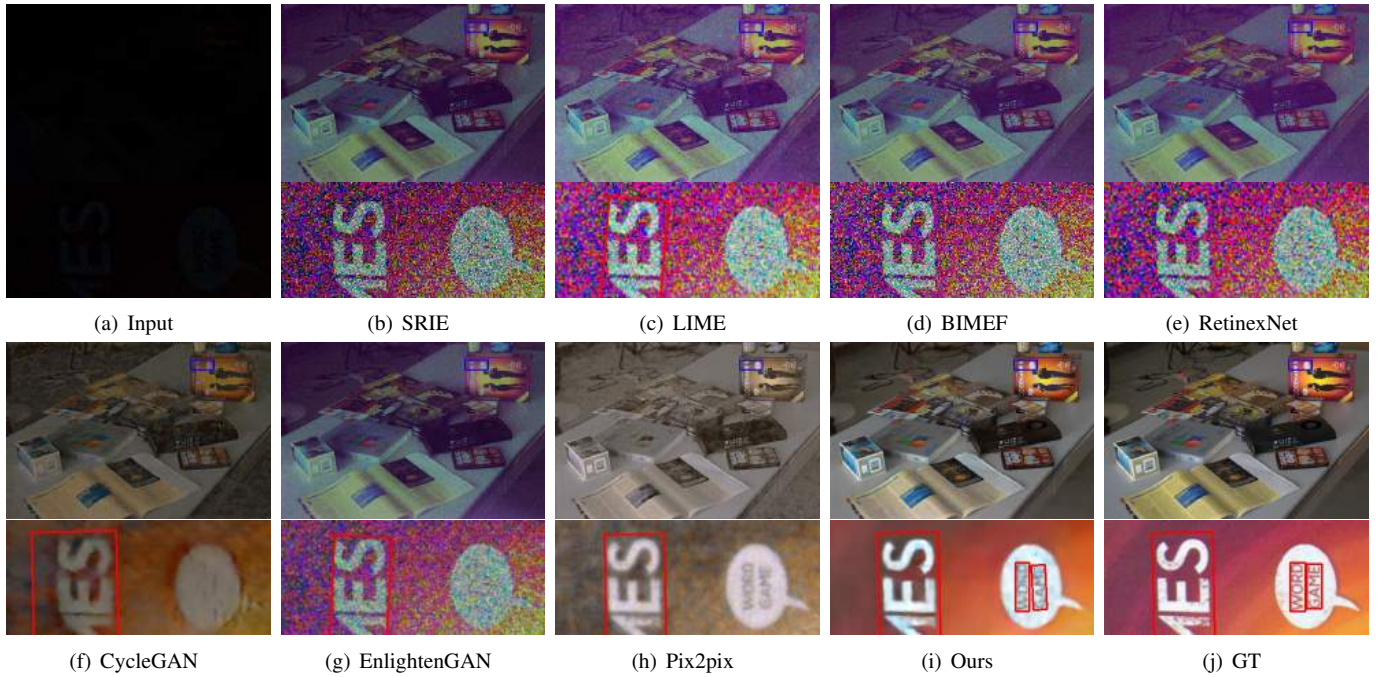


Fig. 3: CRAFT’s results (red boxes) of all methods on the SID Sony dataset, focused on the zoomed-in regions (blue boxes).

C. Qualitative comparison

Figure 3 and 4 demonstrate the image enhancement results on SID Sony and synthetic ICDAR15 datasets, respectively. Figures 3a and 4a are the original low-light input, and (b)-(i) are the images enhanced by SRIE [22], LIME [23], BIMEF [24], RetinexNet [10], CycleGAN [11], EnlightenGAN [12], Pix2pix [13] and our proposed method. The last image in both figures is the ground truth image. We also show the zoomed-in details on selected scene text regions, overlaid with the text detection bounding boxes.

It is noticed that the results obtained with LIME and RetinexNet contain overexposed noises, while SRIE and BIMEF can barely recover the images. Then, CycleGAN and Pix2pix have color distortion and could not restore details. Overall, these methods fail to recover the text regions to the degree that the text detector could detect them accurately.

Scene Text Detection. It can be noticed that the proposed model is able to enhance the low-light image and achieves the best performance in terms of scene text detection. For instance, in Figure 3, only the images enhanced by our method can provide text detection results almost identical to the ground truth. Existing methods either miss certain words or fail to predict each word individually. This happens because only our method can produce enhanced images with less noise and finer details in terms of text detection.

Scene Text Spotting. Upon visualizing the recognition results of ASTER in Figure 4 for the word “BETTER” in the ground truth, only our method is able to recognize the word correctly. While for other methods, recognition results of low-light input, RetinexNet, and EnlightenGAN are “after”. SRIE and LIME leads to the output of “ame” and “bitter” respectively.

BIMEF yields “bmtp”, Pix2pix predicts “am”, and “arms” for CycleGAN. These results show that our method is able to enhance and restore the text at a more refined level, and it can preserve important low-level features such as edges and character strokes to achieve accurate text recognition results.

D. Ablation Study

To understand the effect of each component of our model, we conduct several ablation experiments by either adding or removing them one at a time. We employ U-Net architecture using ℓ_1 loss as our baseline. In Table III and Table IV, the first column shows the performance of the baseline, and the rightmost column is the result of our full model. We can observe that each component contributes to a higher text detection accuracy to a certain extent. For example, the self-attention module guides the network to focus on the dark areas instead of the bright ones, while the edge map demonstrates that the network can extract more information from the sharper edges and fewer artifacts. Multi-scale SSIM loss leads to a better text detection result because the enhanced images are closer to the original ones. The text detection loss encourages the model to restore the texts to the extent of being detectable by the text detector instead of solely improving the overall image quality. As a whole, each component contributes to the overall better text detection and recognition results for the enhanced extremely low-light images.

Also, we test our model further by training it with a mixture of real (SID Sony) and synthetic low-light (ICDAR15) datasets to study if synthetic data can help to improve the image enhancement result. The top 3 models on SID Sony were selected in this experiment. Results in Table V showed that



Fig. 4: CRAFT’s results (red boxes) of all methods on the ICDAR15 dataset, focused on the zoomed-in regions (blue boxes), followed by ASTER’s results based on the ground-truth bounding box of “BETTER”.

Baseline	w/ A	w/ A+E	w/ A+E+M	Full Model
0.235	0.242	0.269	0.291	0.324

TABLE III: Ablation study of our full model on SID Sony with each component in terms of H-Mean. A: attention map; E: edge map; M: MS-SSIM; T: text detection loss.

Baseline	w/o A	w/o E	w/o M	w/o T	Full Model
0.235	0.308	0.302	0.304	0.307	0.324

TABLE IV: Ablation study of our model on SID Sony without each component one at a time in terms of H-Mean. A: attention map; E: edge map; M: MS-SSIM; T: text detection loss.

Model	CRAFT	PAN
EnlightenGAN [12]	0.205	0.146
Pix2pix [13]	0.281	0.224
Ours	0.348	0.278

TABLE V: H-Mean on SID Sony when trained on SID Sony and synthetic ICDAR15 datasets.

there is a significant H-Mean increment when the models are evaluated on the real (SID Sony) test set. We observe that our model gains H-Mean score of 0.024 using CRAFT and 0.012 for PAN. This shows that synthetic low-light data is able to fill up the gap caused by the scarcity of real labeled low-light images and justifies the creation of a synthetic ICDAR15 dataset for such purpose.

V. CONCLUSIONS

This paper presents an attention-guided and edge-awareness CNN model for extremely low-light image enhancement, particularly on the presence of scene texts. The incorporated self-regularized attention map is proven effective for guiding the network to focus on the dark regions. The predicted edge map given the attention map is quantitatively helpful for recovering abundant textures and sharp edges. Most importantly, we propose a novel text detection loss that helps the network attend to those scene text regions to well-recover the scene texts. The experimental results demonstrate that the proposed method consistently outperforms state-of-the-art methods, including low-light image enhancement models and GAN-based ones, in terms of image restoration, text detection, and text spotting on challenging SID Sony and ICDAR15 datasets. Lastly, we also showed that our proposed model trained with additional synthetic extremely low-light images can achieve better results on genuine ones.

VI. FUTURE DIRECTIONS

Extending our proposed model to real-time would be beneficial to real-life applications such as collision avoidance of self-driving cars at nighttime and other object detection applications under extremely low-light settings. Besides, we would also like to create a larger dataset that is able to cover more scenes and texts with a wider variety of styles. We believe that these research directions would bring enormous value to both the industrial and research communities.

REFERENCES

- [1] S. Pizer, E. P. Amburn, J. D. Austin, R. Cromartie, A. Geselowitz, T. Greer, B. T. H. Romeny, and J. B. Zimmerman, "Adaptive histogram equalization and its variations," *Graphical Models graphical Models and Image Processing computer Vision, Graphics, and Image Processing*, vol. 39, pp. 355–368, 1987.
- [2] T. Çelik and T. Tjahjadi, "Contextual and variational contrast enhancement," *IEEE Transactions on Image Processing*, vol. 20, pp. 3431–3441, 2011.
- [3] C. Lee, C. Lee, and C.-S. Kim, "Contrast enhancement based on layered difference representation of 2d histograms," *IEEE Transactions on Image Processing*, vol. 22, no. 12, pp. 5372–5384, 2013.
- [4] D. J. Jobson, Z.-u. Rahman, and G. A. Woodell, "Properties and performance of a center/surround retinex," *IEEE Transactions on Image Processing*, vol. 6, no. 3, pp. 451–462, 1997.
- [5] D. Jobson, Z. Rahman, and G. A. Woodell, "A multiscale retinex for bridging the gap between color images and the human observation of scenes," *IEEE Transactions on Image Processing*, vol. 6, no. 7, pp. 965–976, 1997.
- [6] L. Tao, C. Zhu, J. Song, T. Lu, H. Jia, and X. Xie, "Low-light image enhancement using cnn and bright channel prior," in *ICIP*, 2017, pp. 3215–3219.
- [7] L. Tao, C. Zhu, G. Xiang, Y. Li, H. Jia, and X. Xie, "LLCNN: A convolutional neural network for low-light image enhancement," in *VCIP*, 2017.
- [8] 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.
- [9] M. Gharbi, J. Chen, J. Barron, S. W. Hasinoff, and F. Durand, "Deep bilateral learning for real-time image enhancement," *ACM Transactions on Graphics*, vol. 36, no. 4, pp. 1–12, 2017.
- [10] C. Wei, W. Wang, W. Yang, and J. Liu, "Deep retinex decomposition for low-light enhancement," *arXiv preprint arXiv:1808.04560*, 2018.
- [11] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *ICCV*, 2017.
- [12] 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," *ArXiv*, vol. abs/1906.06972, 2019.
- [13] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," in *CVPR*, 2017.
- [14] M. Zhu, P. Pan, W. Chen, and Y. Yang, "EEMEFN: Low-light image enhancement via edge-enhanced multi-exposure fusion network." in *AAAI*, 2020.
- [15] Y. Shi, X. Wu, and M. Zhu, "Low-light image enhancement algorithm based on retinex and generative adversarial network," *ArXiv*, vol. abs/1906.06027, 2019.
- [16] Y. Baek, B. Lee, D. Han, S. Yun, and H. Lee, "Character region awareness for text detection," in *CVPR*, 2019.
- [17] W. Wang, E. Xie, X. Song, Y. Zang, W. Wang, T. Lu, G. Yu, and C. Shen, "Efficient and accurate arbitrary-shaped text detection with pixel aggregation network," in *ICCV*, 2019.
- [18] J. Baek, G. Kim, J. Lee, S. Park, D. Han, S. Yun, S. J. Oh, and H. Lee, "What is wrong with scene text recognition model comparisons? dataset and model analysis," in *ICCV*, 2019.
- [19] B. Shi, M. Yang, X. Wang, P. Lyu, C. Yao, and X. Bai, "ASTER: An attentional scene text recognizer with flexible rectification," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, no. 9, pp. 2035–2048, 2019.
- [20] C. Chen, Q. Chen, J. Xu, and V. Koltun, "Learning to see in the dark," in *CVPR*, 2018.
- [21] D. Karatzas, L. G. I. Bigorda, A. Nicolaou, S. Ghosh, A. D. Bagdanov, M. Iwamura, J. Matas, L. Neumann, V. Chandrasekhar, S. Lu, F. Shafait, S. Uchida, and E. Valveny, "ICDAR 2015 competition on robust reading," in *ICDAR*, 2015.
- [22] X. Fu, D. Zeng, Y. Huang, X. Zhang, and X. Ding, "A weighted variational model for simultaneous reflectance and illumination estimation," in *CVPR*, 2016.
- [23] 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, 2016.
- [24] Z. Ying, G. Li, and W. Gao, "A bio-inspired multi-exposure fusion framework for low-light image enhancement," *arXiv preprint arXiv:1711.00591*, 2017.
- [25] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. C. Courville, and Y. Bengio, "Generative adversarial nets," in *NIPS*, 2014.
- [26] M. Xue, P. Shivakumara, C. Zhang, Y. Xiao, T. Lu, U. Pal, D. Lopresti, and Z. Yang, "Arbitrarily-oriented text detection in low light natural scene images," *IEEE Transactions on Multimedia*, 2020.
- [27] Y. Liu, M.-M. Cheng, X. Hu, J. Bian, L. Zhang, X. Bai, and J. Tang, "Richer convolutional features for edge detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, pp. 1939–1946, 2019.
- [28] P. Arbelaez, M. Maire, C. Fowlkes, and J. Malik, "Contour detection and hierarchical image segmentation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 33, no. 5, pp. 898–916, 2010.
- [29] Z. Wang, E. Simoncelli, and A. Bovik, "Multiscale structural similarity for image quality assessment," in *ACSSC*, 2003.
- [30] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2015.