

Residual Dilated U-Net with Spatially Adaptive Normalization for the Restoration of Under Display Camera Images

Youngjin Oh*, Gu Yong Park*, Haesoo Chung*, Sunwoo Cho*, and Nam Ik Cho*
 * Department of Electrical & Computer Eng., INMC, Seoul National University, Seoul Korea
 E-mail: {yjymoh0211, benkay, reneeish, etoo33}@ispl.snu.ac.kr, nicho@snu.ac.kr

Abstract—Recent developments in display and imaging devices have prompted manufacturers to place a camera under the display screen, delivering a larger display-to-body ratio that is genuinely full-screen. However, this Under-Display Camera (UDC) imaging system suffers from severe image degradation such as noise, blur, low-light, and color-shift. This is due to the low light transmission and diffraction property of the display panels. To tackle this issue, we present an end-to-end framework based on U-Net to restore the degraded image. Since the point spread function (PSF) of the UDC degradation is known to be spatially dispersed, we utilize dilated convolutions to increase the receptive field of the model. Furthermore, we use spatially adaptive normalization to regularize feature maps to help restore the image efficiently, thereby improving the performance of our model. We show that our method can restore UDC images with fewer artifacts and produce competitive results to state-of-the-art methods.

I. INTRODUCTION

Under-Display Camera (UDC) [1] is an imaging system that places cameras under the display panel, unlike conventional smartphones. This allows manufacturers to provide consumers with a truly bezel-free screen and a larger display-to-body ratio device, enabling a better user experience. The UDC is not confined to smartphones but can also be applied to other devices such as laptops, tablets, and TVs. However, by placing the camera under a screen, image quality is severely degraded due to display panels’ low light transmission rate and diffraction effects caused by pixel arrangement and electronic devices on the display panel.

A typical UDC system places a camera closely underneath the Organic Light-Emitting Diode (OLED) [2] display. As incoming light passes through the OLED display screen, light is diffracted, and its intensity is reduced because of the OLED display pixel arrangement. The degradation process can be modeled as

$$y = k * x + n, \tag{1}$$

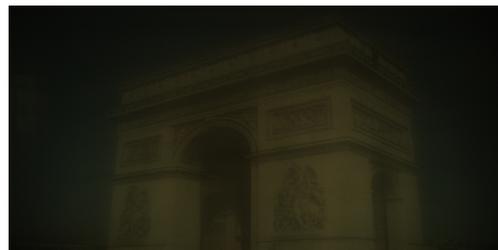
where the input ground-truth image x is convolved with the point spread function (PSF) k of the OLED display, and the noise n added. The output image y captured by the sensor suffers from multiple degradations such as low-light, color-shift, noise, blur, etc. There are two commonly used OLED types (T-OLED, P-OLED) in modern smartphone displays [1],



(a) Ground Truth input.



(b) Image captured under T-OLED.



(c) Image captured under P-OLED.

Fig. 1: Images from the UDC dataset [1]. (a) Ground truth input, (b) image degraded by T-OLED display, (c) by P-OLED. T-OLED images suffer from blur and noise, while P-OLED becomes low-lighted, color-shifted, and hazy.

and some examples of degraded images from UDCs of these systems are shown in Fig. 1 [1].

Our objective of finding x from the observed y belongs to an image restoration problem. Specifically, it can be regarded as deblurring if we consider k as the blur kernel, though not

the same as the conventional out-of-focus or motion blur. A lot of works have been proposed over several decades to alleviate image degradations. Traditional image restoration methods such as Wiener Filter [3] use deconvolution to restore the image convolved with a given degradation kernel (PSF). However, the traditional deconvolution-based methods do not work well due to the ill-posedness of the image restoration problem (inverse of the PSF), especially when the given PSF is spatially wide, which is the case with the UDC systems.

Recent developments in learning-based methods for image restoration have achieved state-of-the-art performances in various areas like deblurring [4], denoising [5]–[7], dehazing [8], light-enhancement [9], high dynamic range imaging (HDR) [10], [11], super-resolution (SR) [12]–[14], and joint HDR and SR [15]. These deep learning-based methods generally work well even with small-sized training image patches. However, using small-sized patches in training is not suitable for UDC image restoration models because the PSF of UDC degradation is spatially very wide. More specifically, the UDC image restoration requires contextual information from wider areas to restore the clean image because the PSF kernel of the display panel is much larger than the kernel of traditional out-of-focus or motion-blur problems.

To address the problem mentioned above, we propose an end-to-end, deep learning-based image restoration model. Our model is a U-Net [16] based architecture that utilizes dilated convolutions [17] and spatially adaptive normalization [18]. Using the dilated convolution, we can suppress the excessive memory increase due to large training patches needed to cope with wide PSF. Also, by using the spatially adaptive normalization scheme, we can exploit additional spatial information for successful restoration.

In summary, our contributions can be summarized as follows.

- We propose an end-to-end convolutional neural network (CNN) for UDC image restoration, which utilizes dilated convolutions and spatially adaptive normalization.
- We show that the spatially adaptive normalization can provide the network with additional spatial information, thereby improving the performance.
- The method can effectively restore images and achieve comparable or often better results than state-of-the-art algorithms.

II. RELATED WORK

A. Under-Display Camera

Since UDC image restoration is a relatively new topic in image restoration fields, there is only a handful of research on this topic. One of the representative works for this topic is [1], where they provide a dataset based on DIV2K using MCIS (Monitor-Camera Imaging System). A challenge [19] was also held to encourage researchers to improve the performance of UDC image restoration, and numerous methods mainly based on deep learning were proposed. Some methods mentioned in [19] use skip or dense connections, while others cascade traditional methods with learning-based algorithms.

B. Image Restoration

Image restoration is an ill-posed problem that aims at restoring a clean image from its degraded version. There are various types of degradations such as blur, noise, down-sampling, etc. In the image degradation process in (1), setting the PSF k as a low-pass filter corresponds to the conventional out-of-focus blur model. Also, setting k as an identity matrix makes it a denoising problem, and a composite operator of blurring and down-sampling makes a super-resolution model. The UDC image restoration problem is the most analogous to deblurring among these image restoration problems, where the blur kernel is a PSF corresponding to the UDC image degradation model. Hence, traditional methods for deblurring, such as Wiener Filter [3] and Richardson-Lucy deconvolution [20], [21], can be used for the restoration of UDC images.

More recently, deep learning-based methods have achieved state-of-the-art performance in image restoration problems as stated previously [4]–[15]. In these works, numerous techniques were introduced, which have been shown to impact restoration performances. For example, dilated convolution [17] was shown to effectively increase the receptive field without increasing parameter size and thus has been successfully used in various computer vision tasks [22], [23]. Wang *et al.* [24] used dilated convolution in denoising to achieve comparable performance to the ones having a larger number of parameters. Also, residual connection [25], [26] and skip connection [27], [28] are the most commonly used techniques in image restoration for better performance and optimization.

C. Normalization

Normalization layers are useful in speeding up and stabilizing the training of neural networks, and the most famous normalization layer is batch normalization [29]. Zhang *et al.* [5] used batch normalization for denoising to speed up the training and boost the performance. Other normalization layers include instance normalization [30], group normalization [31], and layer normalization [32]. These normalization layers do not use external data to normalize the layers, as their main objective is to regularize distributions internally. On the other hand, adaptive instance normalization [33] makes use of external data for normalization. Park *et al.* [18] in particular proposed spatially adaptive normalization, which utilizes semantic layout while performing an affine transformation in normalization layers. Kim *et al.* [34] adopted adaptive instance normalization in denoising for better model generalization and efficient training, resulting in performance improvement.

III. PROPOSED METHOD

To resolve the degradation in UDC images, we propose an end-to-end framework that takes the UDC image as input and outputs a restored image. Our method uses a large receptive field while keeping memory usage at a reasonable level. We also make use of additional spatial information through a spatially adaptive normalization layer. Since the degradation caused by UDC is complex and quite severe, a large receptive field is crucial for improving the training performance. Therefore,

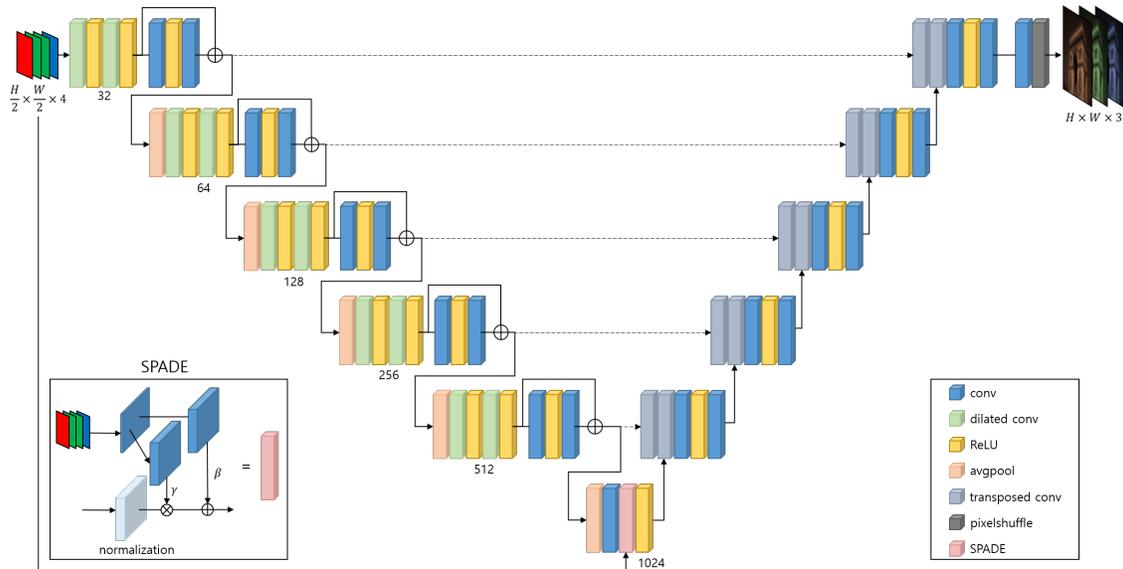


Fig. 2: Proposed network architecture. Based on the structure of U-Net [16], we modify the encoder to use dilated convolutions and residual connections. We also use spatially adaptive normalization [18] in the deepest channel.

the degraded image is processed through a set of encoders which uses dilated convolution to capture features with a larger receptive field without increasing memory consumption. The features are then put through a set of decoders to reconstruct a clean image. Spatially adaptive normalization [18] is applied in the layer of the deepest channel depth to regularize the features and recover the original image using additional spatial information in reconstruction.

A. Network Architecture

The overall architecture of our network is shown in Fig. 2. Based on the well-known architecture of U-Net [16], our network is comprised of a set of encoders and decoders with skip connections. Also, we add dilated convolutions and residual connections at the encoder step of U-Net to maximize the performance. In the layer where the channel depth is deepest, we add a spatially adaptive normalization layer to reconstruct the image more efficiently using additional spatial information. The design of the encoder and decoder is illustrated in Fig. 2.

Encoder The encoder of our network differs from the conventional U-Net in that it has a residual connection and uses dilated convolutions. Residual connections in the encoder help the model to extract better features. We exploit dilated convolution to increase the receptive field without losing spatial resolution so that the model can capture more information to restore a clean image. We use the Average Pool (avgpool) operation as the downsampling operator.

Decoder As the upsampling operator of our decoder, we use transposed convolution. There are several upsampling operations such as pixelshuffle [35], but as we will show in section IV-C, using the pixelshuffle leads to unwanted block and line artifacts in the restored image.

Spatially Adaptive Normalization In the layer where the channel depth is deepest, we add a spatially adaptive normalization layer, also known as SPADE (SPatially Adaptive (DE)normalization) [18], which is illustrated in Fig. 2. According to [18], the SPADE module projects the input mask m with the size of $H \times W \times C$ onto an embedding space, and they are convolved to produce spatially variant learned parameters γ and β . Then, the learned parameters γ and β are multiplied and added element-wise to the normalized activations. Let h^i be the activation value of the i -th layer of a network, and C^i , H^i , and W^i be the number of channels of the layer, height, and width of the i -th activation map in the layer for a batch of N samples, respectively. Then, the activation values at site $(n \in N, c \in C^i, y \in H^i, x \in W^i)$ can be expressed as:

$$\gamma_{c,y,x}^i(m) \frac{h_{n,c,y,x}^i - \mu_c^i}{\sigma_c^i} + \beta_{c,y,x}^i(m), \quad (2)$$

where $h_{n,c,y,x}^i$ is the activation at the site before normalization. Also, μ_c^i and σ_c^i are the mean and standard deviation of activations in channel c , which can be expressed as:

$$\mu_c^i = \frac{1}{NH^iW^i} \sum_{n,y,x} h_{n,c,y,x}^i, \quad (3)$$

$$\sigma_c^i = \sqrt{\frac{1}{NH^iW^i} \sum_{n,y,x} ((h_{n,c,y,x}^i)^2 - (\mu_c^i)^2)}. \quad (4)$$

Adding this spatially adaptive normalization layer assists our model by regularizing the features and providing additional spatial information. In our model, the input mask m is the same input that we put in our restoration process, which is the degraded UDC image that we aim to restore, with the size of

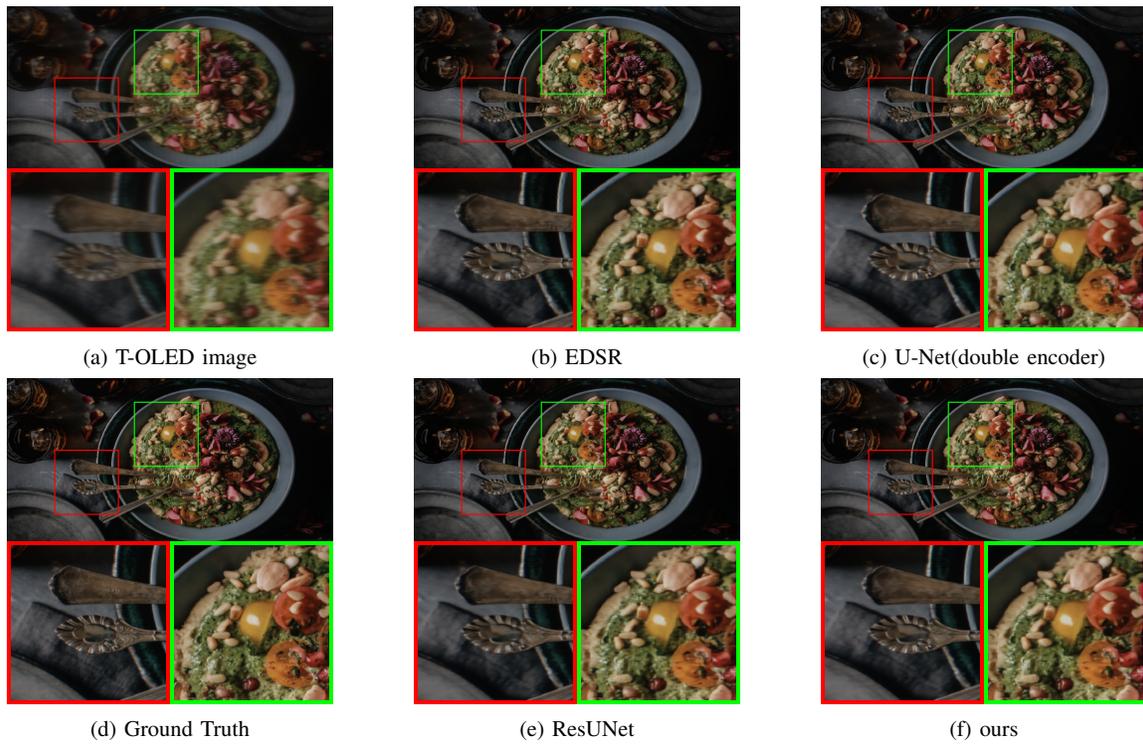


Fig. 3: Visual comparison of restored T-OLED images using methods mentioned in section IV.

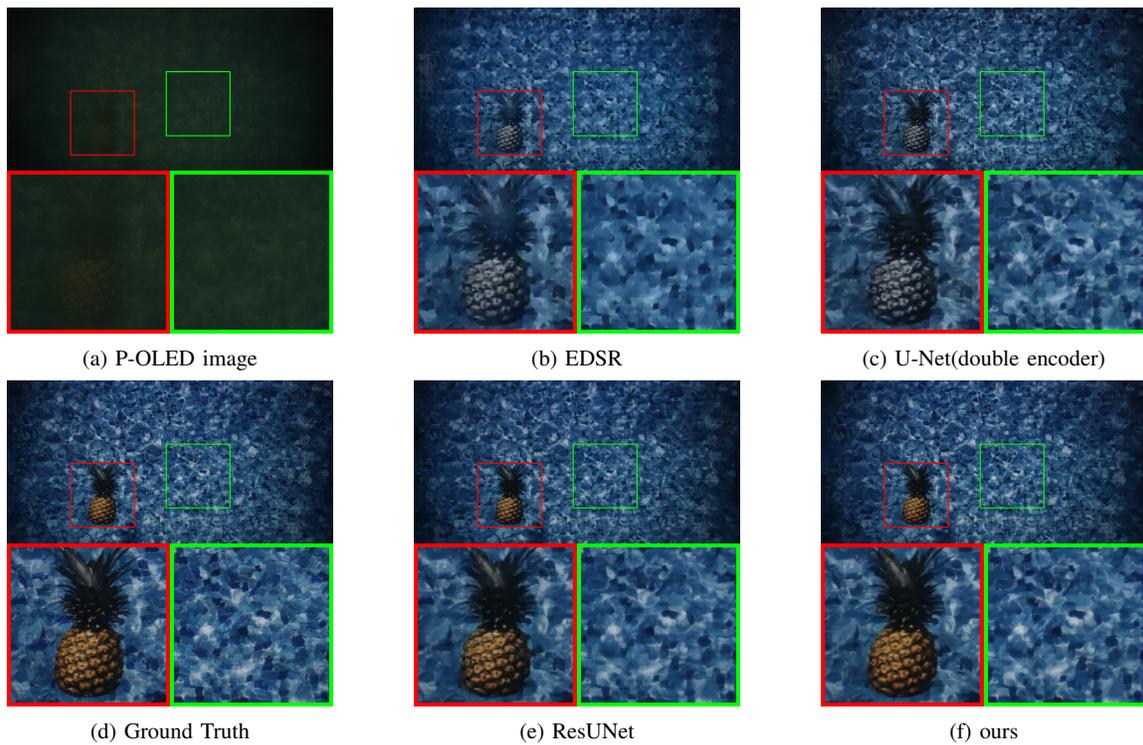


Fig. 4: Visual comparison of restored P-OLED images using methods mentioned in section IV.

TABLE I: Comparison of models for UDC image restoration. The best results are denoted in red and the second best in blue.

Method	# Parameters	T-OLED		P-OLED	
		PNSR	SSIM	PNSR	SSIM
EDSR [12]	1.26M	35.16	0.9674	26.62	0.9104
U-Net(double encoder) [1]	35.72M	36.23	0.9724	29.60	0.9353
ResUNet [36]	33.20M	36.32	0.9708	31.62	0.9471
ours	30.31M	37.05	0.9732	32.20	0.9526

TABLE II: Effects of dilated convolution, spatially adaptive normalization layer, and upsampling layer on performance. Best results are highlighted in bold.

	Method				T-OLED		P-OLED	
	dilated conv	norm	pixelshuffle	convt	PSNR	SSIM	PSNR	SSIM
(a)			✓		36.32	0.9708	31.62	0.9471
(b)	✓		✓		36.41	0.9702	31.84	0.9484
(c)	✓	✓	✓		36.86	0.9722	32.19	0.9519
(d)	✓	✓		✓	37.05	0.9732	32.20	0.9526

$\frac{H}{2} \times \frac{W}{2} \times 4$. We will show that this strategy is effective in improving the performance of our model in the ablation study of section IV-C.

Upscaler The upscaler at the end of our model is a simple convolutional layer with kernel size 1×1 , followed by pixelshuffle with factor 2.

B. Loss Function

We train our network using L1 loss. Our loss function is defined as:

$$\mathcal{L} = \|\mathcal{F}(y) - x\|_1, \tag{5}$$

where y and x are a UDC image and a ground-truth image defined in (1), and \mathcal{F} is our model for restoration.

C. Implementation Details

The proposed model takes 4-channel raw data and outputs a 3-channel RGB image. There are five encoder-decoder pairs, and the kernel size of convolutional layers are all 3×3 except the upscaler. We train the model using Adam optimizer ($\beta_1 = 0.9, \beta_2 = 0.999$ and $\epsilon = 10^{-8}$), with initial learning rate of 10^{-4} and decay factor 0.5.

IV. EXPERIMENTS

A. Experimental Settings

Datasets We use the dataset provided by [1] for both training and testing. The dataset consists of 240 images for training, 30 for validation, and 30 for testing for both T-OLED and P-OLED images. We use 256×256 sized patches for T-OLED images, while we use 1024×1024 sized patches for P-OLED because the degradation is more severe with the P-OLED. As only an 8-bit RGB version of the dataset is available, we make an 8-bit raw version of the dataset and form the input training data as

in [9]. The training data is augmented with random horizontal, vertical flips, and 180-degree rotations.

Evaluation Metrics We evaluate the methods using PSNR and SSIM. Higher PSNR and SSIM mean better performance.

B. Experimental Results

Quantitative Evaluation We compare our results with baseline models (EDSR, U-Net double encoder) introduced in [1] and a model (ResUNet) [36] introduced in the UDC 2020 challenge [19] that inspired our model. Each method is trained according to the training scheme mentioned in [1] and [36]. To make a fair comparison between methods based on U-Net (U-Net with double encoder, ResUNet, ours) regarding parameter numbers, each channel depth in layers of U-Net with double encoder has been doubled from the baseline proposed in [1].

Our method achieves the best performance compared to all the aforementioned methods and especially has significant improvements in P-OLED image recovery. Our method also has the fewest number of parameters except for EDSR [12]. The detailed results are shown in Table I.

Qualitative Evaluation Visual comparisons of the methods are shown in Figs. 3 and 4, for T-OLED and P-OLED images, respectively. Compared to other methods of UDC restoration, our method produces high-quality images in color and texture. Simple models such as EDSR [12] and U-Net (double encoder) [1] do not work well in recovering color in P-OLED images and produce artifacts. ResUNet [36] performs moderately but produces unsatisfactory line and block artifacts.

C. Ablation Study

In this part, we analyze the effect of techniques used to improve the performance of our model by conducting ablative experiments. All experiments are conducted under the same condition as aforementioned. The results are in Table II.

Dilated Convolution By comparing the rows (a) and (b) in Table II, we can see that adding dilated convolution at the encoder blocks of the model results in a noticeable increase in PSNR. This shows that dilated convolution enlarges the receptive field of the model, improving the performance of restoration in return.

Spatially Adaptive Normalization Using spatially adaptive normalization in the deepest block of the model helps the restoration process as the model has additional spatial context. From the rows (b) and (c), we can see that there is an increase of 0.45dB and 0.35dB in T-OLED and P-OLED images, respectively. This demonstrates that adopting spatially adaptive normalization is an effective strategy in improving the performance of the model.

Transposed Convolution Compared to transposed convolutional layers, using the pixelshuffle layer for upsampling generates block and line artifacts that are obviously unnatural and undesirable. These artifacts are more prevalent in P-OLED images, and thus we need a more sophisticated restoration process. Therefore, we choose transposed convolution as our upsampling operator. As shown in Fig. 5, the model using transposed convolutional layers suppresses the artifacts and can produce visually plausible results.

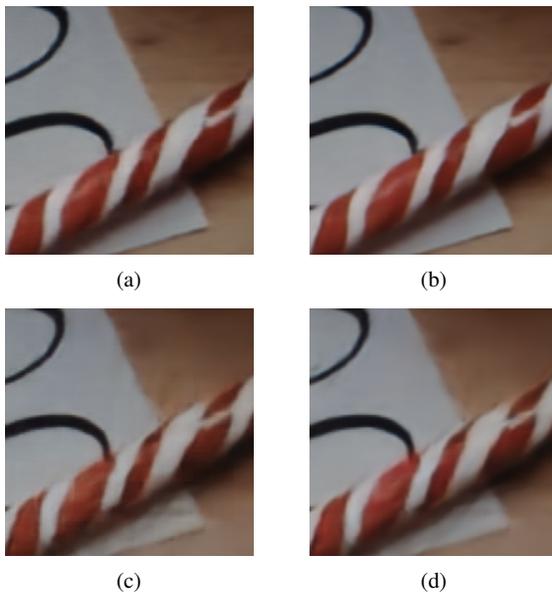


Fig. 5: Comparison of restored images from our baseline model with pixelshuffle and transposed convolution as upsampling layers. (a): T-OLED with pixelshuffle, (b): T-OLED with transposed convolution, (c): P-OLED with pixelshuffle, (d): P-OLED with transposed convolution. Zoom in for a closer view.

V. CONCLUSION

In this paper, we have presented a new method for restoring Under-Display Camera images. The proposed model uses dilated convolutions to increase the receptive field. Also, we

add spatially adaptive normalization to provide the model with extra spatial information. These two techniques improve the performance of UDC image restoration. We have also demonstrated quantitatively and qualitatively that our model can effectively recover degraded images.

ACKNOWLEDGMENT

This paper was result of the research project supported by SK hynix Inc. This work was also partly supported by Institute of Information & communications Technology Planning & Evaluation (IITP) grant funded by the Korea government(MSIT) (2021-0-02068, Artificial Intelligence Innovation Hub).

REFERENCES

- [1] Y. Zhou, D. Ren, N. Emerton, S. Lim, and T. Large, "Image restoration for under-display camera," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 9179–9188.
- [2] I. G. Wenke, "Organic light emitting diode (oled)," *Research gate*, 2016.
- [3] J. S. Goldstein, I. S. Reed, and L. L. Scharf, "A multistage representation of the wiener filter based on orthogonal projections," *IEEE Transactions on Information Theory*, vol. 44, no. 7, pp. 2943–2959, 1998.
- [4] S. Nah, S. Son, R. Timofte, and K. M. Lee, "Ntire 2020 challenge on image and video deblurring," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020, pp. 416–417.
- [5] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang, "Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising," *IEEE transactions on image processing*, vol. 26, no. 7, pp. 3142–3155, 2017.
- [6] J. Liu *et al.*, "Learning raw image denoising with bayer pattern unification and bayer preserving augmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 0–0.
- [7] S. Guo, Z. Yan, K. Zhang, W. Zuo, and L. Zhang, "Toward convolutional blind denoising of real photographs," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1712–1722.
- [8] C. O. Ancuti, C. Ancuti, F-A. Vasluianu, and R. Timofte, "Ntire 2020 challenge on nonhomogeneous dehazing," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020, pp. 490–491.
- [9] C. Chen, Q. Chen, J. Xu, and V. Koltun, "Learning to see in the dark," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3291–3300.
- [10] S.-H. Lee, H. Chung, and N. I. Cho, "Exposure-structure blending network for high dynamic range imaging of dynamic scenes," *IEEE Access*, vol. 8, pp. 117 428–117 438, 2020.
- [11] H. Chung, Y. Kim, J. Jo, S.-H. Lee, and N. I. Cho, "Kernel prediction network for detail-preserving high dynamic range imaging," in *2019 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*. IEEE, 2019, pp. 1589–1594.
- [12] B. Lim, S. Son, H. Kim, S. Nah, and K. Mu Lee, "Enhanced deep residual networks for single image super-resolution," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2017, pp. 136–144.
- [13] A. Lugmayr, M. Danelljan, and R. Timofte, "Ntire 2020 challenge on real-world image super-resolution: Methods and results," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020, pp. 494–495.
- [14] J. W. Soh, G. Y. Park, J. Jo, and N. I. Cho, "Natural and realistic single image super-resolution with explicit natural manifold discrimination," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8122–8131.
- [15] J. W. Soh, J. S. Park, and N. I. Cho, "Joint high dynamic range imaging and super-resolution from a single image," *IEEE Access*, vol. 7, pp. 177 427–177 437, 2019.
- [16] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.

- [17] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," *arXiv preprint arXiv:1511.07122*, 2015.
- [18] T. Park, M.-Y. Liu, T.-C. Wang, and J.-Y. Zhu, "Semantic image synthesis with spatially-adaptive normalization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2337–2346.
- [19] Y. Zhou *et al.*, "Udc 2020 challenge on image restoration of under-display camera: Methods and results," in *European Conference on Computer Vision*. Springer, 2020, pp. 337–351.
- [20] W. H. Richardson, "Bayesian-based iterative method of image restoration," *JoSA*, vol. 62, no. 1, pp. 55–59, 1972.
- [21] L. B. Lucy, "An iterative technique for the rectification of observed distributions," *The astronomical journal*, vol. 79, p. 745, 1974.
- [22] F. Yu, V. Koltun, and T. Funkhouser, "Dilated residual networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 472–480.
- [23] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 4, pp. 834–848, 2017.
- [24] T. Wang, M. Sun, and K. Hu, "Dilated deep residual network for image denoising," in *2017 IEEE 29th international conference on tools with artificial intelligence (ICTAI)*. IEEE, 2017, pp. 1272–1279.
- [25] Y. Zhang, Y. Tian, Y. Kong, B. Zhong, and Y. Fu, "Residual dense network for image super-resolution," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 2472–2481.
- [26] —, "Residual dense network for image restoration," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [27] X.-J. Mao, C. Shen, and Y.-B. Yang, "Image restoration using very deep convolutional encoder-decoder networks with symmetric skip connections," *arXiv preprint arXiv:1603.09056*, 2016.
- [28] T. Tong, G. Li, X. Liu, and Q. Gao, "Image super-resolution using dense skip connections," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 4799–4807.
- [29] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *International conference on machine learning*. PMLR, 2015, pp. 448–456.
- [30] D. Ulyanov, A. Vedaldi, and V. Lempitsky, "Instance normalization: The missing ingredient for fast stylization," *arXiv preprint arXiv:1607.08022*, 2016.
- [31] Y. Wu and K. He, "Group normalization," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 3–19.
- [32] J. L. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," *arXiv preprint arXiv:1607.06450*, 2016.
- [33] X. Huang and S. Belongie, "Arbitrary style transfer in real-time with adaptive instance normalization," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 1501–1510.
- [34] Y. Kim, J. W. Soh, G. Y. Park, and N. I. Cho, "Transfer learning from synthetic to real-noise denoising with adaptive instance normalization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 3482–3492.
- [35] W. Shi *et al.*, "Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 1874–1883.
- [36] Q. Yang, Y. Liu, J. Tang, and T. Ku, "Residual and dense unet for under-display camera restoration," in *European Conference on Computer Vision*. Springer, 2020, pp. 398–408.