# LYT-NET: Lightweight YUV Transformer-based Network for Low-light Image Enhancement

Alexandru Brateanu*, Raul Balmez*, Adrian Avram†, Ciprian Orhei† and Cosmin Ancuti†

*University of Manchester, United Kingdom

†ETcTI, Universitatea Politehnica Timisoara, Romania

*Abstract*—**This letter introduces LYT-Net, a novel lightweight transformer-based model for low-light image enhancement. LYT-Net consists of several layers and detachable blocks, including our novel blocks—Channel-Wise Denoiser (CWD) and Multi-Stage Squeeze & Excite Fusion (MSEF)—along with the traditional Transformer block, Multi-Headed Self-Attention (MHSA). In our method we adopt a dual-path approach, treating chrominance channels $U$ and $V$ and luminance channel $Y$ as separate entities to help the model better handle illumination adjustment and corruption restoration. Our comprehensive evaluation on established LLIE datasets demonstrates that, despite its low complexity, our model outperforms recent LLIE methods. The source code and pre-trained models are available at https://github.com/albrateanu/LYT-Net**

*Index Terms*—**Low-light Image Enhancement, Vision Transformer, Deep Learning**

## I. INTRODUCTION

Low-light image enhancement (LLIE) is an important and challenging task in computational imaging. When images are captured in low-light conditions, their quality often deteriorates, leading to a loss of detail and contrast. This not only makes the images visually unappealing but also affects the performance of many imaging systems. The goal of LLIE is to improve the clarity and contrast of these images, while also correcting distortions that commonly occur in dark environments, all without introducing unwanted artifacts or causing imbalances in color.

Earlier LLIE methods [1] primarily relied on frequency decomposition [2], [3], [4], histogram equalization [5], [6], [7], and Retinex theory [8], [9], [10], [11], [12]. With the rapid advancement of deep learning, various CNN architectures [13], [14], [15], [16], [17], [18], [19], [20], [21], [22] have been shown to outperform traditional LLIE techniques. Based on Retinex theory, Retinex-Net [13] integrates Retinex decomposition with an original CNN architecture, while Diff-Retinex [16] proposes a generative framework to further address content loss and color deviation caused by low light.

The development of Generative Adversarial Networks (GAN) [23] has provided a new perspective for LLIE, where low-light images are used as input to generate their normal-light counterparts. For instance, EnlightenGAN [24] employs a single generator model to directly convert low-light images to normal-light versions, effectively using both global and local discriminators in the transformation process.

More recently, Vision Transformers (ViTs) [25] have demonstrated significant effectiveness in various computer vision tasks [26], [27], [28], [29], [30], largely due to the
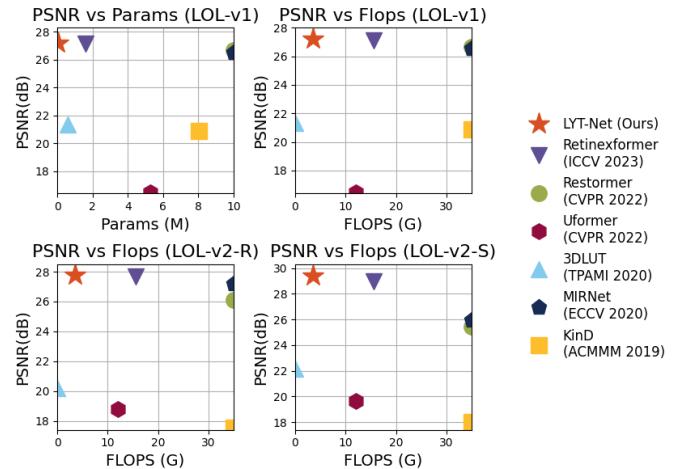


Fig. 1. Our model (LYT-Net) delivers SOTA performance in LLIE task, while maintaining the best computational efficiency (results are plotted on LOL dataset [13]).

self-attention (SA) mechanism. Despite these advancements, the application of ViTs to low-level vision tasks remains relatively underexplored. Only a few LLIE (Low-Light Image Enhancement)-ViT-based strategies have been introduced in the recent literature [31], [32], [33]. Uformer [31] is based on the classical UNet architecture, where the convolution layers are replaced with Transformer blocks while maintaining the hierarchical encoder-decoder structure and skip connections. Restormer [33], on the other hand, introduces a multi-Dconv head transposed attention (MDTA) block, replacing the vanilla multi-head self-attention.

In this letter, we propose a novel lightweight transformer-based approach called LYT-Net. Different from the existing transformer-based methods, our method focuses on computational efficiency while still producing state-of-the-art (SOTA) results. Specifically, we first separate chrominance from luminance employing the YUV color space. The chrominance information (channels $U$ and $V$) is initially processed through a specialized Channel-wise Denoiser (**CWD**) block, which reduces noise while preserving fine details. To minimize computational complexity, the luminance channel $Y$ undergoes convolution and pooling to extract features, which are subsequently enhanced by a traditional Multi-headed Self-Attention (**MHSA**) block. The enhanced channels are then recombined and processed through a novel Multi-stage Squeeze and Excite Fusion (**MSEF**) block. Finally, the chrominance channels $U$ and $V$ channels are concatenated with the luminance $Y$
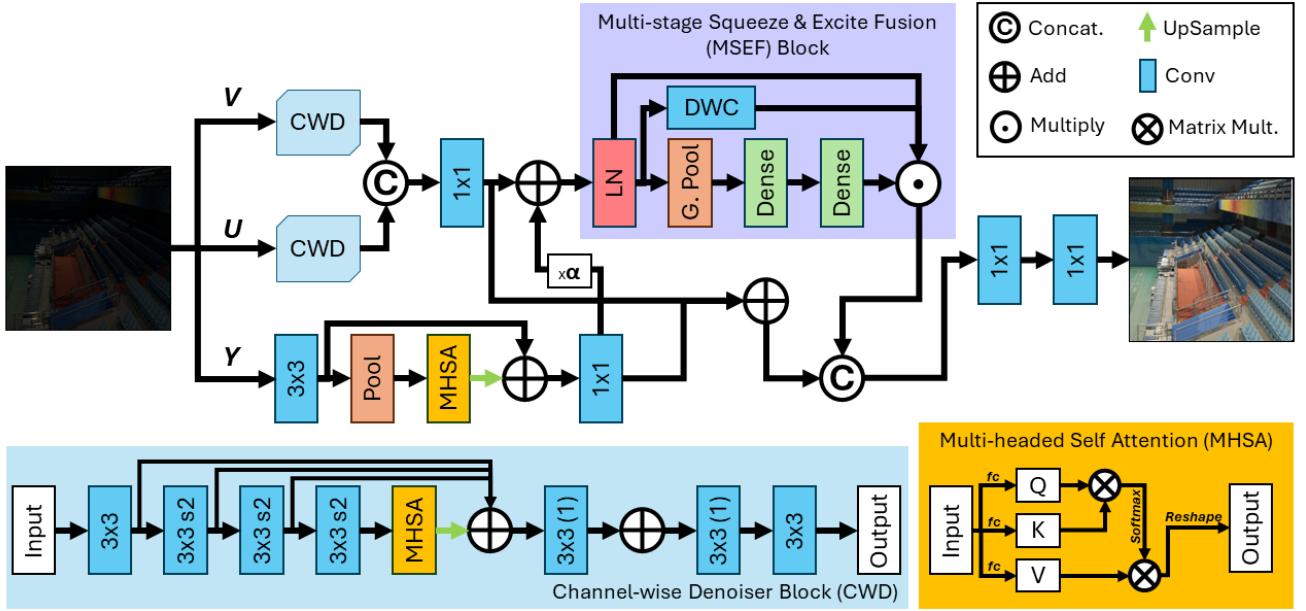
Fig. 2. Overall framework of LYT-Net. The architecture consists of several detachable blocks like Channel-wise Denoiser (CWD), Multi-headed Self-Attention (MHSA), Multi-stage Squeeze and Excite Fusion (MSEF).

channel and passed through a final set of convolutional layers to produce the restored image.

Our method has undergone extensive testing on established LLIE datasets. Both qualitative and quantitative evaluations indicate that our approach achieves highly competitive results. Fig. 1 presents a comparative analysis of performance over complexity between SOTA methods evaluated using the LOL dataset [13]. It can be observed that, despite its lightweight design, our method produces results that are not only comparable to, but often outperform, those of more complex recent deep learning LLIE techniques.

## II. OUR APPROACH

In Fig. 2, we illustrate the overall architecture of LYT-Net, which consists of several layers and detachable blocks, including our novel blocks—Channel-Wise Denoiser (**CWD**) and Multi-Stage Squeeze & Excite Fusion (**MSEF**)—along with the traditional ViT block, Multi-Headed Self-Attention (**MHSA**). We adopt a dual-path approach, treating chrominance and luminance as separate entities to help the model better handle illumination adjustment and corruption restoration. The luminance channel $Y$ undergoes convolution and pooling to extract features, which are then enhanced by the **MHSA** block. Chrominance channels $U$ and $V$ are processed through the **CWD** block to reduce noise while preserving details. The enhanced chrominance channels are then recombined and processed through the **MSEF** block. Finally, the chrominance $U, V$ and luminance $Y$ channels are concatenated and passed through a final set of convolutional layers to produce the output, resulting in a high-quality, enhanced image.

### A. Channel-wise Denoiser Block

The **CWD** Block employs a U-shaped network with **MHSA** as the bottleneck, integrating convolutional and attention-based mechanisms. It includes multiple *conv3×3* layers with varying strides and skip connections, facilitating detailed feature capture and denoising.

It consists of a series of four *conv3×3* layers. The first *conv3×3* has strides of 1 for feature extraction. The other three *conv3×3* layers have strides of 2, helping with capturing features at different scales. The integration of the attention bottleneck enables the model to capture long-range dependencies, followed by upsampling layers and skip connections to reconstruct and facilitate the recovery of spatial resolution.

This approach allows us to apply MHSA on a feature map with reduced spatial dimensions, significantly improving computational efficiency. Additionally, using interpolation-based upsampling instead of transposed convolutions cuts the number of parameters in the CWD by more than half, while preserving performance.

### B. Multi-headed Self-attention Block

In our updated simplified transformer architecture, the input feature $\mathbf{F}_{in} \in \mathbb{R}^{H \times W \times C}$ is first linearly projected into query ($\mathbf{Q}$), key ($\mathbf{K}$), and value ($\mathbf{V}$) components through bias-free fully connected layers. The linear projections maintain the original input dimensionality.

$$\mathbf{Q} = \mathbf{X}\mathbf{W_Q}, \mathbf{K} = \mathbf{X}\mathbf{W_K}, \mathbf{V} = \mathbf{X}\mathbf{W_V}, \ \mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{HW \times C} \quad (1)$$

Next, these projected features are split into $k$ heads:

$$\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2, \cdots, \mathbf{X}_k], \ \mathbf{X}_i \in \mathbb{R}^{HW \times d_k}, \ d_k = \frac{C}{k}, i = \overline{1,k} \quad (2)$$

where each head operates independently with dimensionality $d_k$. The self-attention mechanism is applied to each head, as defined below:

$$\text{Attention}(\mathbf{Q}_i, \mathbf{K}_i, \mathbf{V}_i) = \text{Softmax}\left(\frac{\mathbf{Q}_i\mathbf{K}_i^{\mathsf{T}}}{\sqrt{d_k}}\right) \times \mathbf{V}_i \quad (3)$$
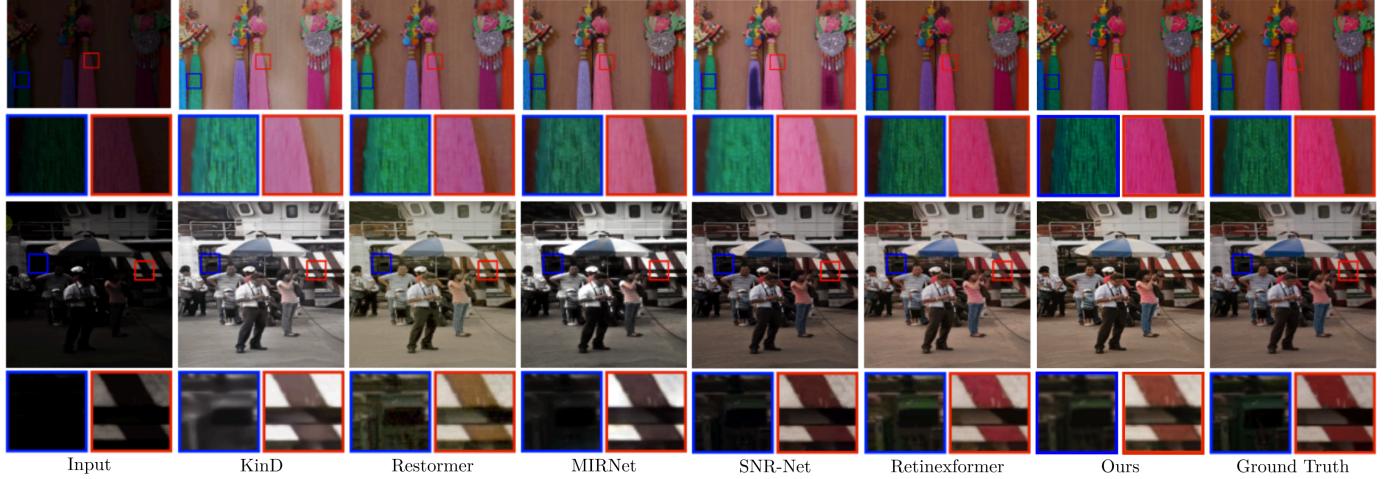
Fig. 3. Qualitative comparison with SOTA LLIE methods on the LOL dataset. Zoom-in regions are used to illustrate differences.

Finally, the attention outputs from all heads are concatenated and the combined output is passed through a linear layer to project it back to the original embedding size. The output tokens $\mathbf{X}_{\text{out}}$ are reshaped back into the original spatial dimensions to form the output feature $\mathbf{F}_{\text{out}} \in \mathbb{R}^{H \times W \times C}$.

### C. Multi-stage Squeeze & Excite Fusion Block

The **MSEF** Block enhances both spatial and channel-wise features of $\mathbf{F}_{\text{in}}$. Initially, $\mathbf{F}_{\text{in}}$ undergoes layer normalization, followed by global average pooling to capture global spatial context and a reduced fully-connected layer with ReLU activation, producing a reduced descriptor $\mathbf{S}_{\text{reduced}}$, as shown in Eq. (4). This descriptor is then expanded back to the original dimensions through another fully-connected layer with Tanh activation, resulting in $\mathbf{S}_{\text{expanded}}$, Eq. (5).

These operations compress the feature map into a reduced descriptor (the **squeezing** operation) to capture essential details, and then re-expand it (the **excitation** operation) to restore the full dimensions while emphasizing the most relevant features.

$$\mathbf{S}_{\text{reduced}} = \text{ReLU}(\mathbf{W}_1 \cdot \text{GlobalPool}(\text{LayerNorm}(\mathbf{F}_{\text{in}}))) \quad (4)$$

$$\mathbf{S}_{\text{expanded}} = \text{Tanh}(\mathbf{W}_2 \cdot \mathbf{S}_{\text{reduced}}) \cdot \text{LayerNorm}(\mathbf{F}_{\text{in}}) \quad (5)$$

A residual connection is added to the fused output to produce the final output feature map $\mathbf{F}_{\text{out}}$, as in Eq. (6).

$$\mathbf{F}_{\text{out}} = \text{DWConv}(\text{LayerNorm}(\mathbf{F}_{\text{in}})) \cdot \mathbf{S}_{\text{expanded}} + \mathbf{F}_{\text{in}} \quad (6)$$

### D. Loss Function

In our approach, a hybrid loss function plays a pivotal role in training our model effectively. The hybrid loss $\mathbf{L}$ is formulated as in Eq. (7), where $\alpha_1$ to $\alpha_5$ are hyperparameters used to balance each constituent loss function.

$$\mathbf{L} = \mathbf{L}_{\text{S}} + \alpha_1 \mathbf{L}_{\text{Perc}} + \alpha_2 \mathbf{L}_{\text{Hist}} + \alpha_3 \mathbf{L}_{\text{PSNR}} + \alpha_4 \mathbf{L}_{\text{Color}} + \alpha_5 \mathbf{L}_{\text{MS-SSIM}}$$
$$(7)$$

The hybrid loss in our model combines several components to enhance image quality and perception. Smooth L1 loss $\mathbf{L}_{\text{S}}$

handles outliers by applying a quadratic or linear penalty based on the difference between predicted and true values. Perceptual loss $\mathbf{L}_{\text{Perc}}$ maintains feature consistency by comparing VGG-extracted feature maps. Histogram loss $\mathbf{L}_{\text{Hist}}$ aligns pixel intensity distributions between predicted and true images. PSNR loss $\mathbf{L}_{\text{PSNR}}$ reduces noise by penalizing mean squared error, while Color loss $\mathbf{L}_{\text{Color}}$ ensures color fidelity by minimizing differences in channel mean values. Lastly, Multiscale SSIM loss $\mathbf{L}_{\text{MS-SSIM}}$ preserves structural integrity by evaluating similarity across multiple scales. Together, these losses form a comprehensive strategy addressing various aspects of image enhancement.

## III. RESULTS AND DISCUSSION

**Implementation details:** The implementation of LYT-Net utilizes the TensorFlow framework. The ADAM Optimizer ($\beta_1 = 0.9$ and $\beta_2 = 0.999$) is employed for training over 1000 epochs. The initial learning rate is set to $2 \times 10^{-4}$ and gradually decays to $1 \times 10^{-6}$ following a cosine annealing schedule, aiding in optimization convergence and avoiding local minima. The hyperparameters of the hybrid loss function are set as: $\alpha_1$=0.06, $\alpha_2$=0.05, $\alpha_3$=0.5, $\alpha_4$=0.0083, and $\alpha_5$=0.25.

LYT-Net is trained and evaluated on three versions of the LOL dataset: LOL-v1, LOL-v2-real, and LOL-v2-synthetic. The corresponding training/testing splits are $485 : 15$ for LOL-v1, $689 : 100$ for LOL-v2-real, and $900 : 100$ for LOL-v2-synthetic.

During training, image pairs undergo random augmentations, including random cropping to $256 \times 256$ and random flipping/rotation, to prevent overfitting. The training is conducted with a batch size of 1. Evaluation metrics include PSNR and SSIM for performance assessment.

**Quantitative results:** The proposed method is compared to SOTA LLIE techniques, as shown in Table I, focusing on two aspects: quantitative performance on the LOL datasets (LOL-v1, LOL-v2-real, LOL-v2-synthetic) and model complexity.

As shown in Table I, LYT-Net consistently outperforms the current SOTA methods across all versions of the LOL dataset in terms of both PSNR and SSIM. Additionally, LYT-Net is highly efficient, requiring only 3.49G FLOPS and

| Methods | Complexity | | LOL-v1 | | LOL-v2-real | | LOL-v2-syn | |
|---|---|---|---|---|---|---|---|---|
| | FLOPS (G) | Params (M) | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| SID [34] | 13.73 | 7.76 | 14.35 | 0.436 | 13.24 | 0.442 | 15.04 | 0.610 |
| 3DLUT [35] | 0.075 | 0.59 | 21.35 | 0.585 | 20.19 | 0.745 | 22.17 | 0.854 |
| DeepUPE [15] | 21.10 | 1.02 | 14.38 | 0.446 | 13.27 | 0.452 | 15.08 | 0.623 |
| DeepLPF [14] | 5.86 | 1.77 | 15.28 | 0.473 | 14.10 | 0.480 | 16.02 | 0.587 |
| UFormer [31] | 12.00 | 5.29 | 16.36 | 0.771 | 18.82 | 0.771 | 19.66 | 0.871 |
| RetinexNet [13] | 587.47 | 0.84 | 18.92 | 0.427 | 18.32 | 0.447 | 19.09 | 0.774 |
| Sparse [36] | 53.26 | 2.33 | 17.20 | 0.640 | 20.06 | 0.816 | 22.05 | 0.905 |
| EnGAN [24] | 61.01 | 114.35 | 20.00 | 0.691 | 18.23 | 0.617 | 16.57 | 0.734 |
| FIDE [37] | 28.51 | 8.62 | 18.27 | 0.665 | 16.85 | 0.678 | 15.20 | 0.612 |
| KinD [17] | 34.99 | 8.02 | 20.86 | 0.790 | 14.74 | 0.641 | 13.29 | 0.578 |
| Restormer [33] | 144.25 | 26.13 | 26.68 | 0.853 | 26.12 | 0.853 | 25.43 | 0.859 |
| MIRNet [20] | 785 | 31.76 | 26.52 | 0.856 | 27.17 | 0.865 | 25.96 | 0.898 |
| SNR-Net [22] | 26.35 | 4.01 | 26.72 | 0.851 | 27.21 | 0.871 | 27.79 | 0.941 |
| Retinexformer [32] | 15.57 | 1.61 | 27.14 | 0.850 | 27.69 | 0.856 | 28.99 | 0.939 |
| LYT-Net | 3.49 | 0.045 | 27.23 | 0.853 | 27.80 | 0.873 | 29.38 | 0.940 |

TABLE I
QUANTITATIVE RESULTS ON LOL DATASETS. BEST RESULTS ARE IN RED, SECOND BEST ARE IN BLUE.
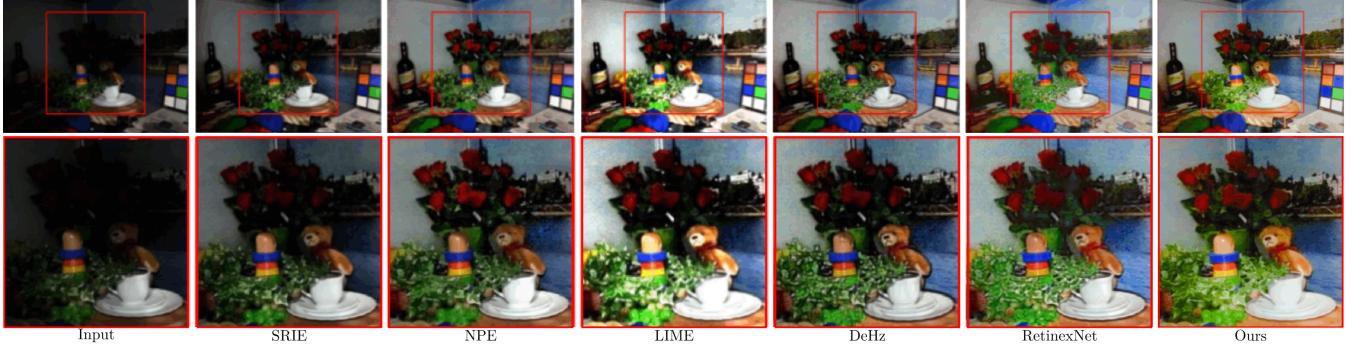


Fig. 4. Qualitative comparison with SOTA LLIE methods on LIME dataset. Zoom-in regions are used to illustrate differences.

utilizing just 0.045M parameters, which gives it a significant advantage over other SOTA methods that are generally much more complex. The only exception is 3DLUT[35], which is comparable to our approach in terms of complexity. However, LYT-Net clearly surpasses the 3DLUT method in both PSNR and SSIM. This combination of strong performance and low complexity highlights the overall effectiveness of LYT-Net.

**Qualitative Results:** The qualitative evaluation of LYT-Net against SOTA LLIE techniques is shown in Fig. 3 on the LOL dataset and in Fig. 4 on LIME [38].

Previous methods, such as KiND[17] and Restormer[33], exhibit color distortion issues, as shown in Fig. 3. Additionally, several algorithms (e.g. MIRNet[20], and SNR-Net[22]) tend to produce over- or under-exposed areas, compromising image contrast while enhancing luminance. Similarly, Fig. 4 demonstrates that SRIE [39], DeHz [40], and NPE [41] result in a loss of contrast. In general, our LYT-Net is highly effective at improving visibility and enhancing low-contrast or poorly lit areas, while efficiently eliminating noise without introducing spots or artifacts.

## IV. ABLATION STUDY

.
The ablation study is conducted on the LOLv1 dataset, using PSNR as the quantitative metric, and evaluates the impact of the **CWD** and **MSEF** blocks. In the YUV decomposition, applying **CWD** to the *Y*-channel (used as the illumination map) results in the retention of lighting artifacts, leading to performance degradation compared to pooling operations and interpolation-based upsampling, which smooth the illumination for better outcomes. However, **CWD** enhances the

| Y-CWD | UV-CWD | MSEF | Params | PSNR |
|---|---|---|---|---|
| ✓ | | | 40238 | 26.62 |
| | ✓ | | 44377 | 26.99 |
| ✓ | ✓ | | 48516 | 26.76 |
| ✓ | | ✓ | 40784 | 26.78 |
| | ✓ | ✓ | 44923 | 27.23 |
| ✓ | ✓ | ✓ | 49062 | 27.02 |

TABLE II
ABLATION STUDY: PERFORMANCE AND PARAMETER IMPACT OF CWD AND MSEF BLOCKS.

chrominance channels (*U* and *V*), preserving detail without introducing noise. Moreover, the **MSEF** block consistently boosts performance across all **CWD** combinations, improving PSNR by 0.16, 0.24, and 0.26 dB, respectively, while increasing the parameter count by only 546.

## V. CONCLUSIONS

We introduce LYT-Net, an innovative lightweight transformer-based model for enhancing low-light images. Our approach utilizes a dual-path framework, processing chrominance and luminance separately to improve the model's ability to manage illumination adjustments and restore corrupted regions. LYT-Net integrates multiple layers and modular blocks, including two unique components—Channel-Wise Denoiser (CWD) and Multi-Stage Squeeze & Excite Fusion (MSEF)—as well as the traditional Vision Transformer (ViT) block with Multi-Headed Self-Attention (MHSA). A comprehensive qualitative and quantitative analysis demonstrates that LYT-Net consistently outperforms SOTA methods on all versions of the LOL dataset in terms of PSNR and SSIM, while maintaining high computational efficiency.

## REFERENCES

[1] W. Wang, X. Wu, X. Yuan, and Z. Gao, "An experiment-based review of low-light image enhancement methods," *IEEE Access*, vol. 8, pp. 87 884–87 917, 2020.

[2] L. Xiao, C. Li, Z. Wu, and T. Wang, "An enhancement method for x-ray image via fuzzy noise removal and homomorphic filtering," *Neurocomputing*, vol. 195, 2016.

[3] A. Loza, D. Bull, and A. Achim, "Automatic contrast enhancement of low-light images based on local statistics of wavelet coefficients," *IEEE Int. Conf. on Image Processing*, 2013.

[4] S. E. Kim, J. J. Jeon, and I. K. Eom, "Image contrast enhancement using entropy scaling in wavelet domain," *Signal Processing*, vol. 127(1), 2016.

[5] T. Arici, S. Dikbas, and Y. Altunbasak, "A histogram modification framework and its application for image contrast enhancement," *IEEE Transactions on Image Processing*, vol. 18(9), 2009.

[6] S.-D. Chen and A. R. Ramli, "Contrast enhancement using recursive mean-separate histogram equalization for scalable brightness preservation," *IEEE Transactions on Consumer Electronics*, vol. 49(4), 2003.

[7] S. Kansal, S. Purwar, and R. K. Tripathi, "Image contrast enhancement using unsharp masking and histogram equalization," *Multimedia Tools Applications*, vol. 77(20), 2018.

[8] E. H. Land, "The retinex theory of color vision," *Scientific american*, vol. 237, no. 6, pp. 108–129, 1977.

[9] X. Fu, Y. Liao, D. Zeng, Y. Huang, X.-P. Zhang, and X. Ding, "A probabilistic method for image enhancement with simultaneous illumination and reflectance estimation," *IEEE Transactions on Image Processing*, vol. 24(12), 2015.

[10] S. Park, S.Yu, B.Moon, S.Ko, and J. Paik, "Low-light image enhancement using variational optimization-based retinex model," *IEEE Transactions on Consumer Electronics*, vol. 63(2), 2017.

[11] Z. Gu, F. Li, F. Fang, and G. Zhang, "A novel retinex-based fractional order variational model for images with severely low light," *IEEE Transactions on Image Processing*, vol. 29, 2019.

[12] J. H. Jang, Y. Bae, and J. B. Ra, "Contrast-enhanced fusion of multisensory images using subband-decomposed multiscale retinex," *IEEE Transactions on Image Processing*, vol. 21(8), 2012.

[13] 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.

[14] S. Moran, P. Marza, S. McDonagh, S. Parisot, and G. Slabaugh, "DeepLPF: Deep local parametric filters for image enhancement," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

[15] R. Wang, Q. Zhang, C.-W. Fu, X. Shen, W.-S. Zheng, and J. Jia, "Underexposed photo enhancement using deep illumination estimation," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

[16] X. Yi, H. Xu, H. Zhang, L. Tang, and J. Ma, "Diff-retinex: Rethinking low-light image enhancement with a generative diffusion model," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.

[17] Y. Zhang, J. Zhang, and X. Guo, "Kindling the darkness: A practical low-light image enhancer," in *Proceedings of ACM international conference on multimedia*, 2019.

[18] Y. Zhang, Y. Tian, Y. Kong, B. Zhong, and Y. Fu, "Residual dense network for image restoration," *In IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.

[19] A. Dudhane, S. Zamir, S. Khan, F. Khan, and M.-H. Yang, "Burst image restoration and enhancement," *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.

[20] S. W. Zamir, A. Arora, S. Khan, M. Hayat, F. S. Khan, M.-H. Yang, and L. Shao, "Learning enriched features for real image restoration and enhancement," in *European Conference on Computer Vision*, 2020.

[21] 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 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.

[22] X. Xu, R. Wang, C.-W. Fu, and J. Jia, "SNR-aware low-light image enhancement," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.

[23] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," *In Advances in neural information processing systems*, 2014.

[24] 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.

[25] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, and S. Gelly, "An image is worth 16x16 words: Transformers for image recognition at scale," *International Conference on Learning Representations (ICLR)*, 2021.

[26] L. Yuan, Y. Chen, TaoWang, W. Yu, Y. Shi, Z. Jiang, F. E. Tay, J. Feng, , and S. Yan, "Tokens-to-token vit: Training vision transformers from scratch on imagenet," *IEEE International Conference on Computer Vision (ICCV)*, 2021.

[27] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, T. Lu, P. Luo, and L. Shao, "Pyramid vision transformer: A versatile backbone for dense prediction without convolutions," *IEEE International Conference on Computer Vision (ICCV)*, 2021.

[28] S. Zheng, J. Lu, H. Zhao, X. Zhu, Z. Luo, Y. Wang, Y. Fu, J. Feng, T. Xiang, and P. H. Torr, "Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers," *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.

[29] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "Hierarchical vision transformer using shifted windows," *In European Conference on Computer Vision (ECCV)*, 2020.

[30] Z. Liu, Y. L. Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," *IEEE International Conference on Computer Vision (ICCV)*, 2021.

[31] Z. Wang, X. Cun, J. Bao, W. Zhou, J. Liu, and H. Li, "Uformer: A general u-shaped transformer for image restoration," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.

[32] Y. Cai, H. Bian, J. Lin, H. Wang, R. Timofte, and Y. Zhang, "Retinexformer: One-stage retinex-based transformer for low-light image enhancement," in *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023.

[33] S. W. Zamir, A. Arora, S. Khan, M. Hayat, F. S. Khan, and M.-H. Yang, "Restormer: Efficient transformer for high-resolution image restoration," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.

[34] C. Chen, Q. Chen, M. N. Do, and V. Koltun, "Seeing motion in the dark," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

[35] H. Zeng, J. Cai, L. Li, Z. Cao, and L. Zhang, "Learning image-adaptive 3d lookup tables for high performance photo enhancement in realtime," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 4, pp. 2058–2073, 2020.

[36] W. Yang, W. Wang, H. Huang, S. Wang, and J. Liu, "Sparse gradient regularized deep retinex network for robust low-light image enhancement," *IEEE Transactions on Image Processing*, vol. 30, pp. 2072–2086, 2021.

[37] K. Xu, X. Yang, B. Yin, and R. W. Lau, "Learning to restore low-light images via decomposition-and-enhancement," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

[38] 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.

[39] X. Fu, D. Zeng, Y. Huang, Y. Liao, X. Ding, and J. Paisley, "A fusion-based enhancing method for weakly illuminated images," *Signal Processing*, vol. 129, pp. 82–96, 2016.

[40] X. Dong, Y. Pang, and J. Wen, "Fast efficient algorithm for enhancement of low lighting video," in *ACM SIGGRAPH 2010 Posters*, 2010, pp. 1–1.

[41] S. Wang, J. Zheng, H.-M. Hu, and B. Li, "Naturalness preserved enhancement algorithm for non-uniform illumination images," *IEEE Transactions on Image Processing*, vol. 22, no. 9, pp. 3538–3548, 2013.