# LLCaps: Learning to Illuminate Low-Light Capsule Endoscopy with Curved Wavelet Attention and Reverse Diffusion

Long Bai[1] *, Tong Chen[2] *, Yanan Wu[1,3], An Wang[1], Mobarakol Islam[4], and Hongliang Ren[1,5] **

[1] Department of Electronic Engineering, The Chinese University of Hong Kong (CUHK), Hong Kong SAR, China
[2] The University of Sydney, Sydney, NSW, Australia
[3] Northeastern University, Shenyang, China
[4] Wellcome/EPSRC Centre for Interventional and Surgical Sciences (WEISS), University College London, London, UK
[5] Shun Hing Institute of Advanced Engineering, CUHK, Hong Kong SAR, China
b.long@link.cuhk.edu.hk, tche2095@uni.sydney.edu.au, yananwu@cuhk.edu.hk, wa09@link.cuhk.edu.hk, mobarakol.islam@ucl.ac.uk, hlren@ee.cuhk.edu.hk

**Abstract.** Wireless capsule endoscopy (WCE) is a painless and non-invasive diagnostic tool for gastrointestinal (GI) diseases. However, due to GI anatomical constraints and hardware manufacturing limitations, WCE vision signals may suffer from insufficient illumination, leading to a complicated screening and examination procedure. Deep learning-based low-light image enhancement (LLIE) in the medical field gradually attracts researchers. Given the exuberant development of the denoising diffusion probabilistic model (DDPM) in computer vision, we introduce a WCE LLIE framework based on the multi-scale convolutional neural network (CNN) and reverse diffusion process. The multi-scale design allows models to preserve high-resolution representation and context information from low-resolution, while the curved wavelet attention (CWA) block is proposed for high-frequency and local feature learning. Moreover, we combine the reverse diffusion procedure to optimize the shallow output further and generate images highly approximate to real ones. The proposed method is compared with eleven state-of-the-art (SOTA) LLIE methods and significantly outperforms quantitatively and qualitatively. The superior performance on GI disease segmentation further demonstrates the clinical potential of our proposed model. Our code is publicly accessible at github.com/longbai1006/LLCaps.

## 1 Introduction

Currently, the golden standard of gastrointestinal (GI) examination is endoscope screening, which can provide direct vision signals for diagnosis and analysis. Benefiting from its characteristics of being non-invasive, painless, and low physical

---

* Long Bai and Tong Chen are co-first authors.
** Corresponding author.

**Fig. 1.** Comparison of normal images with low-light images. Obvious lesions are visible on normal images, but the same lesions can hardly be distinguished by human eyes in the corresponding low-light images.

burden, wireless capsule endoscopy (WCE) has the potential to overcome the shortcomings of conventional endoscopy [23, 35]. However, due to the anatomical complexity, insufficient illumination, and limited performance of the camera, low-quality images may hinder the diagnosis process [3]. Blood vessels and lesions with minor color changes in the early stages can be hard to be screened out [1, 17]. Fig. 1 shows WCE images with low illumination and contrast. The disease features clearly visible in the normal image become challenging to be found in the low-light images. Therefore, it is necessary to develop a low-light image enhancement framework for WCE to assist clinical diagnosis.

Many traditional algorithms (e.g., intensity transformation [9], histogram equalization [16], and Retinex theory [15]) have been proposed for low-light image enhancement (LLIE). For WCE, Long *et al.* [17] discussed adaptive fraction-power transformation for image enhancement. However, traditional methods usually require an ideal assumption or an effective prior, limiting their wider applications. Deep learning (DL) provides novel avenues to solve LLIE problems [8, 12, 18]. Some DL-based LLIE schemes for medical endoscopy have been proposed [7, 20]. Gomez *et al.* [7] offered a solution for laryngoscope low-light enhancement, and Ma *et al.* [20] proposed a medical image enhancement model with unpaired training data.

Recently, denoising diffusion probabilistic model (DDPM) [10] is the most popular topic in image generation, and has achieved success in various applications. Due to its unique regression process, DDPM has a stable training process and excellent output results, but also suffers from its expensive sampling procedure and lack of low-dimensional representation [21]. It has been proved that DDPM can be combined with other existing DL techniques to speed up the sampling process [21]. In our work, we introduce the reverse diffusion process of DDPM into our end-to-end LLIE process, which can preserve image details without introducing excessive computational costs. Our contributions to this work can be summarized as three-fold:

– We design a **L**ow-**L**ight image enhancement framework for **Caps**ule endoscopy (**LLCaps**). Subsequent to the feature learning and preliminary shallow image reconstruction by the convolutional neural network (CNN), the

reverse diffusion process is employed to further promote image reconstruction, preserve image details, and close in the optimization target.

– Our proposed curved wavelet attention (CWA) block can efficiently extract high-frequency detail features via wavelet transform, and conduct local representation learning with the curved attention layer.

– Extensive experiments on two publicly accessible datasets demonstrate the excellent performance of our proposed model and components. The high-level lesion segmentation tasks further show the potential power of LLCaps on clinical applications.
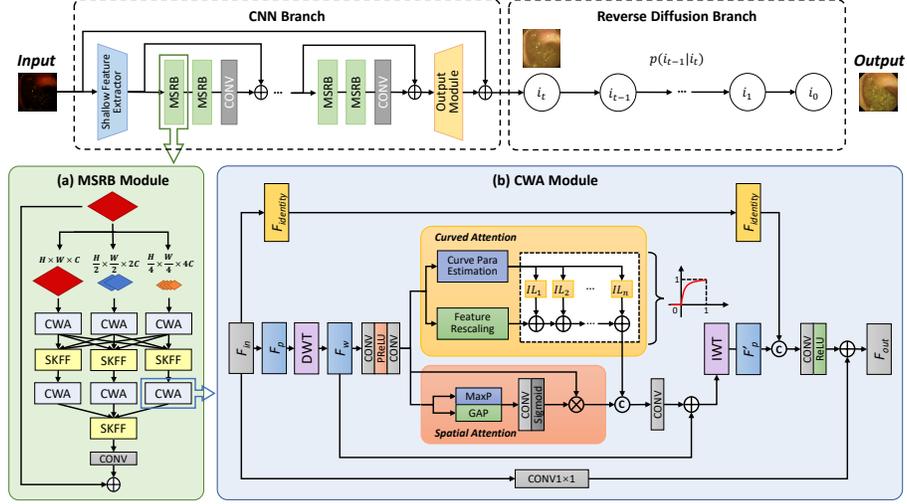
## 2   Methodology

### 2.1   Preliminaries

**Multi-scale Residual Block**  Multi-scale Residual Block (MSRB) [14] constructs a multi-scale neuronal receptive field, which allows the network to learn multi-scale spatial information in the same layer. Therefore, the network can acquire contextual information from the low-resolution features while preserving high-resolution representations. We establish our CNN branch with six stacked multi-scale residual blocks (MSRB), and every two MSRBs are followed by a 2D convolutional layer (Conv2D). Besides, each MSRB shall require feature learning and the multi-scale feature aggregation module. Specifically, we propose our curved wavelet attention (CWA) module to conduct multi-scale feature learning, and employ the selective kernel feature fusion (SKFF) [31] to combine multi-scale features, as shown in Fig. 2 (a).

**Denoising Diffusion Probabilistic Models**  Denoising Diffusion Probabilistic Models (DDPMs) [10] can be summarised as a model consisting of a forward noise addition $q(i_{1:T}|i_0)$ and a reverse denoising process $p(i_{0:T})$, which are both parameterized Markov chains. The forward diffusion process gradually adds noise to the input image until the original input is destroyed. Correspondingly, the reverse process uses the neural network to model the Gaussian distribution and achieves image generation through gradual sampling and denoising.

### 2.2   Proposed Methodology

**Curved Wavelet Attention**  Curved Wavelet Attention (CWA) block is the core component of our CNN branch, which is constructed via a curved dual attention mechanism and wavelet transform, as shown in Fig. 2 (b). Firstly, the input feature map $F_{in}$ is divided into identity feature $F_{identity}$ and processing feature $F_p$. Medical LLIE shall require high image details. In this case, we transform the $F_p$ into wavelet domain $F_w$ to extract high-frequency detail information based on discrete wavelet transform. $F_w$ is then propagated through the feature selector and dual attention module for deep representation learning. Finally, we

**Fig. 2.** The overview of our proposed LLCaps. The CNN branch shall extract the shallow image output while the DDPM branch further optimizes the image via Markov chain inference. (a) represents the multi-scale residual block (MSRB), which allows the model to learn representation on different resolutions. (b) denotes our curved wavelet attention (CWA) block for attention learning and feature restoration. In (b), DWT and IWT denote discrete wavelet transform and inverse wavelet transform, respectively. The PReLU with two convolutional layers constructs the feature selector. 'MaxP' denotes max pooling, and 'GAP' means global average pooling.

conduct reverse wavelet transform (IWT) to get $F'_p$, and concatenate it with $F_{identity}$ before the final output convolution layer.

We construct our curved dual attention module with parallel spatial and curved attention blocks. The spatial attention (SA) layer exploits the inter-spatial dependencies of convolutional features [31]. The SA layer performs the global average pooling and max pooling on input features respectively, and concatenates the output $F_{w(mean)}$ and $F_{w(max)}$ to get $F_{cat}$. Then the feature map will be dimensionally reduced and passed through the activation function.

However, literature [8,36] has discussed the problem of local illumination in LLIE. If we simply use a global computing method such as the SA layer, the model may not be able to effectively understand the local illumination/lack of illumination. Therefore, in order to compensate for the SA layer, we design the Curved Attention (CurveA) layer, which is used to model the high-order curve of the input features. Let $IL_{n(c)}$ denote the curve function, $c$ denote the feature location coordinates, and $Curve_{(n-1)}$ denote the pixel-wise curve parameter, we can obtain the curve estimation equation as:

$$\frac{IL_{n(c)}}{IL_{n-1(c)}} = Curve_{n-1}(1 - IL_{n-1(c)}) \tag{1}$$

The detailed CurveA layer is presented in the top of Fig. 2 (b), and the Equ. (1) is related to the white area. The Curve Parameter Estimation module consists of a Sigmoid activation and several Conv2D layers, and shall estimate the pixel-wise curve parameter at each order. The Feature Rescaling module will rescale the input feature into [0, 1] to learn the concave down curves. By applying the CurveA layer to the channels of the feature map, the CWA block can better estimate local areas with different illumination.

**Reverse Diffusion Process** Some works [21,28] have discussed combining diffusion models with other DL-based methods to reduce training costs and be used for downstream applications. In our work, We combine the reverse diffusion process of DDPM in a simple and ingenious way, and use it to optimize the shallow output by the CNN branch. Various experiments shall prove the effectiveness of our design in improving image quality and assisting clinical applications.

In our formulation, we assume that $i_0$ is the learning target $Y^*$ and $i_T$ is the output shallow image from the CNN branch. Therefore, we only need to engage the reverse process in our LLIE task. The reverse process is modeled using a Markov chain:

$$p_\theta \left( i_{0:T} \right) = p \left( i_T \right) \prod_{t=1}^{T} p_\theta \left( i_{t-1} \mid i_t \right)$$
$$p_\theta \left( i_{t-1} \mid i_t \right) = \mathcal{N} \left( i_{t-1}; \boldsymbol{\mu}_\theta \left( i_t, t \right), \boldsymbol{\Sigma}_\theta \left( i_t, t \right) \right) \tag{2}$$

$p_\theta \left( i_{t-1} \mid i_t \right)$ are parameterized Gaussian distributions whose mean $\boldsymbol{\mu}_\theta \left( i_t, t \right)$ and variance $\boldsymbol{\Sigma}_\theta \left( i_t, t \right)$ are given by the trained network. Meanwhile, we simplify the network and directly include the reverse diffusion process in the end-to-end training of the entire network. Shallow output is therefore optimized by the reverse diffusion branch to get the predicted image $Y$. We further simplify the optimization function and only employ a pixel-level loss on the final output image, which also improves the training and convergence efficiency.

**Overall Network Architecture** An overview of our framework can be found in Fig. 2. Our LLCaps contains a CNN branch (including a shallow feature extractor (SFE), multi-scale residual blocks (MSRBs), an output module (OPM)), and the reverse diffusion process. The SFE is a Conv2D layer that maps the input image into the high-dimensional feature representation $F_{SFE} \in \mathbb{R}^{C \times W \times H}$ [29]. Stacked MSRBs shall conduct deep feature extraction and learning. OPM is a Conv2D layer that recovers the feature space into image pixels. A residual connection is employed here to optimize the end-to-end training and converge process. Hence, given a low-light image $x \in \mathbb{R}^{3 \times W \times H}$, where $W$ and $H$ represent the width and height, the CNN branch can be formulated as:

$$F_{SFE} = H_{SFE}(x)$$
$$F_{MSRBs} = H_{MSRBs}(F_{SFE}), \tag{3}$$
$$F_{OPM} = H_{OPM}(F_{MB}) + x$$

The shallow output $F_{OPM} \in \mathbb{R}^{3 \times W \times H}$ shall further be propagated through the reverse diffusion process and achieve the final enhanced image $Y \in \mathbb{R}^{3 \times W \times H}$. The whole network is constructed in an end-to-end mode and optimized by Charbonnier loss [2]. The $\varepsilon$ is set to $10^{-3}$ empirically.

$$\mathcal{L}\left(x, x^{*}\right)=\sqrt{\|Y-Y^{*}\|^{2}+\varepsilon^{2}} \tag{4}$$

in which $Y$ and $Y^{*}$ denote the input and ground truth images, respectively.

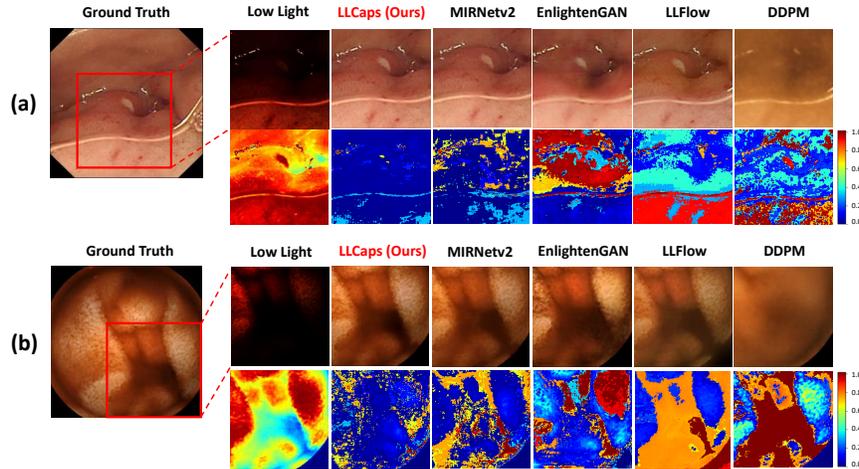## 3    Experiments

### 3.1    Dataset

We conduct our experiments on two publicly accessible WCE datasets, the Kvasir-Capsule [24] and the Red Lesion Endoscopy (RLE) dataset [5].

**Kvasir-Capsule dataset** [24] is a WCE classification dataset with three anatomy classes and eleven luminal finding classes. By following [4], we randomly select 2400 images from the Kvasir-Capsule dataset, of which 2000 are used for training and 400 for testing. To create low-light images, we adopt random Gamma correction and illumination reduction following [13, 18]. Furthermore, to evaluate the performance on real data, we add an external validation on 100 real images selected from the Kvasir-Capsule dataset. These images are with low brightness and are not included in our original experiments.

**Red Lesion Endoscopy dataset** [5] (RLE) is a WCE dataset for red lesion segmentation tasks (e.g., angioectasias, angiodysplasias, and bleeding). We randomly choose 1283 images, of which 946 images are used for training and 337 for testing. We adopt the same method in the Kvasir-Capsule dataset to generate low-light images. Furthermore, we conduct a segmentation task on the RLE test set to investigate the effectiveness of the LLIE models in clinical applications.

### 3.2    Implementation Details

We compare the performance of our LLCaps against the following state-of-the-art (SOTA) LLIE methodologies: LIME [9], DUAL [33], Zero-DCE [8], EnlightenGAN [12], LLFlow [26], HWMNet [6], MIRNet [31], SNR-Aware [30], Still-GAN [19], MIRNetv2 [32], and DDPM [10]. Our models are trained using Adam optimizer for 200 epochs with a batch size of 4 and a learning rate of $1 \times 10^{-4}$. For evaluation, we adopt three commonly used image quality assessment metrics: Peak Signal-to-Noise Ratio (PSNR) [11], Structural Similarity Index (SSIM) [27], and Learned Perceptual Image Patch Similarity (LPIPS) [34]. For the external validation set, we evaluate with no-reference metrics LPIPS [34] and Perception-based Image Quality Evaluator (PIQE) [25] due to the lack of ground truth images. To verify the usefulness of the LLIE methods for downstream medical tasks, we conduct red lesion segmentation on the RLE test set and evaluate the performance via mean Intersection over Union (mIoU), Dice similarity coefficient

**Fig. 3.** The quantitative results for LLCaps compared with SOTA approaches on (a) Kvasir-Capsule dataset [24] and (b) RLE dataset [5]. The first row visualizes the enhanced images from different LLIE approaches, and the second row contains the reconstruction error heat maps. The blue and red represent low and high error, respectively.

(Dice), and Hausdorff Distance (HD). We train UNet [22] using Adam optimizer for 20 epochs. The batch size and learning rate are set to 4 and $1 \times 10^{-4}$, respectively. All experiments are implemented by Python PyTorch and conducted on NVIDIA RTX 3090 GPU. Results are the average of 3-fold cross-validation.

### 3.3 Results

We compare the performance of our LLCaps to the existing approaches, as demonstrated in Table 3 and Fig. 3 quantitatively and qualitatively. Compared with other methods, our proposed method achieves the best performance among all metrics. Specifically, our method surpasses MIRNetv2 [32] by 3.57 dB for the Kvasir-Capsule dataset and 0.33 dB for the RLE dataset. The SSIM of our method has improved to 96.34% in the Kvasir-Capsule dataset and 93.34% in the RLE dataset. Besides that, our method also performs the best in the no-reference metric LPIPS. The qualitative results of the comparison methods and our method on the Kvasir-Capsule and RLE datasets are visualized in Fig. 3 with the corresponding heat maps. Firstly, we can see that directly performing LLIE training on DDPM [10] cannot obtain good image restoration, and the original structures of the DDPM images are largely damaged. EnlightenGAN [12] also does not perform satisfactorily in structure restoration. Our method successfully surpasses LLFlow [26] and MIRNetv2 [32] in illumination restoration. The error heat maps further reflect the superior performance of our method in recovering the illumination and structure from low-light images. Moreover, our solution yields the best on the real low-light dataset during the external validation, proving the superior performance of our solution in real-world applications.

**Table 1.** Image quality comparison with existing methods on Kvasir-Capsule [24] and RLE dataset [5]. The 'External Val' denotes the external validation experiment conducted on 100 selected real low-light images from the Kvasir-Capsule dataset [24]. The red lesion segmentation experiment is also conducted on RLE test set [5].

| Models | Kvasir-Capsule | | | RLE | | | External Val | | RLE Segmentation | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR ↑ | SSIM ↑ | LPIPS ↓ | PSNR ↑ | SSIM ↑ | LPIPS ↓ | LPIPS ↓ | PIQE ↓ | mIoU ↑ | Dice ↑ | HD ↓ |
| LIME [9] | 12.07 | 29.66 | 0.4401 | 14.21 | 15.93 | 0.5144 | 0.3498 | 26.41 | 60.19 | 78.42 | 56.20 |
| DUAL [33] | 11.61 | 29.01 | 0.4532 | 14.64 | 16.11 | 0.4903 | 0.3305 | 25.47 | 61.89 | 78.15 | 55.70 |
| Zero-DCE [8] | 14.03 | 46.31 | 0.4917 | 14.86 | 34.18 | 0.4519 | 0.6723 | 21.47 | 54.77 | 71.46 | 56.24 |
| EnlightenGAN [12] | 27.15 | 85.03 | 0.1769 | 23.65 | 80.51 | 0.1864 | 0.4796 | 34.75 | 61.97 | 74.15 | 54.89 |
| LLFlow [26] | 29.69 | 92.57 | 0.0774 | 25.93 | 85.19 | 0.1340 | 0.3712 | 35.67 | 61.06 | **78.55** | 60.04 |
| HWMNet [6] | 27.62 | 92.09 | 0.1507 | 21.81 | 76.11 | 0.3624 | 0.5089 | 35.37 | 56.48 | 74.17 | 59.90 |
| MIRNet [31] | 31.23 | 95.77 | 0.0436 | 25.77 | 86.94 | 0.1519 | 0.3485 | 34.28 | 59.84 | 78.32 | 63.10 |
| StillGAN [19] | 28.28 | 91.30 | 0.1302 | 26.38 | 83.33 | 0.1860 | 0.3095 | 38.10 | 58.32 | 71.56 | 55.02 |
| SNR-Aware [30] | 30.32 | 94.92 | 0.0521 | 27.73 | 88.44 | 0.1094 | 0.3992 | 26.82 | 58.95 | 70.26 | 57.73 |
| MIRNetv2 [32] | 31.67 | 95.22 | 0.0486 | 32.85 | 92.69 | 0.0781 | 0.3341 | 41.24 | 63.14 | 75.07 | 53.71 |
| DDPM [10] | 25.17 | 73.16 | 0.4098 | 22.97 | 70.31 | 0.4198 | 0.5069 | 43.64 | 54.09 | 75.10 | 67.54 |
| **LLCaps (Ours)** | **35.24** | **96.34** | **0.0374** | **33.18** | **93.34** | **0.0721** | **0.3082** | **20.67** | **66.47** | 78.47 | **44.37** |

**Table 2.** Ablation experiments of our LLCaps on the Kvasir-Capsule Dataset [24]. In order to observe the performance changes, we (i) remove the wavelet transform, (ii) degenerate the CurveA layer, and (iii) remove the reverse diffusion branch.

| Wavelet Transform | Curve Attention | Reverse Diffusion | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
|---|---|---|---|---|---|
| ✗ | ✗ | ✗ | 31.12 | 94.96 | 0.0793 |
| ✓ | ✗ | ✗ | 32.78 | 96.26 | 0.0394 |
| ✗ | ✓ | ✗ | 32.08 | 96.27 | 0.0415 |
| ✗ | ✗ | ✓ | 33.10 | 94.53 | 0.0709 |
| ✓ | ✓ | ✗ | 33.92 | 96.20 | 0.0381 |
| ✓ | ✗ | ✓ | 34.07 | 95.61 | 0.0518 |
| ✗ | ✓ | ✓ | 33.41 | 95.03 | 0.0579 |
| ✓ | ✓ | ✓ | **35.24** | **96.34** | **0.0374** |

Furthermore, a downstream red lesion segmentation task is conducted to investigate the usefulness of our LLCaps on clinical applications. As illustrated in Table 3, LLCaps achieve the best lesion segmentation results, manifesting the superior performance of our LLCaps model in lesion segmentation. Additionally, LLCaps surpasses all SOTA methods in HD, showing LLCaps images perform perfectly in processing the segmentation boundaries, suggesting that our method possesses better image reconstruction and edge retention ability.

Besides, an ablation study is conducted on the Kvasir-Capsule dataset to demonstrate the effectiveness of our design and network components, as shown in Table 2. To observe and compare the performance changes, we try to (i) remove the wavelet transform in CWA blocks, (ii) degenerate the curved attention (CurveA) layer in CWA block to a simple channel attention layer [31], and (iii) remove the reverse diffusion branch. Experimental results demonstrate that the absence of any component shall cause great performance degradation. The significant improvement in quantitative metrics is a further testament to the effectiveness of our design for each component.

## 4 Conclusion

We present LLCaps, an end-to-end capsule endoscopy LLIE framework with multi-scale CNN and reverse diffusion process. The CNN branch is constructed by stacked MSRB modules, in which the core CWA block extracts high-frequency detail information through wavelet transform, and learns the local representation of the image via the Curved Attention layer. The reverse diffusion process further optimizes the shallow output, achieving the closest approximation to the real image. Comparison and ablation studies prove that our method and design bring about superior performance improvement in image quality. Further medical image segmentation experiments demonstrate the reliability of our method in clinical applications. Potential future works include extending our model to various medical scenarios (e.g., surgical robotics, endoscopic navigation, augmented reality for surgery) and clinical deep learning model deployment.
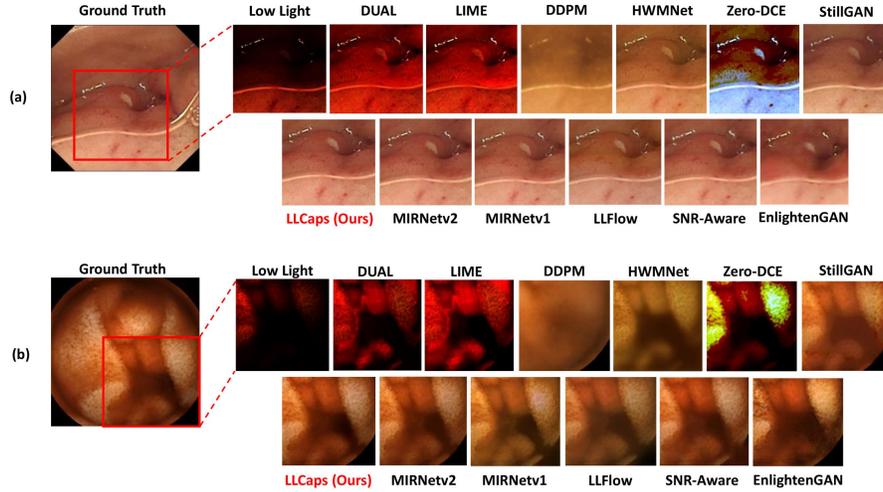
## References

1. Bai, L., Wang, L., Chen, T., Zhao, Y., Ren, H.: Transformer-based disease identification for small-scale imbalanced capsule endoscopy dataset. Electronics **11**(17), 2747 (2022)
2. Charbonnier, P., Blanc-Feraud, L., Aubert, G., Barlaud, M.: Two deterministic half-quadratic regularization algorithms for computed imaging. In: Proceedings of 1st international conference on image processing. vol. 2, pp. 168–172. IEEE (1994)
3. Che, H., Chen, S., Chen, H.: Image quality-aware diagnosis via meta-knowledge co-embedding. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 19819–19829 (2023)
4. Chen, W., Liu, Y., Hu, J., Yuan, Y.: Dynamic depth-aware network for endoscopy super-resolution. IEEE Journal of Biomedical and Health Informatics **26**(10), 5189–5200 (2022)
5. Coelho, P., Pereira, A., Leite, A., Salgado, M., Cunha, A.: A deep learning approach for red lesions detection in video capsule endoscopies. In: Image Analysis and Recognition: 15th International Conference, ICIAR 2018, Póvoa de Varzim, Portugal, June 27–29, 2018, Proceedings 15. pp. 553–561. Springer (2018)
6. Fan, C.M., Liu, T.J., Liu, K.H.: Half wavelet attention on m-net+ for low-light image enhancement. In: 2022 IEEE International Conference on Image Processing (ICIP). pp. 3878–3882. IEEE (2022)
7. Gómez, P., Semmler, M., Schützenberger, A., Bohr, C., Döllinger, M.: Low-light image enhancement of high-speed endoscopic videos using a convolutional neural network. Medical & biological engineering & computing **57**, 1451–1463 (2019)
8. Guo, C., Li, C., Guo, J., Loy, C.C., Hou, J., Kwong, S., Cong, R.: Zero-reference deep curve estimation for low-light image enhancement. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 1780–1789 (2020)

9. Guo, X., Li, Y., Ling, H.: Lime: Low-light image enhancement via illumination map estimation. IEEE transactions on image processing **26**(2), 982–993 (2016)

10. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. Advances in Neural Information Processing Systems **33**, 6840–6851 (2020)

11. Huynh-Thu, Q., Ghanbari, M.: Scope of validity of psnr in image/video quality assessment. Electronics letters **44**(13), 800–801 (2008)

12. Jiang, Y., Gong, X., Liu, D., Cheng, Y., Fang, C., Shen, X., Yang, J., Zhou, P., Wang, Z.: Enlightengan: Deep light enhancement without paired supervision. IEEE transactions on image processing **30**, 2340–2349 (2021)

13. Li, C., Guo, C., Han, L., Jiang, J., Cheng, M.M., Gu, J., Loy, C.C.: Low-light image and video enhancement using deep learning: A survey. IEEE transactions on pattern analysis and machine intelligence **44**(12), 9396–9416 (2021)

14. Li, J., Fang, F., Mei, K., Zhang, G.: Multi-scale residual network for image super-resolution. In: Proceedings of the European conference on computer vision (ECCV). pp. 517–532 (2018)

15. Li, M., Liu, J., Yang, W., Sun, X., Guo, Z.: Structure-revealing low-light image enhancement via robust retinex model. IEEE Transactions on Image Processing **27**(6), 2828–2841 (2018)

16. Liu, Y.F., Guo, J.M., Yu, J.C.: Contrast enhancement using stratified parametric-oriented histogram equalization. IEEE Transactions on Circuits and Systems for Video Technology **27**(6), 1171–1181 (2016)

17. Long, M., Li, Z., Xie, X., Li, G., Wang, Z.: Adaptive image enhancement based on guide image and fraction-power transformation for wireless capsule endoscopy. IEEE transactions on biomedical circuits and systems **12**(5), 993–1003 (2018)

18. Lore, K.G., Akintayo, A., Sarkar, S.: Llnet: A deep autoencoder approach to natural low-light image enhancement. Pattern Recognition **61**, 650–662 (2017)

19. Ma, Y., Liu, J., Liu, Y., Fu, H., Hu, Y., Cheng, J., Qi, H., Wu, Y., Zhang, J., Zhao, Y.: Structure and illumination constrained gan for medical image enhancement. IEEE Transactions on Medical Imaging **40**(12), 3955–3967 (2021)

20. Ma, Y., Liu, Y., Cheng, J., Zheng, Y., Ghahremani, M., Chen, H., Liu, J., Zhao, Y.: Cycle structure and illumination constrained gan for medical image enhancement. In: Medical Image Computing and Computer Assisted Intervention–MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part II 23. pp. 667–677. Springer (2020)

21. Pandey, K., Mukherjee, A., Rai, P., Kumar, A.: Diffusevae: Efficient, controllable and high-fidelity generation from low-dimensional latents. arXiv preprint arXiv:2201.00308 (2022)

22. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18. pp. 234–241. Springer (2015)

23. Sliker, L.J., Ciuti, G.: Flexible and capsule endoscopy for screening, diagnosis and treatment. Expert review of medical devices **11**(6), 649–666 (2014)

24. Smedsrud, P.H., Thambawita, V., Hicks, S.A., Gjestang, H., Nedrejord, O.O., Næss, E., Borgli, H., Jha, D., Berstad, T.J.D., Eskeland, S.L., et al.: Kvasir-capsule, a video capsule endoscopy dataset. Scientific Data **8**(1), 142 (2021)

25. Venkatanath, N., Praneeth, D., Bh, M.C., Channappayya, S.S., Medasani, S.S.: Blind image quality evaluation using perception based features. In: 2015 twenty first national conference on communications (NCC). pp. 1–6. IEEE (2015)

26. Wang, Y., Wan, R., Yang, W., Li, H., Chau, L.P., Kot, A.: Low-light image enhancement with normalizing flow. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 36, pp. 2604–2612 (2022)
27. Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P.: Image quality assessment: from error visibility to structural similarity. IEEE transactions on image processing **13**(4), 600–612 (2004)
28. Wu, J., Fu, R., Fang, H., Zhang, Y., Xu, Y.: Medsegdiff-v2: Diffusion based medical image segmentation with transformer. arXiv preprint arXiv:2301.11798 (2023)
29. Xiao, T., Singh, M., Mintun, E., Darrell, T., Dollár, P., Girshick, R.: Early convolutions help transformers see better. Advances in Neural Information Processing Systems **34**, 30392–30400 (2021)
30. Xu, X., Wang, R., Fu, C.W., Jia, J.: Snr-aware low-light image enhancement. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 17714–17724 (2022)
31. Zamir, S.W., Arora, A., Khan, S., Hayat, M., Khan, F.S., Yang, M.H., Shao, L.: Learning enriched features for real image restoration and enhancement. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXV 16. pp. 492–511. Springer (2020)
32. Zamir, S.W., Arora, A., Khan, S., Hayat, M., Khan, F.S., Yang, M.H., Shao, L.: Learning enriched features for fast image restoration and enhancement. IEEE transactions on pattern analysis and machine intelligence **45**(2), 1934–1948 (2022)
33. Zhang, Q., Nie, Y., Zheng, W.S.: Dual illumination estimation for robust exposure correction. In: Computer Graphics Forum. vol. 38, pp. 243–252 (2019)
34. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 586–595 (2018)
35. Zhang, Y., Bai, L., Liu, L., Ren, H., Meng, M.Q.H.: Deep reinforcement learning-based control for stomach coverage scanning of wireless capsule endoscopy. In: 2022 IEEE International Conference on Robotics and Biomimetics (ROBIO). pp. 01–06. IEEE (2022)
36. Zhou, S., Li, C., Change Loy, C.: Lednet: Joint low-light enhancement and deblurring in the dark. In: Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings. pp. 573–589. Springer (2022)
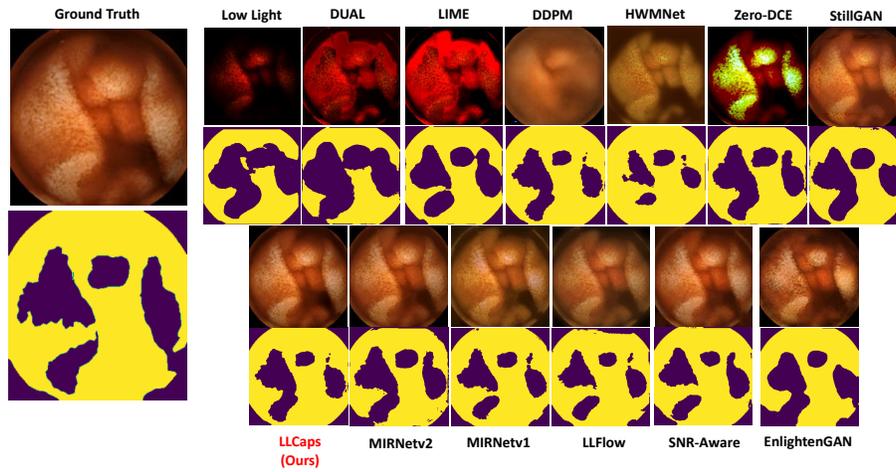
# Supplementary Materials for "LLCaps: Learning to Illuminate Low-Light Capsule Endoscopy with Curved Wavelet Attention and Reverse Diffusion"



**Fig. 4.** The full visualization results for LLCaps compared with SOTA approaches on (a) Kvasir-Capsule dataset and (b) RLE dataset.

**Table 3.** Ablation study of the frequency analysis on the CWA block using average gradient. Higher mean and variance of AG denote richer details. When removing the wavelet transform or the CWA block, the mean and var drop greatly, showing the effectiveness of our proposed CWA block in extracting features.

| Average Gradient | Mean $\times 10^5$ ↑ | Variance $\times 10^{10}$ ↑ |
|---|---|---|
| DDPM [9] | 2.19 | 0.50 |
| MIRNetv2 [30] | 3.51 | 1.86 |
| LLCaps w/o CWA | 3.54 | 0.99 |
| LLCaps w/o Curved Attention | 3.85 | 2.01 |
| LLCaps w/o Wavelet Transform | 3.55 | 1.93 |
| LLCaps | 3.87 | 2.09 |
| GT | 3.95 | 2.19 |

**Fig. 5.** Visualization of the red lesion segmentation comparison experiments on the RLE dataset.