

# IMAGE SUPER-RESOLUTION VIA DEEP DILATED CONVOLUTIONAL NETWORKS

Zehao Huang<sup>1,2</sup>, Lingfeng Wang<sup>1,2</sup>, Gaofeng Meng<sup>2</sup> and Chunhong Pan<sup>2</sup>

1. Hunan Provincial Key Laboratory of Network Investigational Technology, Hunan Police Academy

2. NLPR, Institute of Automation, Chinese Academy of Sciences

{zehaohuang18@gmail.com, lfwang@nlpr.ia.ac.cn, gfmeng@nlpr.ia.ac.cn, and chpan@nlpr.ia.ac.cn}

## ABSTRACT

Deep learning techniques have been successfully applied in single image super-resolution (SR). Recently, researches have shown that increasing the depth of network can significantly improve SR performance. Very deep networks for SR achieved a large improvement than former methods. However, simply increasing depths basically introduce more parameters and this lead to cumbersome computational cost. In this paper, we present a general and effective method to accelerate very deep networks for single image SR. Our method is based on dilated convolution operation, which support exponential expansion of the receptive field without increasing filter size. With the help of dilated convolution, shallow networks can achieve large receptive field and exploit contextual information in an efficient way. Based on a very deep network, we propose a 12 layers dilated convolutional network for SR (DCNSR). While accelerating 2x speed, our shallow network achieves better performance than original deep networks and shows state-of-the-art reconstructed results.

**Index Terms**— Super-Resolution, Deep Networks, Dilated Convolution, Acceleration

## 1. INTRODUCTION

As a classical problem in computer vision, single image super-resolution (SR) aims at recovering a visually pleasing high-resolution (HR) image from a given low-resolution (LR) one. Since multiple HR image patches could correspond to the same LR image patches, SR is an inherently ambiguous and highly ill-posed problem. To address this problem, a large number of single image SR methods have been proposed, such as interpolation-based methods [1, 2, 3], reconstruction-based methods [4, 5, 6] and learning-based methods [7, 8, 9, 10, 11, 12, 13, 14, 15]. Among them, learning-based methods have attracted more attention from the community recently. Through learning a mapping function from corresponding pairs of LR-HR image patches, learning-based methods delivered superior performance in image SR. More recently, inspired by the great success achieved by deep learning, deep convolutional neural networks (CNNs) have been successfully used for image SR and obtained large improvements in accuracy [9, 10, 11, 12, 13, 14].

To the best of our knowledge, Dong *et al.* [9] first proposed super-resolution convolutional neural network (SRCNN) and demonstrated that CNNs can be used to learn a mapping from LR to HR space in an end-to-end manner. Lately, in order to investigate whether domain expertise can be used to design better deep network

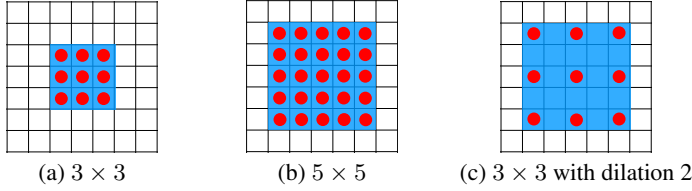
architectures, Wang *et al.* [10] proposed sparse coding based network (SCN). Based on a set of sparse coding sub-networks, they successfully added sparse prior into deep network and their SCN outperforms SRCNN with a smaller model size. Furthermore, deeper structures have been explored in [11, 12]. Inspired by VGG-net used for ImageNet classification [16], Kim *et al.* [11] proposed a 20 weight layers CNN termed VDSR and obtained a significant improvement in accuracy. Benefited from larger image contextual information and stronger learning capacity, deeper network significantly boosts performance in image SR. However, increasing depth basically introduces more parameters and this will cause two problems. First, expensive computation will slow down the speed of reconstruction. Second, more training data is required to prevent over-fitting. To handle the problem of increasing parameters while keeping performance, Kim *et al.* [12] used a deeply-recursive convolutional network (DRCN). Nonetheless, the computation of DRCN is as costly as VDSR. Therefore, the speed of SR is still slow for real time application.

Recently, there are a number of CNN acceleration studies. However, the exploration of accelerating deep networks for SR is very limited. The general procedure of SR consists of two decoupled steps: upsampling the original LR image to the desired size using bicubic interpolation; taking the upscaled LR image as input and reconstructing the HR image by learned mapping function. Because of the first step, the computation cost grows quadratically with the spatial size of the HR image. Thus, handling the upsampling operation at the end of network is an effective way for acceleration. Inspired by this idea, Shi *et al.* [13] proposed a sub-pixel convolution layer to store the information of HR image into smaller size feature maps with multiple channels. Similarly, Dong *et al.* [14] adopted deconvolutional layer to replace the bicubic interpolation. Besides, they re-design the original SRCNN structure into a more compact one. Their new fast SRCNN achieves a significant speed up. However, much domain knowledge and cumbersome experiments are needed for designing a faster architecture. A simple and adaptive method to accelerate standard CNN networks for image SR is still absent.

In this work, we propose a simple yet effective way to accelerate very deep convolutional networks for SR without losing performance. In [11], they compared the SR performance of different networks with depth ranging from 5 to 20 and indicated that deeper networks produced better performance because of large receptive field and high nonlinearities. Besides this, we further argue that the size of receptive field is more important than the depth of network. In SR problem, high receptive field means more contextual information used for reconstruction. While holding the same size of receptive field, networks with different depths will produce similar reconstructed results. So we introduce dilated convolution operation into the framework of deep learning based SR methods. Compared to normal convolution, dilated convolution support exponential expansion

This work is supported by the National Natural Science Foundation of China (Grant No. 61403376 and 61370039), the Beijing Nature Science Foundation (Grant No. 4162064) and the Open Research Fund of Hunan Provincial Key Laboratory of Network Investigational Technology (Grant No. 2015HNWLFZ055).

sion of the receptive field without increasing filter size. Therefore, shallower networks with less parameters can obtain the same size of receptive field as very deep networks. With the help of dilated convolution, we show a 12 layers CNN can produce better results than VDSR, with nearly half parameters and computations.



**Fig. 1.** Dilated convolution operation supports larger receptive field than traditional convolution. The receptive field of (c) is the same as (b), but filter size in (c) is only  $3 \times 3$ .

Specifically, the details and contributions of this work are mainly in three aspects:

- We firstly introduce dilated convolution into deep learning based SR methods and proposed a state-of-the-art dilated convolutional network for SR (DCNSR).
- We demonstrate receptive field is a significant important factor in SR task. While keeping the same size of receptive field, networks with different depths will produce similar results.
- With the help of dilated convolution, DCNSR yields better performance with less parameters and faster speed. The idea of network design can be easily applied into existing SR architectures and this strategy benefits other accelerating approaches.

In the following, we will first review the related work of VDSR. Then, in Section 3 we will describe dilated convolution in detail and present the new DCNSR. Section 4 shows our implementation details and experiments, in which we compare the performance of our method with the state-of-the-art approaches. Finally, we conclude this paper in Section 5.

## 2. REVIEW OF VERY DEEP CONVOLUTIONAL NETWORKS FOR SR

Inspired by the successful of deep and thin network in image classification task [16], Kim *et al.* [11] designed a very deep convolutional networks for SR (VDSR). Compared to 3 weight layers in SRCNN [9], VDSR used 20 layers and achieved significant improvement. All layers except the first and the last are of the same type: 64 filters of the size  $3 \times 3 \times 64$ . Both of the first and last layer consisted of a single filter of size  $3 \times 3 \times 64$ . Rectified linear unit (ReLU) is used as activation function. The receptive field size of VDSR is  $41 \times 41$ , which is much larger than SRCNN ( $13 \times 13$  in SRCNN). With larger receptive field size, VDSR can use more context information to predict image details.

In addition, for speeding up training convergence, they suggested a residual learning network structure. Instead of reconstructing HR results directly, they used VDSR to learn the residual of HR and LR images. The advantage of residual learning can be explained in two sides. For SR problem, LR images and HR images are largely similar. Modelling the difference between HR and LR images can reserve the information of LR images more effectively. This idea is

similar to former classical SR methods [7] which aimed at recovering the high-frequency information of HR images. For deep neural networks, residual learning make very deep networks easier to optimize [17]. With the help of this strategy and high learning rate, they trained VDSR over 80 epochs and achieved a new state-of-the-art performance in single image SR. Unfortunately, increasing depth basically introduces more parameters and this will cause to slow down the speed of reconstruction. In this work, we mainly focus on this problem, and propose the new DCNSR model.

## 3. THE PROPOSED MODEL

In this section, we first introduce dilated convolution operation. Then, we describe the importance of receptive field in CNN based single image SR. Due to these analyses, we proposed DCNSR by introducing dilated convolution operation into VDSR, and ensure the receptive field is same with VDSR.

### 3.1. Dilated Convolution

Let  $f : \mathbb{Z}^2 \rightarrow \mathbb{R}$  be a discrete function. Let  $\Omega_r = [-r, r]^2 \cap \mathbb{Z}^2$  and let  $k : \Omega_r \rightarrow \mathbb{R}$  be a discrete filter of size  $(2r+1)^2$ . The discrete (full) convolution operator  $*$  can be defined as

$$(f * k)(\mathbf{p}) = \sum_{\mathbf{s}+\mathbf{t}=\mathbf{p}} f(\mathbf{s})k(\mathbf{t}). \quad (1)$$

Recently, Yu *et al.* [18] generalized this operator and proposed a new one named dilated convolution:

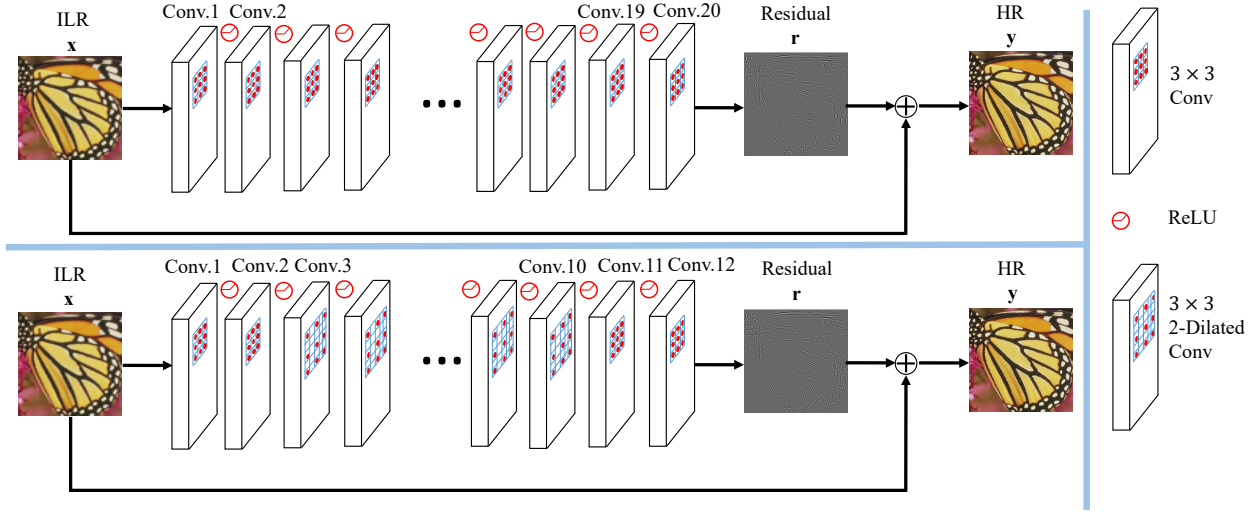
$$(f *_l k)(\mathbf{p}) = \sum_{\mathbf{s}+l\mathbf{t}=\mathbf{p}} f(\mathbf{s})k(\mathbf{t}). \quad (2)$$

where  $l$  is a dilation factor and  $*_l$  is referred as a dilated convolution or an  $l$ -dilated convolution. The familiar full convolution  $*$  can be regarded as a specific version of  $l$ -dilated convolution when  $l = 1$ .

Compared to traditional convolution, the dilated convolution operation can apply the same filter at different ranges using different dilation factors. Thus, it can support exponential expansion of the receptive field without increasing filters size or network depth. For example, Fig. 1 illustrates the receptive field of different filter size and dilation. As shown in this figure, the receptive field of 2-dilated convolution is same with  $5 \times 5$  convolution, while its parameters are same with the  $3 \times 3$  convolution.

### 3.2. The Importance of Receptive Field

Yu *et al.* [18] had demonstrated the effectiveness of dilated convolution for dense prediction. Similarly, we argue that image SR can benefit from dilated convolution operation since receptive field is also an important factor in SR problem. In the task of SR, the size of receptive field means the amount of contextual information that can be exploited to infer high-frequency components. As SR is a highly ill-posed problem, collecting and analyzing more context can afford the network more clues to predict image details. Thus, we consider networks with different depths but the same size of receptive field will produce similar reconstructed results. There are three strategies to raise the size of receptive field: (1) holding on filters size, adding more weight layers; (2) using large filters and (3) replacing traditional convolution by  $l$ -dilated convolution, where  $l > 1$ . In order to show the importance of receptive field and the effectiveness of dilated convolution, strategy (2) and (3) are used to get large receptive field. While holding the same size of receptive field, we compare the performances of different network settings.



**Fig. 2.** Network architectures of VDSR (top) and VDSR\_12.Dilated (bottom). VDSR has 20 layers, and VDSR\_12.Dilated has 12 layers.

**Table 1.** Network architecture of VDSR\_12.Dilated. Benefiting from dilated convolution operation, VDSR\_12.Dilated can achieve  $41 \times 41$  receptive field with 12 layers and  $3 \times 3$  filters.

Layer	1	2	3	4	5	6	7	8	9	10	11	12
Convolution	$3 \times 3$	$3 \times 3$	$3 \times 3$	$3 \times 3$	$3 \times 3$	$3 \times 3$	$3 \times 3$	$3 \times 3$	$3 \times 3$	$3 \times 3$	$3 \times 3$	$3 \times 3$
Dilation	1	1	2	2	2	2	2	2	2	2	1	1
Receptive Filed	$3 \times 3$	$5 \times 5$	$9 \times 9$	$13 \times 13$	$17 \times 17$	$21 \times 21$	$25 \times 25$	$29 \times 29$	$33 \times 33$	$37 \times 37$	$39 \times 39$	$41 \times 41$
Output Channels	64	64	64	64	64	64	64	64	64	64	64	1

### 3.3. DCNSR Model

VDSR achieves the new state-of-the-art and it does not need specific network structure design, therefore, we choose it as our baseline network. Holding on the same size of receptive field as VDSR ( $41 \times 41$ ), we design five networks with different filters size or dilations for comparison. (1) VDSR\_20 is the baseline 20 layers network from [11]. The filter size of all layers is  $3 \times 3$ . (2) VDSR\_12 is a 12 layers network with filter size  $3 \times 3$  and  $5 \times 5$ . (3) VDSR\_12.Dilated is also a 12 layers network but with filter size  $3 \times 3$ . All layers with filter size  $5 \times 5$  in (2) are replaced by  $3 \times 3$  filters and 2-dilated convolution. (4) VDSR\_10 and (5) VDSR\_10.Dilated are similar to (2) and (3) respectively. Fig. 2 illustrates the difference between VDSR and VDSR\_12.Dilated in detail. Table 1 explains how a 12 layers network achieves  $41 \times 41$  receptive filed.

## 4. EXPERIMENTAL RESULTS

In this section, we evaluate the performance of our method on several datasets. We first describe datasets used for training and testing. Next, we describe implementation details. Finally, we evaluate different dilated convolution based SR models, and compare our method with several state-of-the-art SR methods.

### 4.1. Datasets

**Training Dataset:** The 91-image dataset from Yang *et al.* [7] is widely used as the training set in learning-based SR methods. However, 91 images are not enough for deep network to gain best performance because of the increasing parameters and complexity of deep model. Therefore, following Kim *et al.* [11], we used 200 images

from Berkeley Segmentation Dataset [19] additionally. Thus 291 images are used for benchmark with other methods. Data augmentation technique is used to getting more training data. Scale augmentation used in [11] is also adopted in our training. So we can handle multiple scales SR task with only one model.

Notedly, for results in Section 4.3, we only used 91 images to train networks fast, so performances can be slightly different between Section 4.3 and Section 4.4.

**Testing Dataset:** For a fair comparison, we use Set5 [20], Set14 [21] and BSD100 [19] for testing. The original images are first down-sampled by bicubic interpolation and then up-sampled to desired size to generate LR-HR image pairs for both training and evaluation. Zero padding is applied in all convolutional layers to keep the size of output maps as the same as input.

### 4.2. Implementation Details

The network structure and parameters setting are described in Table 2. We use Adam [22] with a mini-batch size of 64.  $\beta_1$  and weight decay are set to 0.9 and 0.0001, respectively. Following [11], we initialize the weights by the method described in He *et al.* [23]. All experiments are trained over 80 epochs (11698 iterations with batch size 64) with a learning rate of  $10^{-4}$ . Gradient clipping used in [11] is not necessary in our training since we use Adam instead of SGD.

In testing, we only process the luminance channel with our method. After reconstructing, we shave the image border in the same way as [9] for objective evaluations to ensure fair comparison. In our implementation, all the experiments are implemented using the caffe package [24] on a GTX 980Ti GPU. Since dilated convolution operation has not been supported by CuDNN [25], all the

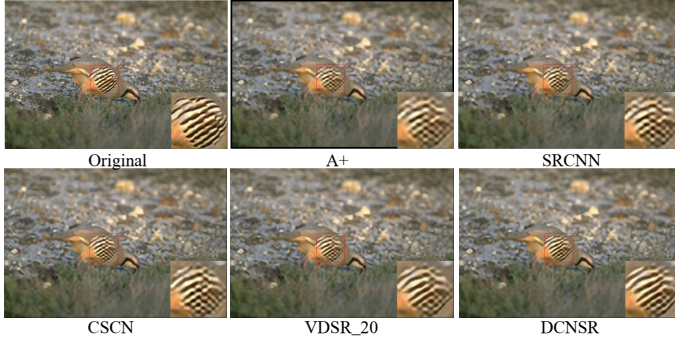
**Table 2.** Performance of different network settings. All the five networks have  $41 \times 41$  receptive field and scale factor is 2.

	VDSR_20	VDSR_12	VDSR_12_Dilated	VDSR_10	VDSR_10_Dilated
First Part	Conv(3,64,1,D1)	Conv(3,64,1,D1)	Conv(3,64,1,D1)	Conv(3,64,1,D1)	Conv(3,64,1,D1)
Mid Part	18Conv(3,64,64,D1)	Conv(3,64,64,D1)- 8Conv(5,64,64,D1)- -Conv(3,64,64,D1)	Conv(3,64,64,D1)- 8Conv(3,64,64,D2)- -Conv(3,64,64,D1)	Conv(3,64,64,D1)- Conv(5,64,64,D1)- 4Conv(7,64,64,D1)- -Conv(5,64,64,D1)- -Conv(3,64,64,D1)	Conv(3,64,64,D1)- Conv(3,64,64,D2)- 4Conv(3,64,64,D3)- -Conv(3,64,64,D2)- -Conv(3,64,64,D1)
Last Part	Conv(3,1,64,D1)	Conv(3,1,64,D1)	Conv(3,1,64,D1)	Conv(3,1,64,D1)	Conv(3,1,64,D1)
PSNR(Set5)	<b>37.46</b> dB	37.45 dB	37.44 dB	37.40 dB	37.38 dB
PSNR(Set14)	32.83 dB	32.80 dB	<b>32.87</b> dB	32.65 dB	32.80 dB
PSNR(BSD100)	31.65 dB	31.62 dB	<b>31.72</b> dB	31.45 dB	31.65 dB
Parameters	665921	895873	370368	1082496	296064
Speedup	1.0×	1/13.2×	2.0×	1/10.4×	2.6×

**Table 3.** PSNR and SSIM comparison on three test datasets among different SR methods.

Dataset	Scale	Bicubic PSNR/SSIM	A+ PSNR/SSIM	SRCNN PSNR/SSIM	CSCN PSNR/SSIM	VDSR_20 PSNR/SSIM	DCNSR PSNR/SSIM
Set5	×2	33.66/0.9299	36.57/0.9545	36.48/0.9542	36.93/0.9552	<b>37.58/0.9591</b>	37.46/0.9585
	×3	30.39/0.8682	32.67/0.9093	32.57/0.9090	33.10/0.9144	33.68/0.9218	<b>33.74/0.9219</b>
	×4	28.42/0.8104	30.36/0.8617	30.31/0.8628	30.86/0.8732	33.33/0.8828	<b>31.37/0.8831</b>
Set14	×2	30.23/0.8688	32.47/0.9063	32.45/0.9067	32.56/0.9074	<b>33.00/0.9125</b>	32.91/0.9116
	×3	27.54/0.7742	29.29/0.8203	29.30/0.8215	29.41/0.8231	29.75/0.8305	<b>29.76/0.8312</b>
	×4	26.00/0.7027	27.47/0.7514	27.50/0.7513	27.64/0.7578	27.95/0.7647	<b>27.99/0.7661</b>
BSD100	×2	29.56/0.8431	30.77/0.8756	31.36/0.8879	31.40/0.8884	<b>31.86/0.8956</b>	31.81/0.8947
	×3	27.21/0.7385	28.18/0.7791	28.41/0.7863	28.50/0.7875	<b>28.80/0.7964</b>	<b>28.80/0.7972</b>
	×4	25.96/0.6675	26.74/0.7065	26.90/0.7103	27.03/0.7161	27.24/0.7230	<b>27.26/0.7241</b>

experiments are tested without CuDNN acceleration for comparing pure computational loads.

**Fig. 3.** The ‘8023’ image from BSD100 dataset ( $4\times$  upscaling).

#### 4.3. Evaluation of Different Dilated Convolution Based Models

In Table 2, we give the architectures and performances of these 5 networks. All these networks expect VDSR\_10 achieve similar PSNR performance. This is corresponding to our assumption that different networks with the same size of receptive field will produce similar results. With the same depths, dilated networks (3) and (5) show better performance, less parameters and faster speed than (2) and (4). This is caused by the drawback of large filter size. While filters with bigger size can obtain large receptive field, they also bring more noise into the network learning procedure. In addition, compared to  $3 \times 3$  filters, there are much redundancy in learned  $5 \times 5$  and  $7 \times 7$  kernels because the values in these big kernels are highly correlated. Lastly, we find that the speed of (2) and (4) are much

slower than (3) and (5), even slower than (1). To sum up, receptive field is an important factor in SR task and dilated convolution is a better technique for achieving large receptive field.

#### 4.4. Comparisons with State-of-the-Art Methods

Since VDSR\_12\_Dilated achieves the best result with fast speed, we compare VDSR\_12\_Dilated with other state-of-the-art SR methods and named it as dilated convolutional network for SR (DCN-SR). Compared methods are A+[8], SRCNN [26], CSCN [10] and our VDSR\_20 implementation. The implementations are all from the publicly available codes provided by the authors. In Table 3, we provide a summary of quantitative evaluation on testing datasets. Our DCNSR yields the highest average PSNR and SSIM in all these datasets. In Fig. 3, the reconstructed images of our DCNSR is much sharper and clearer than other results. Reconstructed results obtained with DCNSR are available online for all three datasets<sup>1</sup>.

## 5. CONCLUSION

In this paper, we introduce dilated convolution to accelerate the speed of very deep networks for SR. We first show that receptive field is an important factor in image SR. Networks with different depths but the same receptive field will produce similar HR results. Second, we propose dilated convolution to replace full convolution operation. Dilated convolution operation is a much better technique for gathering large receptive field. Based on a 20 layers very deep network, we design five different networks setting and show the effectiveness of dilated convolution operation both in performance and speed. Without specific network design, our strategy is a general method for deep learning based SR acceleration and it benefits other accelerating approaches.

<sup>1</sup><https://drive.google.com/open?id=0ByMcIJq3Oj8peGplamJQNjM1TIU>

## 6. REFERENCES

- [1] Philippe Thévenaz, Thierry Blu, and Michael Unser, “Image interpolation and resampling,” *Handbook of medical imaging, processing and analysis*, pp. 393–420, 2000.
- [2] Lingfeng Wang, Shiming Xiang, Gaofeng Meng, Huaiyu Wu, and Chunhong Pan, “Edge-directed single-image super-resolution via adaptive gradient magnitude self-interpolation,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 23, no. 8, pp. 1289–1299, 2013.
- [3] Lingfeng Wang, Huaiyu Wu, and Chunhong Pan, “Fast image upsampling via the displacement field,” *IEEE Transactions on Image Processing*, vol. 23, no. 12, pp. 5123–5135, 2014.
- [4] Hong Chang, Dit-Yan Yeung, and Yimin Xiong, “Super-resolution through neighbor embedding,” in *CVPR*, 2004.
- [5] Daniel Glasner, Shai Bagon, and Michal Irani, “Super-resolution from a single image,” in *ICCV*, 2009.
- [6] Matan Protter, Michael Elad, Hiroyuki Takeda, and Peyman Milanfar, “Generalizing the nonlocal-means to super-resolution reconstruction,” *IEEE Transactions on Image Processing*, vol. 18, no. 1, pp. 36–51, 2009.
- [7] Jianchao Yang, John Wright, Thomas S Huang, and Yi Ma, “Image super-resolution via sparse representation,” *IEEE Transactions on Image Processing*, vol. 19, no. 11, pp. 2861–2873, 2010.
- [8] Radu Timofte, Vincent De Smet, and Luc Van Gool, “A+: Adjusted anchored neighborhood regression for fast super-resolution,” in *ACCV*, 2014.
- [9] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang, “Learning a deep convolutional network for image super-resolution,” in *ECCV*, 2014.
- [10] Zhaowen Wang, Ding Liu, Jianchao Yang, Wei Han, and Thomas Huang, “Deep networks for image super-resolution with sparse prior,” in *ICCV*, 2015.
- [11] Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee, “Accurate image super-resolution using very deep convolutional networks,” in *CVPR*, 2016.
- [12] Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee, “Deeply-recursive convolutional network for image super-resolution,” in *CVPR*, 2016.
- [13] Wenzhe Shi, Jose Caballero, Ferenc Huszar, Johannes Totz, Andrew P. Aitken, Rob Bishop, Daniel Rueckert, and Zehan Wang, “Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network,” in *CVPR*, 2016.
- [14] Chao Dong, Chen Change Loy, and Xiaoou Tang, “Accelerating the super-resolution convolutional neural network,” in *ECCV*, 2016.
- [15] Lingfeng Wang, Zehao Huang, Yongchao Gong, and Chunhong Pan, “Ensemble based deep networks for image super-resolution,” *Pattern Recognition*, vol. 68, pp. 191–198, 2017.
- [16] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” in *ICLR*, 2015.
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, “Deep residual learning for image recognition,” in *CVPR*, 2016.
- [18] Fisher Yu and Vladlen Koltun, “Multi-scale context aggregation by dilated convolutions,” in *ICLR*, 2016.
- [19] 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 *ICCV*, 2001.
- [20] Marco Bevilacqua, Aline Roumy, Christine Guillemot, and Marie Line Alberi-Morel, “Low-complexity single-image super-resolution based on nonnegative neighbor embedding,” in *BMVC*, 2012.
- [21] Roman Zeyde, Michael Elad, and Matan Protter, “On single image scale-up using sparse-representations,” in *ICCS*, 2010.
- [22] Diederik Kingma and Jimmy Ba, “Adam: A method for stochastic optimization,” in *ICLR*, 2015.
- [23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, “Delving deep into rectifiers: Surpassing human-level performance on imagenet classification,” in *ICCV*, 2015.
- [24] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell, “Caffe: Convolutional architecture for fast feature embedding,” in *ACM MM*, 2014.
- [25] Sharan Chetlur, Cliff Woolley, Philippe Vandermersch, Jonathan Cohen, John Tran, Bryan Catanzaro, and Evan Shelhamer, “cudnn: Efficient primitives for deep learning,” *arXiv preprint arXiv:1410.0759*, 2014.
- [26] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang, “Image super-resolution using deep convolutional networks,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 38, no. 2, pp. 295–307, 2016.