# RAISING THE LIMIT OF IMAGE RESCALING USING AUXILIARY ENCODING

*Chenzhong Yin*[1,*]     *Zhihong Pan*[2,*]     *Xin Zhou*[2]     *Le Kang*[2]     *Paul Bogdan*[1]

[1]University of Southern California, Los Angeles, CA, USA
[2]Baidu Research (USA), Sunnyvale, CA, 94089, USA

## ABSTRACT

Normalizing flow models using invertible neural networks (INN) have been widely investigated for successful generative image super-resolution (SR) by learning the transformation between the normal distribution of latent variable $z$ and the conditional distribution of high-resolution (HR) images gave a low-resolution (LR) input. Recently, image rescaling models like IRN utilize the bidirectional nature of INN to push the performance limit of image upscaling by optimizing the downscaling and upscaling steps jointly. While the random sampling of latent variable $z$ is useful in generating diverse photo-realistic images, it is not desirable for image rescaling when accurate restoration of the HR image is more important. Hence, in places of random sampling of $z$, we propose auxiliary encoding modules to further push the limit of image rescaling performance. Two options to store the encoded latent variables in downscaled LR images, both readily supported in existing image file format, are proposed. One is saved as the alpha-channel, the other is saved as meta-data in the image header, and the corresponding modules are denoted as suffixes -A and -M respectively. Optimal network architectural changes are investigated for both options to demonstrate their effectiveness in raising the rescaling performance limit on different baseline models including IRN and DLV-IRN.

*Index Terms*— Super-Resolution, Image Rescaling

## 1. INTRODUCTION

Currently, ultra-high resolution (HR) images are often needed to be reduced from their original resolutions to lower ones due to various limitations like display or transmission. Once resized, there could be subsequent needs of scaling them up so it is useful to restore more high-frequency details [1]. While deep learning super-resolution (SR) models [2, 3, 4] are powerful tools to reconstruct HR images from low-resolution (LR) inputs, they are often limited to pre-defined image downscaling methods. Additionally, due to memory and speed constraints, HR images or videos are also commonly resized to lower resolution for downstream computer vision tasks like image classification and video understanding.

Similarly, they rely on conventional resizing methods which are subject to information loss and have negative impact on downstream tasks [5]. Hence, learned image downscaling techniques with minimum loss in high-frequency information are quite indispensable for both scenarios. Lastly, it is known that SR models optimized for upscaling only are subject to model stability issues when multiple downscaling-to-upscaling cycles are applied [6] so it further validates the necessity of learning downscaling and upscaling jointly.

To overcome these challenges and utilize the relationship between upscaling and downscaling steps, recent works designed the encoder-decoder framework to unite these two independent tasks together. Kim *et al*. [7] utilized autoencoder (AE) architecture, where the encoder is the downscaling network and the decoder is the upscaling network, to find the optimal LR result that maximizes the restoration performance of the HR image. Sun *et al*. [8] designed a learned content adaptive image downscaling model in which an SR model is trained simultaneously to best recover the HR images. Later on, Li *et al*. [9] proposed a learning approach for image compact-resolution using a convolutional neural network (CNN-CR) where the image SR problem is formulated to jointly minimize the reconstruction loss and the regularization loss. Although the above models can efficiently improve the quality of HR images recovered from corresponding LR images, these works only optimize downscaling and SR separately, while ignoring the potential mutual intension between downscaling and inverse upscaling.

More recently, a jointly optimized rescaling model was proposed by Xiao *et al*. [10] to achieve significantly improved performance. An Invertible Rescaling Net (IRN) was designed to model the reciprocal nature of the downscaling and upscaling processes. For downscaling, IRN was trained to convert HR input to visually-pleasing LR output and a latent variable $z$. As $z$ is trained to follow an input-agnostic Gaussian distribution, the HR image can be accurately reconstructed during the inverse up-scaling procedure although $z$ is randomly sampled from a normal distribution. Nevertheless, the model's performance can be further improved if the high-frequency information remaining in $z$ is efficiently stored.

To resolve the above difficulties and take full potential of the IRN, here we propose two approaches, namely the IRN-meta (IRN-M) and IRN-alpha (IRN-A), respectively, to ef-
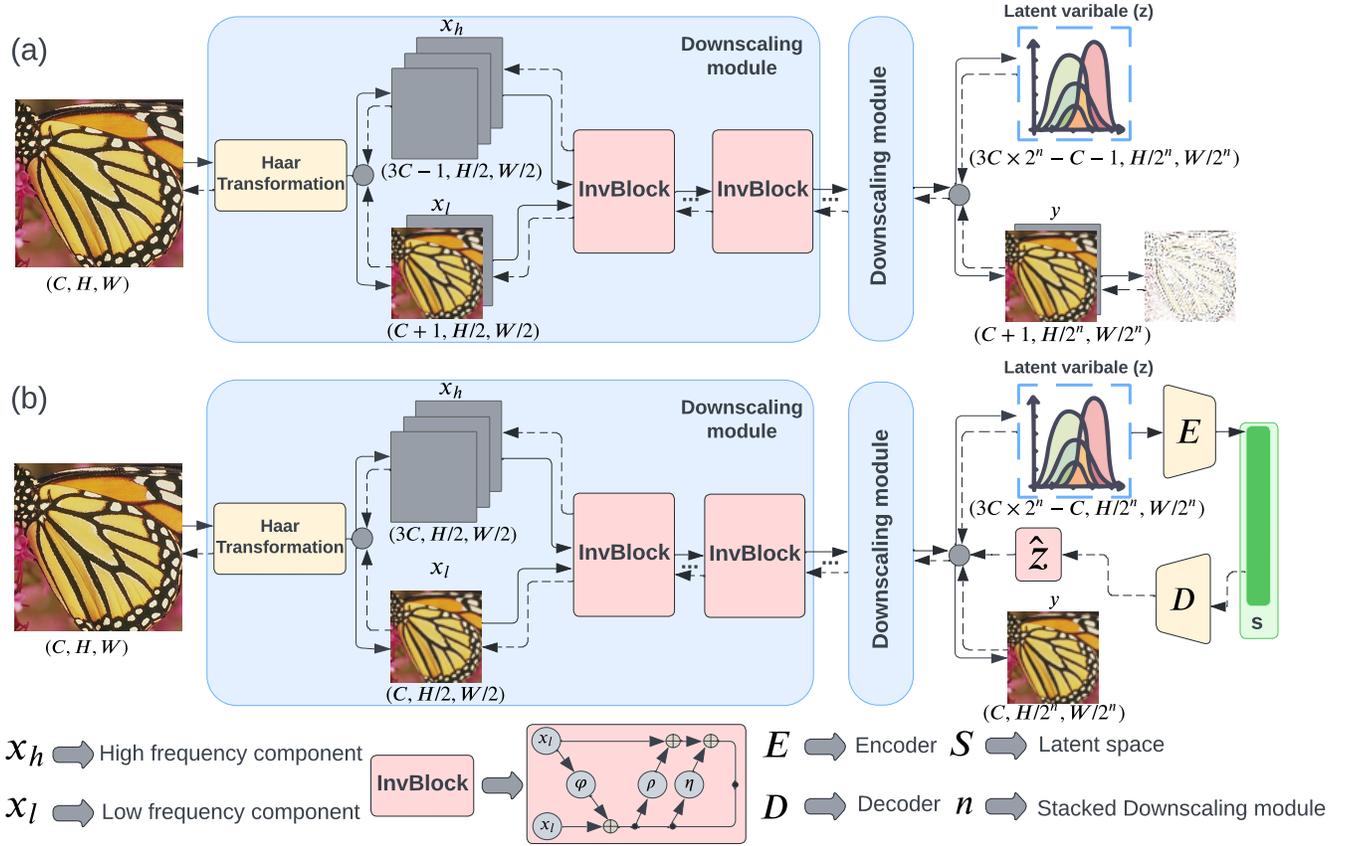
---

**Fig. 1**: Illustration of invertible image rescaling network architecture: (a) RGBA approach and (b) metadata approach.

ficiently compress the high frequency information stored in $z$, which can be used to recover $z$ and help restore HR image consequently during the inverse up-scaling. For IRN-A, we train the model to extract a fourth LR channel in addition to the LR RGB channels. It represents essential high frequency information which was lost in the IRN baseline due to random sampling of $z$, and is saved as the alpha-channel of saved LR output. For IRN-M approach, an AE module is trained to compress $z$ as a compact latent variable, which can be saved as metadata of the LR output. In the inverse upscaling process, $z$ is restored from the latent space by utilizing the well-trained decoder. Both modules are also successfully applied to the state-of-the-art (SOTA) rescaling model DLV-IRN [11]. In summary, the main contribution of this paper is that we are the first to compress the high-frequency information in $z$, which is not fully utilized in current invertible image rescaling models, to improve the restored HR image quality in upscaling progress.

## 2. PROPOSED METHOD

### 2.1. IRN-A

Fig. 1 (a) shows the IRN-A network architectures, where the invertible neural network blocks (InvBlocks) are referenced from previous work IRN [10]. In the new model, the input HR

image is resized via Haar transformation before splitting to a lower branch $x_l$ and a higher branch $x_h$. More specifically, Haar transformation converts the input HR image $(C, H, W)$ into a matrix of shape $(4C, H/2, W/2)$, where $C$, $H$, and $W$ represent image color channels, height and width respectively. The first $C$ channels represent low-frequency components of the input image in general and the remaining $3C$ channels represent the high-frequency information on vertical, horizontal and diagonal directions respectively. Different from the IRN baseline, which uses only the $C$ low-frequency channels in the lower branch, we add 1 additional channel, denoted as alpha-channel for convenience as it would be stored as the alpha-channel in RGBA format, in the lower branch $x_l$ to store the compressed high-frequency information. After the first Haar transformation, the alpha-channel is initialized with the average value across all $3C$ high-frequency channels, and only $3C - 1$ channels are included in $x_h$ as the first channel is removed to make the total number of channels remain constant.

After channel splitting, $x_l$ and $x_h$ are fed into cascaded InvBlocks and transformed to an LR RGBA image $y$ and an auxiliary latent variable $z$. First three channels of $y$ consist of the visual RGB channels and the fourth channel contains the compressed high-frequency components transformed along the InvBlocks. The alpha-channel was normalized via a

*sigmoid* function, $S(\alpha) = \frac{1}{1+e^{-\alpha}}$, to help quantization of the alpha-channel and maintain training stability.

For the inverse upscaling process, the model needs to recover $z$, denoted as $\hat{z}$ as it is not stored. In previous work, $\hat{z}$ was randomly drawn from normal Gaussian distribution. While this helps creating diverse samples in generative models, it is not optimal for tasks like image rescaling which aims to restore one HR image instead of diverse variations. Therefore, we set $\hat{z}$ as 0, the mean value of the normal distribution, for the inverse up-scaling process. This technique was also validated in previous works like FGRN [12] and DLV-IRN [11]. Of note, at the end of inverse process, the deleted high frequency channel needs to be recovered as

$$x_m = 3C \times x_\alpha - \sum_{i=1}^{3C-1} x_h^i \quad (1)$$

where $x_m$ represents the channel removed from $x_h$ and $x_\alpha$ represents the alpha-channel in $x_l$.

## 2.2. IRN-M

Besides storing the compressed high-frequency information in a separate alpha-channel, we also propose an alternative space-saving approach to store the extracted information as metadata of the image file. Image metadata is text information pertaining to an image file that is embedded into the image file or contained in a separate file in a digital asset management system. Metadata is readily supported by existing image format so this proposed method could be easily integrated with current solutions.

The network architecture of our metadata approach is shown in Fig. 1 (b). Here $x_l$ and $x_h$, same as the IRN baseline, are split from Haar transformed $4C$ channels to $C$ and $3C$ channels respectively. Unlike the RGBA approach, the metadata method uses an encoder at the end to compress the $z$ and save the latent vector $S$ as metadata, rather than saving as the alpha-channel of the output. $S$ will be decompressed by the decoder for the inverse upscaling step. In our AE architecture, the encoder compacts the number of $z$ channels from $3C \times n^2 - C$ to 4 via 2D convolution layers and compresses the $z$'s height and width from $(H/2^n, W/2^n)$ to $(H/2^{n+2}, W/2^{n+2})$ by using max-pooling layers. Here $n$ is 1 or 2 depending on the scale factor of $2\times$ or $4\times$. Of note, the AE was pre-trained with MSE loss before being embedded into the model structure.

After placing the well-trained AE in the IRN architecture, the entire structure was trained to minimize the following mixture loss function:

$$L = \lambda_1 L_r + \lambda_2 L_g + \lambda_3 L_d + \lambda_4 L_{mse} \quad (2)$$

where $L_r$ is the L1 loss for reconstructing HR image; $L_g$ is the L2 loss for the generated LR image; $L_d$ is the distribution matching loss; and $L_{mse}$ is the MSE loss between the input of the encoder and the output of the decoder.

**Table 1**: Comparison of $4\times$ upscaling results using different IRN-A hyperparameters and settings. The best results are highlighted in <span style="color:red">red</span>.

| IRN-A | $\alpha_{avg}$ | BSD100 PSNR/SSIM↑ | Urban100 PSNR/SSIM↑ | DIV2K PSNR/SSIM↑ |
|---|---|---|---|---|
| Post-split | ✗ | 32.66 / 0.9083 | 32.50 / 0.9328 | 36.19 / 0.9464 |
| Pre-split | ✗ | 33.02 / 0.9132 | 32.17 / 0.9186 | 36.60 / 0.9495 |
| | ✓ | 33.12 / 0.9150 | 33.10 / 0.9384 | 36.67 / 0.9504 |

**Table 2**: Comparison of $4\times$ upscaling results using different IRN-M hyperparameters and settings. The best results are highlighted in <span style="color:red">red</span>.

| IRN-M | $AE_p$ | $AE_f$ | BSD100 PSNR/SSIM↑ | Urban100 PSNR/SSIM↑ | DIV2K PSNR/SSIM↑ |
|---|---|---|---|---|---|
| 2layers | ✗ | ✗ | 31.41 / 0.8771 | 30.79 / 0.9074 | 34.79 / 0.9283 |
| | ✓ | ✓ | 31.58 / 0.8793 | 31.30 / 0.9123 | 35.06 / 0.9306 |
| | ✓ | ✗ | 31.65 / 0.8804 | 31.34 / 0.9154 | 35.09 / 0.9306 |
| 4layers | ✗ | ✗ | 28.15 / 0.7765 | 25.82 / 0.7989 | 30.72 / 0.8591 |
| | ✓ | ✗ | 31.69 / 0.8812 | 31.44 / 0.9143 | 35.15 / 0.9314 |

## 3. EXPERIMENTS

Following the same training strategy and hyperparameters in IRN baseline, our models were trained on the DIV2K [13] dataset, which includes 800 HR training images. IRN-M and IRN-A were trained with 500,000 and 250,000 iterations respectively. Both models were evaluated across five benchmark datasets: Set5 [14], Set14 [15], BSD100 [16], Urban100 [17] and the validation set of DIV2K. The upscaled images quality across different models were assessed via the peak noise-signal ratio (PSNR) and SSIM on the Y channel of the YCbCr color space. Following previous works [12, 11], as it is not beneficial to add randomness in restoring HR images, we set $\hat{z}$ as 0 during the inverse up-scaling process for both training and validation steps in all experiments.

### 3.1. Ablation study

As the transformed alpha-channel is the key innovation for improved performance for IRN-A, the pre-splitting and initial settings of the alpha-channel before the forward transformation process are very important. For better analysis of their effects, Table 1 shows an ablation study that compares the results for different settings of the alpha-channel, where "post-split" and "pre-split" refer to splitting the alpha-channel after the downscaling module or before the InvBlock respectively, and $\alpha_{avg}$ represents presetting the average value of high-frequency information in the pre-split alpha-channel. From Table 1, we notice that using the $\alpha_{avg}$ with pre-split architecture performs best across all options.

The IRN-M model constructs the HR image by decoding the latent space $s$ saved in the metadata file. Table 2 shows another ablation study for determining the optimal AE structure, where $AE_p$ represents that AE, before training as part of IRN-M, is pre-trained using MSE loss with standalone random $z$;
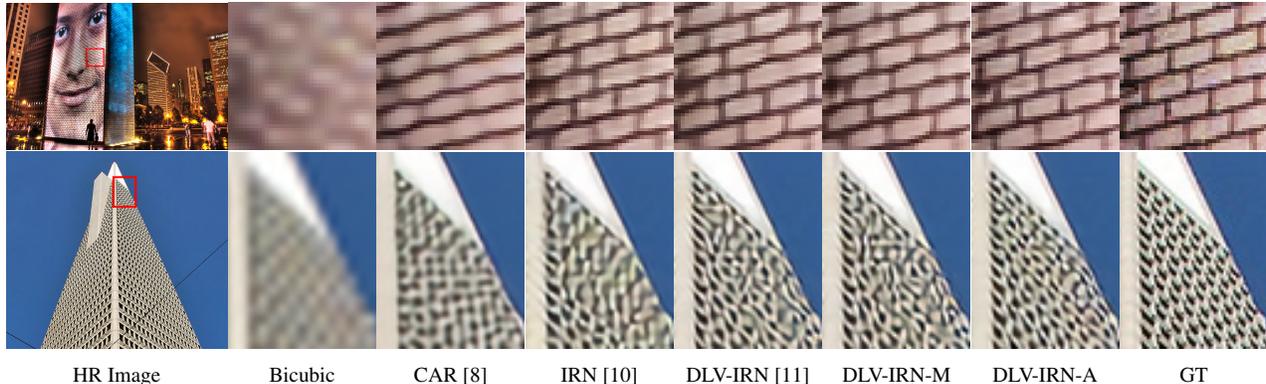
| HR Image | Bicubic | CAR [8] | IRN [10] | DLV-IRN [11] | DLV-IRN-M | DLV-IRN-A | GT |

**Fig. 2**: Visual examples from Urban100 test set (Best viewed in online version with zoom-in).

**Table 3**: Quantitative results of upscaled $\times 4$ images of 5 datasets across different bidirectional rescaling approaches. The best two results highlighted in red and blue respectively.

| Method | Scale | Set5 [14] PSNR/SSIM↑ | Set14 [15] PSNR/SSIM↑ | BSD100 [16] PSNR/SSIM↑ | Urban100 [17] PSNR/SSIM↑ | DIV2K [18] PSNR/SSIM↑ |
|---|---|---|---|---|---|---|
| CAR [8] | 2 | 38.94 / 0.9658 | 35.61 / 0.9404 | 33.83 / 0.9262 | 35.24 / 0.9572 | 38.26 / 0.9599 |
| IRN [10] | 2 | 43.99 / 0.9871 | 40.79 / 0.9778 | 41.32 / 0.9876 | 39.92 / 0.9865 | 44.32 / 0.9908 |
| FGRN [12] | 2 | 44.15 / 0.9902 | 42.28 / 0.9840 | 41.87 / 0.9887 | 41.71 / 0.9904 | 45.08 / 0.9917 |
| DLV-IRN [11] | 2 | 45.42 / 0.9910 | 42.16 / 0.9839 | 42.91 / 0.9916 | 41.29 / 0.9904 | 45.58 / 0.9934 |
| **DLV-IRN-M** | 2 | 45.83 / 0.9916 | 42.47 / 0.9850 | 43.38 / 0.9925 | 41.77 / 0.9911 | 45.91 / 0.9939 |
| **DLV-IRN-A** | 2 | 47.81 / 0.9937 | 44.96 / 0.9884 | 47.15 / 0.9967 | 45.07 / 0.9953 | 48.94 / 0.9968 |
| CAR [8] | 4 | 33.88 / 0.9174 | 30.31 / 0.8382 | 29.15 / 0.8001 | 29.28 / 0.8711 | 32.82 / 0.8837 |
| IRN [10] | 4 | 36.19 / 0.9451 | 32.67 / 0.9015 | 31.64 / 0.8826 | 31.41 / 0.9157 | 35.07 / 0.9318 |
| HCFlow [19] | 4 | 36.29 / 0.9468 | 33.02 / 0.9065 | 31.74 / 0.8864 | 31.62 / 0.9206 | 35.23 / 0.9346 |
| FGRN [12] | 4 | 36.97 / 0.9505 | 33.77 / 0.9168 | 31.83 / 0.8907 | 31.91 / 0.9253 | 35.15 / 0.9322 |
| DLV-IRN [11] | 4 | 36.62 / 0.9484 | 33.26 / 0.9093 | 32.05 / 0.8893 | 32.26/ 0.9253 | 35.55/ 0.9363 |
| **DLV-IRN-M** | 4 | 36.67 / 0.9490 | 33.33 / 0.9105 | 32.12 / 0.8909 | 32.33 / 0.9264 | 35.63 / 0.9373 |
| **DLV-IRN-A** | 4 | 37.56 / 0.9566 | 34.12 / 0.9246 | 33.12 / 0.9150 | 33.10 / 0.9384 | 36.67 / 0.9504 |

$AE_f$ represents fixing the AE during training the IRN-M; and "2layers" and "4layers" represent two and four convolutional layers used in AE respectively. As shown in Table 2, using the IRN-M with pre-trained 4 layers AE and not fixing the AE during training has the best performance. Of all three settings, pre-training of AE is the most critical factor in maximizing performance.

### 3.2. Image rescaling

The quantitative comparison results for HR image reconstruction are shown in Table 3. Rather than choosing SR models which only optimize upscaling steps, we consider SOTA bidirectional (jointly optimizing downscaling and upscaling steps) models for fair comparison [8, 10, 11, 12, 19]. As shown in Table 3, DLV-IRN-A is efficient at storing high-frequency information in the alpha-channel and consequently outperforms its baseline DLV-IRN, as well as other models, including HCFlow and IRN models, which randomly samples $\hat{z}$ for the upscaling step. For DLV-IRN-M, while not as good as the -A variant, it still performs better than all other models, only trailing behind FGRN for two small test sets at

$4\times$. Hence we conclude that both -M and -A modules can improve the modeling of the high-frequency information and help restore the HR image consequently. Visual examples of the $4\times$ test in Fig 2 also validate the improved performance from our models.

## 4. CONCLUSIONS

To fully mine the potential of image rescaling models based on INN, two novel modules are proposed to store otherwise lost high-frequency information $z$. The IRN-M model utilizes an autoencoder to compress $z$ and save as metadata in native image format so it can be decoded to an approximate of $z$, while IRN-A adds an additional channel to store crucial high-frequency information, which can be quantized and stored as the alpha-channel, in addition to the RGB channels, in existing RGBA format. With carefully designed autoencoder and alpha-channel pre-split, it is shown that both modules can improve the upscaling performance significantly comparing to the IRN baseline. The proposed modules are also applicable to newer baseline models like DLV-IRN and DLV-IRN-A is by far the best, which further pushes the limit of image rescaling performance with a significant margin.

# 5. REFERENCES

[1] Jinbo Xing, Wenbo Hu, and Tien-Tsin Wong, "Scale-arbitrary invertible image downscaling," *arXiv preprint arXiv:2201.12576*, 2022.

[2] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang, "Learning a deep convolutional network for image super-resolution," in *European conference on computer vision*. Springer, 2014, pp. 184–199.

[3] Yulun Zhang, Kunpeng Li, Kai Li, Lichen Wang, Bineng Zhong, and Yun Fu, "Image super-resolution using very deep residual channel attention networks," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 286–301.

[4] Chenzhong Yin, Phoebe Imms, Mingxi Cheng, Anar Amgalan, Nahian F Chowdhury, Roy J Massett, Nikhil N Chaudhari, Xinghe Chen, Paul M Thompson, Paul Bogdan, et al., "Anatomically interpretable deep learning of brain age captures domain-specific cognitive impairment," *Proceedings of the National Academy of Sciences*, vol. 120, no. 2, pp. e2214634120, 2023.

[5] Yuan Tian, Guo Lu, Xiongkuo Min, Zhaohui Che, Guangtao Zhai, Guodong Guo, and Zhiyong Gao, "Self-conditioned probabilistic learning of video rescaling," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 4490–4499.

[6] Zhihong Pan, Baopu Li, Dongliang He, Mingde Yao, Wenhao Wu, Tianwei Lin, Xin Li, and Errui Ding, "Towards bidirectional arbitrary image rescaling: Joint optimization and cycle idempotence," pp. 17389–17398, 2022.

[7] Heewon Kim, Myungsub Choi, Bee Lim, and Kyoung Mu Lee, "Task-aware image downscaling," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 399–414.

[8] Wanjie Sun and Zhenzhong Chen, "Learned image downscaling for upscaling using content adaptive resampler," *IEEE Transactions on Image Processing*, vol. 29, pp. 4027–4040, 2020.

[9] Yue Li, Dong Liu, Houqiang Li, Li Li, Zhu Li, and Feng Wu, "Learning a convolutional neural network for image compact-resolution," *IEEE Transactions on Image Processing*, vol. 28, no. 3, pp. 1092–1107, 2018.

[10] Mingqing Xiao, Shuxin Zheng, Chang Liu, Yaolong Wang, Di He, Guolin Ke, Jiang Bian, Zhouchen Lin, and Tie-Yan Liu, "Invertible image rescaling," in *European Conference on Computer Vision*. Springer, 2020, pp. 126–144.

[11] Min Zhang, Zhihong Pan, Xin Zhou, and C-C Jay Kuo, "Enhancing image rescaling using dual latent variables in invertible neural network," *arXiv preprint arXiv:2207.11844*, 2022.

[12] Shang Li, Guixuan Zhang, Zhengxiong Luo, Jie Liu, Zhi Zeng, and Shuwu Zhang, "Approaching the limit of image rescaling via flow guidance," *arXiv preprint arXiv:2111.05133*, 2021.

[13] Sefi Bell-Kligler, Assaf Shocher, and Michal Irani, "Blind super-resolution kernel estimation using an internal-gan," *Advances in Neural Information Processing Systems*, vol. 32, 2019.

[14] Marco Bevilacqua, Aline Roumy, Christine Guillemot, and Marie line Alberi Morel, "Low-complexity single-image super-resolution based on nonnegative neighbor embedding," in *Proceedings of the British Machine Vision Conference*. 2012, pp. 135.1–135.10, BMVA Press.

[15] Roman Zeyde, Michael Elad, and Matan Protter, "On single image scale-up using sparse-representations," in *International conference on curves and surfaces*. Springer, 2010, pp. 711–730.

[16] David Martin, Charless Fowlkes, Doron Tal, and Jitendra Malik, "A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics," in *Proceedings Eighth IEEE International Conference on Computer Vision. ICCV 2001*. IEEE, 2001, vol. 2, pp. 416–423.

[17] Jia-Bin Huang, Abhishek Singh, and Narendra Ahuja, "Single image super-resolution from transformed self-exemplars," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 5197–5206.

[18] Eirikur Agustsson and Radu Timofte, "NTIRE 2017 challenge on single image super-resolution: Dataset and study," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2017, pp. 126–135.

[19] Jingyun Liang, Andreas Lugmayr, Kai Zhang, Martin Danelljan, Luc Van Gool, and Radu Timofte, "Hierarchical conditional flow: A unified framework for image super-resolution and image rescaling," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 4076–4085.