

OverNet: Lightweight Multi-Scale Super-Resolution with Overscaling Network

Parichehr Behjati
Computer vision center
Barcelona, Catalonia Spain
pbehjati@cvc.uab.es

Pau Rodríguez
Element AI
Montreal, Canada

Armin Mehri
Computer Vision Center
Barcelona, Catalonia Spain

Isabelle Hupont
Herta Security
Barcelona, Catalonia Spain

Carles Fernández Tena
Herta Security
Barcelona, Catalonia Spain

Jordi González
Computer Vision Center
Barcelona, Catalonia Spain

Abstract

Super-resolution (SR) has achieved great success due to the development of deep convolutional neural networks (CNNs). However, as the depth and width of the networks increase, CNN-based SR methods have been faced with the challenge of computational complexity in practice. Moreover, most SR methods train a dedicated model for each target resolution, losing generality and increasing memory requirements. To address these limitations we introduce OverNet, a deep but lightweight convolutional network to solve SISR at arbitrary scale factors with a single model. We make the following contributions: first, we introduce a lightweight feature extractor that enforces efficient reuse of information through a novel recursive structure of skip and dense connections. Second, to maximize the performance of the feature extractor, we propose a model agnostic reconstruction module that generates accurate high-resolution images from overscaled feature maps obtained from any SR architecture. Third, we introduce a multi-scale loss function to achieve generalization across scales. Experiments show that our proposal outperforms previous state-of-the-art approaches in standard benchmarks, while maintaining relatively low computation and memory requirements.

1. Introduction

Single image super-resolution (SISR) is the task of reconstructing a high resolution image (HR) from its low-resolution version (LR). As obtaining a HR image from LR is an ill-posed problem, the model needs to *learn* the original data distribution to produce the most likely solutions.

Convolutional neural networks (CNNs) have recently become the main workhorse to tackle SISR [5]. Thanks to the increase in capacity of CNNs in depth [23] and width [36],

their performance has greatly improved. Despite their remarkable performance, most deep networks still have some drawbacks. Firstly, increase in depth and width has also raised computational demands and memory consumption. This makes modern architectures less applicable in practice, such as in mobile and embedded applications. Secondly, as the network depth increases, low-level feature information gradually disappears in the successive non-linear operations to produce the output. However, these low-level features are crucial for the network to reconstruct high quality images.

Aside from the aforementioned problems, another desired ability is to upsample images to arbitrary scales using a single model. Current state-of-the-art SISR models such as RDN [41], ESPCNN [28] and EDSR [23], only consider SR at certain integer scale factors ($\times 2$, $\times 3$, $\times 4$) and treat each super-resolution scale as an independent task. They then train a different specialized model for each, which is not practical for mobile applications.

To address these problems, we propose Overscaling Network (OverNet), a novel lightweight method for SISR. OverNet consists of two main parts: a lightweight feature extractor and an Overscaling module (OSM) for reconstruction. The feature extractor follows a novel recursive framework of skip and dense connections to reduce low-level feature degradation. The OSM is a new inductive bias which generates accurate SR image by internally constructing an overscaled intermediate representation of the output features. Finally, to solve the problem of reconstruction at arbitrary scale factors, we introduce a novel multi-scale loss by downsampling the output at multiple super resolution factors and we minimize the reconstruction error in all of them. Our main contributions can be summarized as follows:

- A lightweight recursive feature extractor, which results in improved performance over state-of-the-art models, even those having an order of magnitude more parameters.

- An Overscaling Module (OSM) that generates overscaled maps from which HR images can be accurately recovered at arbitrary scales. This module boosts the reconstruction accuracy efficiently with respect to its number of parameters. Additionally, we demonstrate that integrating this module into existing state-of-the-art models improves on their original performance.
- A novel multi-scale loss function for SISR, that allows the simultaneous training of all scale factors using a single model. As a result, the model is able to maintain accurate reconstruction results across scales.

2. Related Work

Recently, deep learning models have dramatically improved the SISR task. Dong et al. [5] first presented SRCNN, a CNN to predict super-resolved images. SRCNN has a large number of operations compared to its depth, since the network operates by initially upsampling LR images and subsequently refining them. In contrast to the SRCNN, FSCRNN [6] and ESPCN [28] only upsample images at the output of the network, which leads to a reduction in the number of operations compared to SRCNN.

Despite the higher capacity of deep neural networks, the aforementioned methods have settled for shallow models because of the difficulty in training. VDSR [15] and IRCNN [38] improved the performance by increasing the network depth, using stacked convolutions with residual connections. Lim et al. [23] further expanded the network size and improved the residual block by removing batch normalization layers. Ahn et al. [1] proposed a cascading residual network using ResNet blocks [9] to learn the relationship between LR input and HR output. Later, Ledig et al. [20] introduced the SRResNet and further improved in [34] and [33] by introducing dense connections. More recently, Zhang et al. [41] and Liu et al. [24] also used dense and residual connections in RDN and RFANet to utilize information from all the feature hierarchy. DBPN [8] and SRFBN [22] architectures comprise of a series of up and down sampling layers densely connected with each other. These methods achieved significant improvement over conventional SR methods and indicate the effectiveness of residual learning.

Another issue of deep learning-based SR is how to reduce the parameters and number of operations to make it effective in mobile applications. For instance, DRCN [16] was the first to apply recursive algorithm to SISR to reduce the number of parameters by reusing them multiple times. Tai et al. [31] improved DRCN by combining the recursive and residual network schemes in order to achieve better performance with even fewer parameters. They also introduced a deep memory network to solve the problem of long-term dependencies [32]. On the other hand, Lap-

SRN [18] employs a pyramidal framework to increase the image size gradually. By doing so, LapSRN effectively performs SR on extremely low-resolution cases. More recently, Muqet et al. [25] proposed stacked multi-attention blocks to further improve the performance. However, these methods use very deep networks to compensate for the loss of parameters and hence, they require heavy computing resources. Therefore, we focus on developing a lightweight model to maximize the performance of existing networks as well as minimize their computational cost.

One of the most important stages of SISR is reconstruction, which consists of generating HR images based on high-level features extracted from a low-dimensional space. Interpolation is a commonly used method in SR networks, such as SRCNN [5], VDSR [15] and DRRN [31], to resize the LR image to the target size as the input of a CNN model for SR reconstruction. However, computational operations are greatly increased due to the large input image size. Thus, FSRCNN [6] and SRDenseNet [20] directly adopted the LR image as input, in which a transposed convolution layer was added to implement the final upsampling reconstruction [33]. This method greatly reduces unnecessary computational overhead. Furthermore, EPSCN [28] proposed a method called pixelshuffle [2] to overcome the problem of the checkerboard effect in transposed convolution. Pixelshuffle has been widely used in recent SR models, such as EDSR [23], WDSR [36] and RCAN [40]. However, these methods cannot manage multi-scale training.

Few works tackle SR at different scale factors, and those that do treat the problem as independent tasks, i.e. a model is trained for each scale. Lim et al. [23] proposed the first multi-scale SR model, which has different image processing blocks and upsampling modules for each integer scale factor. Later, Li et al. [21] proposed a multi-scale residual network. They use multi-path convolution layers with different kernel sizes to extract multi-scale spatial features. Grm et al. [7] proposed to upsample the image progressively by $\times 2$ using a series of so-called SR modules and compute the loss of generated SR results by each module. Thus, these methods require vast amounts of computational resources. Recently, Meta-SR [12] introduced an upsampling module based on meta-learning to solve SR at arbitrary scale factors with a single model through a weight prediction technique. However, this method must predict a large number of convolution weights for each target pixel, the prediction is inefficient, and the results may be unstable [35].

3. Proposed Overscaling Network

This section describes the main components of our architecture as shown in Figure 1, and the novel loss function.

Problem formulation. Algorithm 1 formulates the main pipeline steps. Given a set of HR images and their down-scaled versions $\{I^{HR}, I^{LR}\}$, the goal of SISR is to find a

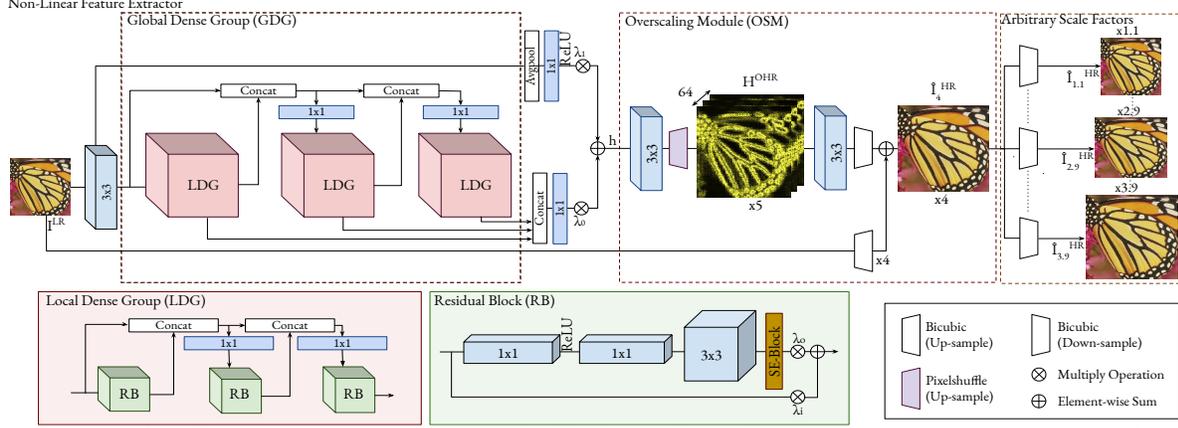


Figure 1: Demonstration of our proposed overscaling network with short and long skip connections. As the maximum scale factor in this particular example is set to $N = 4$, the required overscaling is $\times 5$.

Algorithm 1 Overscaling network forward step. Given a LR image and a set of output scales, OverNet produces an HR reconstruction for each scale. Learnable parameters are omitted to improve readability.

```

function OVERNET(LR image  $I^{LR}$ , target scales  $S$ )
  # Compute features with the CNN
   $\mathbf{h} = \mathcal{H}(I^{LR})$ 
  # Overscaling module
   $\hat{I}^{HR} = \mathcal{O}(\mathbf{h})$ 
  # Output
  for  $s$  in  $S$  do
     $\hat{I}_s^{HR} = \text{bicubic}_{\downarrow}(\hat{I}^{HR}, \text{scale} = s)$ 
  end for
  return  $\{\hat{I}_s^{HR}, s \in S\}$ 
end function

```

function $\mathcal{F} : LR \rightarrow HR$ that maps LR images to their original HR version. The problem is ill-posed, since there are multiple possible HR images corresponding to a single LR image. However, it is possible to *learn* the most likely reconstruction by parametrizing \mathcal{F} over a set of parameters θ , and finding the most likely θ given some criterion \mathcal{L} :

$$\theta^* = \arg \min_{\theta} \sum \mathcal{L}(\mathcal{F}(I^{LR}, \theta), I^{HR}) \quad (1)$$

We chose \mathcal{L} to be the $L1$ distance, since we empirically obtained superior PSNR results compared to $L2$. In this work \mathcal{F} is composed of two parts: (i) a feature extractor \mathcal{H} :

$$\mathbf{h} = \mathcal{H}(I^{LR}, \theta_h) \quad (2)$$

with parameters θ_h , and (ii) the overscaling module \mathcal{O} :

$$\hat{I}^{HR} = \mathcal{O}(\mathbf{h}, \theta_o) \quad (3)$$

with θ_o the parameters used in this operation, and \hat{I}^{HR} the reconstructed image. These two parts are described next.

3.1. Feature Extractor

The feature extractor computes useful representations of the LR patch in order to infer its HR version. Concretely, we propose a recursive structure based on Residual Blocks (RBs) assembled into Local Dense Groups (LDGs) and LDGs into Global Dense Group (GDG), see Figure 1.

Residual blocks. We use a modified version of WDSR [36] with wide low-rank convolutions instead of using standard residual blocks [41]. These convolutions widen the activation space before the non-linearity to let more information pass through it and lose less detail, while using the same amount of computation as standard 3×3 residual blocks. In order to make the network focus on more informative features, we exploit the inter-dependencies among feature channels using squeeze-and-excitation (SE) operations [14] after these convolutions, see Figure 1.

Inspired by [29, 30], the model learns a scalar multiplier λ to balance the amount of information that should be carried by the identity and activation operations within the residual blocks (RBs) of the network. Let \mathbf{x}_i and \mathbf{x}_o be the input and output vectors of the k -th RB, and WA the wide activation operation [36]. Then, the RB proceeds as:

$$\mathbf{x}_o = \lambda_o SE(WA(\mathbf{x}_i)) + \lambda_i \mathbf{x}_i \quad (4)$$

Local and global dense groups. RBs are grouped into the so-called Local Dense Groups (LDGs). The input of a RB is concatenated with the output of the all the previous RBs in the group and merged with a 1×1 convolution. This recursion is repeated for all RBs within the LDG. In this way, we gather all local information progressively by 1×1 convolution layers.

To increase the network capacity, a similar recursion is applied to the Global Dense Group (GDG), but this time incorporating skip connections between LDGs. We repeat this procedure while integrating the recursive concatenations through the LDGs into a single output. The output

of each LDG is concatenated to the input of the next one. In order to facilitate access to local information, the final output of the network receives the concatenation of the outputs of all the LDGs. Therefore, the model incorporates features from multiple layers. This strategy makes information propagation efficient due to the multi-level representation and many shortcut connections. Inspired by MemNet [32], we then introduce a 1×1 convolutional layer to adaptively merge the output information, as directly using these concatenated features would greatly increase computational complexity. The output of these hierarchical features can be formulated as

$$\mathbf{f}_D = \text{conv}_{1 \times 1}([\mathbf{f}_0, \dots, \mathbf{f}_{D-1}]) \quad (5)$$

where $[\mathbf{f}_0, \dots, \mathbf{f}_{D-1}]$ refers to the concatenation of feature maps produced by LDGs.

To make sure that no information is lost before the reconstruction step, we incorporate a long-range skip connection to grant access to the original information, and encourage back-propagation of gradients from the output of the feature extractor to the first 3×3 convolution layer. We also include a global average pooling followed by a 1×1 convolution, to fully capture channel-wise dependencies from the aggregated information. The final output before the reconstruction step is then,

$$\mathbf{h} = \lambda_0 \mathbf{f}_D + \lambda_1 \sigma(\text{conv}_{1 \times 1}(\text{GAP}(\text{conv}_{3 \times 3}(I^{LR})))) \quad (6)$$

where σ denotes the ReLU activation, GAP denotes global average pooling, and λ_0 and λ_1 are learned parameters.

3.2. Overscaling Module

In this work we introduce a new inductive bias in SISR architectures so as to generate images that are more accurate and present fewer artifacts. We hypothesize that, since overscaling produces multiple values for the same pixel, these values act as an ensemble of predictions thus reducing noise when combined to produce the final image.

Let us consider N the maximum scale factor addressed by the network. We first generate an intermediate representation of the final image consisting of overscaled maps H^{OHR} , with an overscale factor $(N+1)$ times larger. Thus, given the features \mathbf{h} extracted from I^{LR} , we use a 3×3 convolutional layer followed by the strided sub-pixel convolution proposed in [2] to upscale the features \mathbf{h} to H^{OHR} :

$$H^{OHR} = \text{pixelshuffle}(\text{conv}_{3 \times 3}(\mathbf{h})) \quad (7)$$

To obtain the final output of the overscaling module, we further include a second long-range skip connection from the original I^{LR} image. The final HR image is obtained by adjusting the overscaled maps and incorporating them into the naïve upscaling of the original LR image:

$$\hat{I}^{HR} = \text{bicubic}_{\downarrow}(\text{conv}_{3 \times 3}(H^{OHR})) + \text{bicubic}_{\uparrow}(I^{LR}) \quad (8)$$

Hence, we could think of the whole network as learning how to *refine* or *correct* a naïve bicubic upscaling of the low-resolution input, in order to bring it closer to the actual high-resolution counterpart. Since the final \hat{I}^{HR} images are obtained with an efficient non-parametric interpolation, we are able to produce multiple scales with negligible computational cost, and only using differentiable operations.

3.3. Multi-Scale Loss

We propose the minimization of a multi-scale loss to optimize the network. We choose a finite set of scale factors $S = \{s_1 \dots s_n\}$, all within the interval of scales targeted by the network. Once the network has reconstructed the HR image, images at the target scales are obtained through a bank of bicubic interpolators, $\hat{I}_s^{HR} = \text{bicubic}_{\downarrow}(\hat{I}^{HR}, s)$. Then, we minimize the following loss function:

$$\mathcal{L} = \sum_{s \in S} |\hat{I}_s^{HR} - \text{bicubic}_{\downarrow}(I^{HR}, s)| \quad (9)$$

Training with this multi-scale loss at different target scales simultaneously provides additional supervision to the model, compared to a single-scale training. As a result, the model is enforced to learn how to generate highly representative overscaled maps, from which HR images at arbitrary scales can be recovered accurately, hence enforcing the generalization capability of the network across scales.

4. Experimental Results

Datasets and metrics. We use the DIV2K dataset for training, a high-quality image dataset containing 800 images for training, 100 for validation and 100 for testing. Several benchmark datasets are used for testing, namely Set5 [4], Set14 [37], B100 [3], and Urban100 [13]. SR results are evaluated with two commonly used metrics: PSNR (peak-to-peak signal-to-noise ratio) and SSIM (structural similarity index), on the Y channel of the YCbCr space.

Degradation models. To comprehensively illustrate the efficacy of the proposed method, three degradation models are used to simulate LR images, following [38, 39, 41]. The first one, denoted by **BI**, consists of generating LR images by bicubic-downsampling ground truth HR images with $\times 2$, $\times 3$, $\times 4$. The second one, denoted by **DB**, first performs bicubic downsampling on HR images with $\times 3$, and then blurs the images with a Gaussian kernel of size 7×7 and standard deviation 1.6. Finally, we further produce LR images in a third challenging way, denoted by **DN**, by carrying out bicubic downsampling followed by additive Gaussian noise, with noise level of 30.

Implementation details. We denote our original model as OverNet and further introduce OverNet w/o OSM (OverNet without overscaling module). We used 64×64 RGB input patches from the LR images for training. LR patches

Table 1: Effects of skip connections (SCs) in local and global dense groups (LDG, GDG) measured on Urban100 with $\times 3$. The best result is **highlighted**.

Config	1	2	3	4
SCs in LDGs	×	✓	×	✓
SCs in GDG	×	×	✓	✓
#Params	695K	806K	732K	943K
PSNR	28.23	28.21	28.29	28.37

were sampled randomly and augmented with random horizontal flips and 90° rotation. The number of LDGs and RBs was set to 3 in all experiments. We trained our models with the ADAM optimizer [17]. The mini-batch size was set to 64, and the learning rate to the maximum convergent value 10^{-3} , applying weight normalization in all convolutional layers [36]. The learning rate was decreased by half every 2×10^5 back-propagation iterations. We implemented our networks using the PyTorch framework [26] and trained them on NVIDIA 1080 Ti GPUs.

4.1. Ablation Studies

To investigate the performance behaviour of the proposed methods we first show how skip connections inside the proposed local and global dense groups affect the performance of OverNet. Next, we analyze the effect of OSM and the multi-scale loss.

Feature extractor ablation. Table 1 presents the ablation study on the effect of skip connections (SCs) inside the local and global dense groups (LDG, GDG). In this work, SCs contain concatenation and 1×1 convolutions. The small changes in number of parameters between columns is due to the removal of SCs with 1×1 convolutions.

It can be observed that the model which used SCs only in GDG attains better performance than the one without SCs (config 1 which is ResNet+OSM) because the short connections inside the GDG effectively carry the information from intermediate to higher layers. Furthermore, by gathering all features before the upscaling module, the model can better leverage multi-level representations.

On the other hand, as discussed in [10], multiplicative manipulations such as 1×1 convolutions on the shortcut connection can hamper information propagation, and complicate optimization. Similarly, SCs in LDGs behave as shortcut connections inside the residual blocks. Thus, it is natural to expect performance degradation when the global SCs are deactivated. This is because the global SCs ease the information propagation while the local connections are being learned. Therefore, when OverNet uses SCs in both LDGs and GDG, it outperforms all three models.

In detail, information propagates globally via SCs used in GDG, and information flows in the LDGs are fused with the ones that come through global connections. By doing so, information is transmitted by multiple shortcuts and thus

mitigates the vanishing gradient problem: the advantage of multi-level representation is leveraged by the SCs in GDG, which help the information to propagate to higher layers.

Effect of the OSM across scales. Here we analyze the benefits of incorporating the OSM module, and also explore the influence of different interpolation methods on the reconstruction. We run the following experiments: (i) directly using pixelshuffle to generate the images without overscaling feature maps, followed by bicubic interpolation to downscale to arbitrary scales; (ii) downscaling with bilinear interpolation the overscaled feature maps produced by pixelshuffle and (iii) doing the same as (ii) with bicubic interpolation. As shown in Table 2, superior results are achieved by a large margin when the proposed overscaling method is applied. These experiments suggest that, contrary to common practice in the field, the addition of overscaling strongly increases reconstruction accuracy. Best results are achieved using OSM with bicubic interpolation, which in turn yields better results than bilinear.

In addition, we compare our results with Meta-RDN [12], the only method in the literature (to our knowledge) able to carry out SISR at non-integer scales. Meta-RDN is a heavier state-of-the-art model with 22M parameters. For fair comparison, we trained Meta-RDN by replacing its meta-upscale module with OSM (RDN-OSM), while applying their original training settings. RDN-OSM achieves better or comparable performance.

OSM across architectures. The aim of this section is to demonstrate that the benefits of our OSM hold across architectures. To this end, we use state-of-the-art networks including CARN [1], EDSR[23], RDN[41], Meta-RDN [12] and RCAN[40] as references. We replaced their typical upsample modules with our overscaling module (CARN-OSM, EDSR-OSM, RDN-OSM and RCAN-OSM in Table 3 and trained them on DIV2K for all scale factors while applying their original training settings.

It can be observed that all the methods with OSM have higher PSNR than the corresponding baselines at all scale factors. This shows that OSM is robust and orthogonal to the feature extractor chosen, and it increases PSNR.

Generalization across scales. By construction, the overscaling factor in our architecture is always $(N+1)$ when targeting a maximum scale of N , c.f. Section 3.2. The following experiments investigate the generalization capability of models that target a maximum scale N across lower scales $M \leq N$. To this end, we trained models for $N \in \{2, 3, 4\}$ and evaluated them across scales. Table 4 illustrates the experimental results. It can be observed that models trained to target larger scales yield better PSNR scores for all scale factors. This demonstrates the generalization capabilities of the proposed architecture across scales, as it is not necessary to train a dedicated model for each scale. Instead, training a larger scale seems to be always beneficial for lower scales.

Table 2: PSNR results of different OSM upscaling methods trained for arbitrary scales. The test dataset is B100. Best results are **highlighted**, second best underlined.

Experiment	Scale									
	$\times 1.1$	$\times 1.2$	$\times 1.3$	$\times 1.4$	$\times 1.5$	$\times 1.6$	$\times 1.7$	$\times 1.8$	$\times 1.9$	$\times 2.0$
Pixelshuffle	42.40	39.71	38.10	36.75	35.60	34.70	33.96	33.30	33.65	32.22
OSM-bilinear	42.63	39.89	38.15	36.83	35.70	34.78	34.05	33.37	32.76	32.31
OSM-bicubic	42.74	39.95	38.19	36.87	35.74	34.80	34.10	33.42	32.81	32.34
Meta-RDN	<u>42.82</u>	<u>40.40</u>	<u>38.28</u>	<u>36.95</u>	<u>35.86</u>	<u>34.90</u>	<u>34.13</u>	<u>33.45</u>	<u>32.86</u>	<u>32.35</u>
OSM-RDN	42.93	40.48	38.42	37.06	36.01	35.02	34.25	35.53	32.95	32.46
	$\times 2.1$	$\times 2.2$	$\times 2.3$	$\times 2.4$	$\times 2.5$	$\times 2.6$	$\times 2.7$	$\times 2.8$	$\times 2.9$	$\times 3.0$
Pixelshuffle	31.60	31.22	30.75	30.50	30.27	29.95	29.73	29.42	29.17	29.14
OSM-bilinear	<u>31.71</u>	<u>31.29</u>	<u>30.84</u>	<u>30.55</u>	<u>30.37</u>	<u>30.02</u>	<u>29.77</u>	<u>29.52</u>	<u>29.30</u>	<u>29.26</u>
OSM-bicubic	31.75	31.34	30.86	30.65	30.42	30.11	29.83	29.64	29.36	29.30
Meta-RDN	31.82	31.41	31.06	30.62	30.45	30.13	29.82	29.67	29.40	29.30
RDN-OSM	<u>31.75</u>	<u>31.46</u>	<u>31.10</u>	<u>30.60</u>	<u>30.48</u>	<u>30.15</u>	<u>29.79</u>	<u>29.71</u>	<u>29.35</u>	<u>29.38</u>
	$\times 3.1$	$\times 3.2$	$\times 3.3$	$\times 3.4$	$\times 3.5$	$\times 3.6$	$\times 3.7$	$\times 3.8$	$\times 3.9$	$\times 4.0$
Pixelshuffle	28.78	28.70	28.50	28.30	28.14	28.10	28.72	27.74	27.60	27.65
OSM-bilinear	28.81	<u>28.77</u>	28.62	28.49	28.23	28.22	<u>28.90</u>	<u>27.82</u>	<u>27.79</u>	<u>27.75</u>
OSM-bicubic	28.90	28.81	28.66	28.51	28.26	28.25	28.96	27.84	27.83	27.80
Meta-RDN	28.87	28.79	28.68	28.54	28.32	28.27	28.04	27.92	27.82	27.75
RDN-OSM	28.96	<u>28.70</u>	28.80	28.64	28.41	<u>28.23</u>	<u>28.00</u>	<u>27.97</u>	27.89	27.83

Table 3: Average PSNR of SoA methods using OSM instead of their typical upsampling module. The best results are **highlighted**.

Dataset	Scale	CARN[1]	CARN-OSM	EDSR[23]	EDSR-OSM	RDN[41]	Meta-RDN[12]	RDN-OSM	RCAN[40]	RCAN-OSM
Set5	$\times 2$	37.76	37.90	38.20	38.28	38.24	-	38.31	38.27	38.36
	$\times 3$	34.29	34.35	34.76	34.80	34.71	-	34.77	34.74	34.81
	$\times 4$	32.13	32.15	32.62	32.66	32.47	-	32.58	32.63	32.70
Set14	$\times 2$	33.52	33.60	34.02	34.08	34.01	34.04	34.11	34.12	34.19
	$\times 3$	30.29	30.36	30.66	30.71	30.57	30.55	30.63	30.65	30.74
	$\times 4$	28.60	28.68	28.94	29.01	28.81	28.84	28.91	28.87	28.93
Urban100	$\times 2$	31.92	32.01	33.10	33.15	32.89	-	32.96	33.34	33.40
	$\times 3$	28.06	28.12	29.02	29.09	28.80	-	28.91	29.09	29.15
	$\times 4$	26.07	26.13	26.86	26.91	26.61	-	26.70	26.82	26.90

Table 4: Average PSNR to show the performance of OverNet across scales. The test dataset is Set5. Best results are **highlighted**.

Overscaling factor	Parameters	Scales		
		$\times 2$	$\times 3$	$\times 4$
$\times 3$	927K	38.11	-	-
$\times 4$	943K	38.12	34.49	-
$\times 5$	1079K	38.14	34.54	32.32
$\times 8$	955K	38.15	34.56	32.36

Table 5: Effect of multi-scale loss. OverNet-S uses single-scale loss, OverNet-M multi-scale loss. Best results are **highlighted**.

Dataset	OverNet-S			OverNet-M		
	$\times 2$	$\times 3$	$\times 4$	$\times 2$	$\times 3$	$\times 4$
Set5	38.11	34.49	32.32	38.23	34.60	32.45
B100	32.24	29.17	27.67	32.34	29.30	27.80
Urban100	32.44	28.37	26.31	32.59	28.45	26.42

Moreover, the cost to pay in terms of additional parameters is low. Note that $\times 4$ and $\times 8$ are composed of multiple consecutive $\times 2$ operations, thus introducing less parameters.

Overscaling to higher scales slightly improves the PSNR at the expense of more computation. For the rest of experiments, we overscale to $N + 1$ since it still provides significant improvement at slightly higher computational cost.

Effect of multi-scale loss. Multi-scale learning can process multiple scales with a single trained model, while most of the state-of-the-art algorithms require to train separate models for each supported scale. This property targets real-world applications, where the output size is usually fixed but the input LR scale can vary. Moreover, the multi-scale loss acts as a regularizer, enforcing the generalization of the network across scales and improving performance. As a result, the model is able to maintain accurate reconstruction results across scales. Table 5 shows experimental results, where the model trained with multi-scale loss achieves better performance with a large margin.

4.2. Comparison with State-of-the-art Methods

4.2.1 Results with BI degradation models

We compare the proposed OverNet with nine lightweight state-of-the-art SISR methods [1, 11, 15, 16, 19, 22, 25, 31,

Table 6: Average PSNR/SSIM values for models with the same order of magnitude of parameters. Performance is shown for scale factors $\times 2$, $\times 3$ and $\times 4$ with **BI** degradation. The number of parameters and multi-adds of each method are indicated under their name. The best performance is shown **highlighted** and the second best underlined.

Dataset	Scale	VDSR[15] 0.7M / 0.6T	DRCN[16] 1.7M / 18T	LapSRN[19] 0.8M / 30G	DRNN[31] 0.3M / 6.8T	MemNet[32] 0.7M / 2.6T	SRFBN_S[22] 0.3M / 50G	OISR_LF_s[11] 1.4M / 0.2T	CARN[1] 1.6M / 0.2T	MAFFSRN-L[25] 0.8M / 0.1T	OverNet w/o OSM 0.9M / 0.2T	OverNet 0.9M / 0.2T
Set5	$\times 2$	37.53/0.9587	37.63/0.9588	37.52/0.9591	37.74/0.9591	33.78/0.9597	37.78/0.9156	38.02/0.9605	37.76/0.9590	38.07/0.9607	<u>38.08/0.9609</u>	38.11/0.9610
	$\times 3$	33.66/0.9213	33.82/0.9226	33.82/0.9227	34.03/0.9244	34.09/0.9245	34.20/0.9255	34.39/0.9272	34.29/0.9255	<u>34.45/0.9267</u>	34.43/0.9265	34.49/0.9267
	$\times 4$	31.35/0.8838	31.53/0.8838	31.54/0.886	31.68/0.8888	31.74/0.8893	31.98/0.8923	32.14/0.8947	32.13/0.8932	32.20/0.8953	<u>32.23/0.8954</u>	32.32/0.8956
Set14	$\times 2$	33.05/0.9127	33.06/0.9121	32.99/0.9124	33.23/0.9136	33.28/0.9142	33.35/0.9156	33.62/0.9178	33.52/0.9166	33.59/0.9177	<u>33.64/0.9176</u>	33.71/0.9179
	$\times 3$	29.78/0.8318	29.77/0.8314	29.79/0.8320	29.96/0.8349	30.00/0.8350	30.10/0.8372	30.35/0.8426	30.29/0.8407	30.40/0.8432	<u>30.41/0.8433</u>	30.47/0.8436
	$\times 4$	28.02/0.7678	28.03/0.7673	28.09/0.7994	28.21/0.7720	28.26/0.7723	28.45/0.7779	28.63/0.7819	28.60/0.7806	28.62/0.7822	<u>28.64/0.7823</u>	28.71/0.7826
B100	$\times 2$	31.90/0.8960	31.85/0.8942	31.80/0.8949	32.05/0.8973	32.00/0.8970	32.20/0.9000	32.08/0.8984	32.09/0.8978	<u>32.23/0.9005</u>	32.18/0.8988	32.24/0.9007
	$\times 3$	28.83/0.7976	28.80/0.7963	28.82/0.7973	28.95/0.8004	28.96/0.8010	29.11/0.8085	28.97/0.8025	29.06/0.8034	<u>29.13/0.8061</u>	29.09/0.8033	29.17/0.8063
	$\times 4$	27.29/0.7252	27.24/0.7233	27.32/0.7264	27.38/0.7284	27.44/0.7313	27.60/0.7369	27.44/0.7325	27.58/0.7349	27.59/0.7370	<u>27.60/0.7371</u>	27.67/0.7373
Urban100	$\times 2$	30.77/0.9141	30.76/0.9133	30.41/0.9101	31.23/0.9188	31.31/0.9195	31.41/0.9207	32.21/0.9290	31.92/0.9256	<u>32.38/0.9308</u>	32.35/0.9305	32.44/0.9311
	$\times 3$	27.14/0.8279	27.15/0.8277	27.07/0.8271	27.53/0.8377	27.56/0.8376	26.41/0.8064	28.24/0.8544	28.06/0.8493	28.26/0.8552	<u>28.27/0.8553</u>	28.37/0.8572
	$\times 4$	25.18/0.7525	25.14/0.7511	25.21/0.7553	25.44/0.7638	25.50/0.7630	24.60/0.7258	26.17/0.7888	26.07/0.7837	26.16/0.7887	<u>26.22/0.7920</u>	26.31/0.7923

Table 7: Average PSNR/SSIM for models with the same order of magnitude of parameters (RDN included as a high-capacity reference model). Scores shown for scale factor $\times 3$ using **BD** and **DN** degradation models. Best performance is **highlighted**, second best underlined.

DB	Degrad.	Bicubic	SPMSR[27]	SRCNN[5]	FSRCNN[6]	VDSR[15]	IRCNN_G[38]	IRCNN_C[38]	SRMD(NF)[33]	OverNet w/o OSM	OverNet	RDN[41]
Set5	BD	28.34/0.8161	32.21/0.9001	31.75/0.8988	26.25/0.8130	33.78/0.9198	33.38/0.9182	29.55/0.8246	34.09/0.9242	34.50/0.9270	34.59/0.9287	34.58/0.9280
	DN	24.14/0.5445	-/-	28.10/0.7783	24.24/0.6992	27.81/0.7901	24.85/0.7205	26.18/0.7430	27.74/0.8026	28.37/0.8166	28.49/0.8200	28.47/0.8151
Set14	BD	26.12/0.7106	28.97/0.8205	28.72/0.8024	25.63/0.7312	29.90/0.8369	29.73/0.8292	27.33/0.7135	30.11/0.8304	30.35/0.8307	<u>30.46/0.8310</u>	30.53/0.8447
	DN	23.14/0.4828	-/-	25.55/0.6610	23.10/0.5869	25.92/0.6786	23.84/0.6091	24.68/0.6300	26.13/0.6974	26.56/0.7088	26.62/0.7116	26.60/0.7101
B100	BD	26.02/0.6733	28.13/0.7740	27.97/0.7921	24.88/0.6850	28.70/0.8003	28.65/0.7922	26.46/0.6572	28.98/0.8009	29.06/0.8043	<u>29.13/0.8060</u>	29.23/0.8079
	DN	22.94/0.4461	-/-	25.31/0.6351	23.70/0.5856	25.60/0.6455	23.89/0.5688	24.52/0.5850	25.64/0.6495	25.89/0.6566	25.95/0.6602	25.93/0.6573
Urban100	BD	23.20/0.6661	25.84/0.7856	25.50/0.7812	22.14/0.6815	26.80/0.8191	26.81/0.8189	24.89/0.7172	27.50/0.8370	28.16/0.8471	<u>28.24/0.8485</u>	28.46/0.8582
	DN	21.63/0.4701	-/-	23.40/0.6590	21.15/0.5682	24.01/0.6802	21.96/0.6018	22.63/0.6205	24.28/0.7092	24.84/0.7321	24.93/0.7365	24.92/0.7364

32]. We also train OverNet by replacing its OSM with the typical pixelshuffle upsampling (OverNet w/o OSM). For fair comparison, we train our models individually for each scale factor, including $\times 2$, $\times 3$ and $\times 4$. We test our models on different benchmarks with PSNR and SSIM.

Table 6 shows quantitative evaluation results, including the number of parameters and the number of multiplications and additions (multi-adds), for a more informative comparison (under the method name). Multi-adds were calculated with 1280×720 SR images at all scales. Note that, in this table we only compare models that have a roughly similar number of parameters as ours¹. OverNet exceeds all the previous methods on numerous benchmark dataset. OverNet w/o OSM also achieves comparable or better results. Results show that both OSM and the proposed feature extractor independently increase PSNR when compared to other SR methods. Finally, combining the proposed feature extractor and OSM together further increases performance.

In addition, we present qualitative results in Figure 2. Our proposal produces high-quality image structures. For image *Img_073*, we observe that, unlike OverNet, most of the compared methods fail to recover the definition and orientation of the lines of the blue building. For image *Img_076*, the texture of the predicted SR images for all compared methods contains blur or aliasing. In contrast, our proposal partially recovers the brick pattern, resulting in a more faithful SR image.

¹Additional analyses and qualitative results can be found as supplementary material.

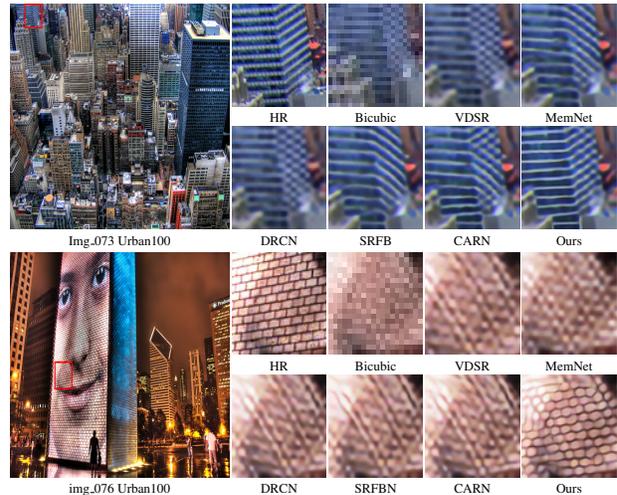


Figure 2: Visual results of **BI** degradation model for $\times 4$.

4.2.2 Results with **BD** and **DN** degradation models

Following [41], we show the results obtained after applying **BD** and **DN** degradation models, and compare to seven SR methods [5, 6, 15, 27, 33, 38], see Table 7. We included the RDN [41] high-capacity model for reference. Because of the mismatch of degradation setups, SRCNN [5], FSRCNN [6], and VDSR [15] have been re-trained for both **BD** and **DN**. Our models achieve the best PSNR and SSIM scores over other SR methods with similar capacity. It can be observed that RDN performs slightly better in some **BD**

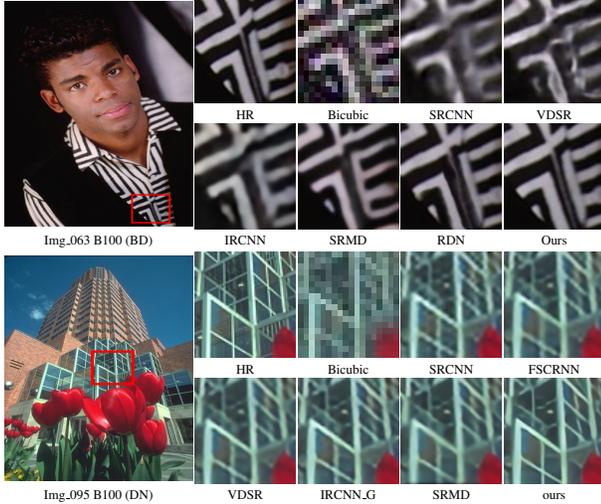


Figure 3: Visual results of **BD** and **DN** degradation models for $\times 3$.

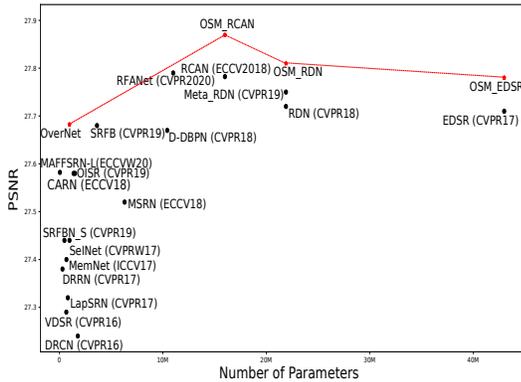


Figure 4: Comparative capacity and performance of state-of-the-art SISR models. The **red stars** represents our methods.

datasets but not in DN datasets. Thanks to OSM, OverNet is able to reduce the DN degradation to obtain better results when compared to RDN. It is worth noting that while RDN has 22M parameters, OverNet only has 0.9M parameters.

In Figure 3 we show two sets of visual results with the **BD** and **DN** degradation models from the standard benchmark datasets. For **BD** degradation, other methods were unable to remove blurring artifacts. In contrast, OverNet could alleviate distortions and generate more accurate details in the SR images. Regarding **DN** degradation, we observe that it is difficult to recover the details with the other methods. However, our method can deliver good results by removing more noise and enhancing details.

4.2.3 Memory complexity and running time analysis

In Figure 4, we compare OverNet against various benchmark algorithms in terms of network parameters and re-

Table 8: Average running time comparison on Urban100 for $\times 4$.

Model	Parameters	Running Time (s)	PSNR
MemNet	0.6M	0.481	25.54
EDSR	43M	1.218	26.64
SRFBN_S	0.4M	0.006	25.71
D-DBPN	10M	0.015	26.38
RDN	22M	1.268	26.61
Meta-RDN	22M	1.350	26.65
Ours	0.9M	0.004	26.31

construction PSNR, using the B100 dataset with a scale of $\times 4$. OverNet achieves the best SR results among all the lightweight SR networks with fewer parameters. In comparison with the networks with a large number of parameters, our proposed OverNet achieves better or competitive results. This demonstrates our method can well balance the number of parameters and the reconstruction performance. We also replace the original upsample modules from different SR methods with OSM: RDN, EDSR and RCAN (RDN+OSM, EDSR+OSM and RCAN+OSM). It can be observed that all the methods with OSM have higher PSNR than the corresponding baselines.

We compare the running time of OverNet on Urban100 with five other state-of-the-art networks, namely MemNet [32], EDSR [23], SRFBN [22], D-DBPN [8], and Meta-RDN [12], using a scale factor $\times 4$. The running time of each network is evaluated using its official code, on the same machine with a NVIDIA 1080 Ti GPU. OverNet is the fastest (see Table 8), reflecting its efficiency.

5. Conclusion

We introduced OverNet, a novel efficient architecture for image super-resolution at arbitrary scales using a single model. OverNet outperforms state-of-the-art algorithms with a reduced number of parameters and low computational requirements. The main contributions are: (i) a lightweight feature extractor that enhances the flow of information to preserve details; (ii) an Overscaling Module that helps to generate accurate SR images at different scaling factors, and (iii) a multi-scale loss that improves training compared to dedicated single-scale models. Thanks to the OSM, we can train a single model for super-resolution at arbitrary scale factors. We proved that the overscaling head can be flexibly applied to other SR models by simply replacing their upsampling module, thus improving their original performance. The provided evidence suggests that the proposed overscaling method may help with other low-level image restoration tasks, such as denoising and dehazing.

Acknowledgements. Authors acknowledge the funding received by the Spanish project TIN2015-65464-R (MINECO/FEDER).

References

- [1] N. Ahn, B. Kang, and K.-A. Sohn. Fast, accurate, and lightweight super-resolution with cascading residual network. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 252–268, 2018.
- [2] A. Aitken, C. Ledig, L. Theis, J. Caballero, Z. Wang, and W. Shi. Checkerboard artifact free sub-pixel convolution: A note on sub-pixel convolution, resize convolution and convolution resize. *arXiv preprint arXiv:1707.02937*, 2017.
- [3] P. Arbelaez, M. Maire, C. Fowlkes, and J. Malik. Contour detection and hierarchical image segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 33(5): 898–916, 2010.
- [4] M. Bevilacqua, A. Roumy, C. Guillemot, and M. L. Alberi-Morel. Low-complexity single-image super-resolution based on nonnegative neighbor embedding. 2012.
- [5] C. Dong, C. C. Loy, K. He, and X. Tang. Learning a deep convolutional network for image super-resolution. In *European conference on computer vision*, pages 184–199. Springer, 2014.
- [6] C. Dong, C. C. Loy, and X. Tang. Accelerating the super-resolution convolutional neural network. In *European conference on computer vision*, pages 391–407. Springer, 2016.
- [7] K. Grm, W. J. Scheirer, and V. Štruc. Face hallucination using cascaded super-resolution and identity priors. *IEEE Transactions on Image Processing*, 29(1):2150–2165, 2019.
- [8] M. Haris, G. Shakhnarovich, and N. Ukita. Deep back-projection networks for super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1664–1673, 2018.
- [9] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [10] K. He, X. Zhang, S. Ren, and J. Sun. Identity mappings in deep residual networks. In *European conference on computer vision*, pages 630–645. Springer, 2016.
- [11] X. He, Z. Mo, P. Wang, Y. Liu, M. Yang, and J. Cheng. Ode-inspired network design for single image super-resolution. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1732–1741, 2019.
- [12] X. Hu, H. Mu, X. Zhang, Z. Wang, T. Tan, and J. Sun. Meta-sr: A magnification-arbitrary network for super-resolution. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1575–1584, 2019.
- [13] J.-B. Huang, A. Singh, and N. Ahuja. Single image super-resolution from transformed self-exemplars. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5197–5206, 2015.
- [14] F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally, and K. Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and 0.5 mb model size. *arXiv preprint arXiv:1602.07360*, 2016.
- [15] J. Kim, J. Kwon Lee, and K. Mu Lee. Accurate image super-resolution using very deep convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1646–1654, 2016.
- [16] J. Kim, J. Kwon Lee, and K. Mu Lee. Deeply-recursive convolutional network for image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1637–1645, 2016.
- [17] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *ICLR*, 2014.
- [18] W. Lai, J. Huang, N. Ahuja, and M. Yang. Fast and accurate image super-resolution with deep laplacian pyramid networks. corr abs/1710.01992. *arXiv preprint arXiv:1710.01992*, 2017.
- [19] W.-S. Lai, J.-B. Huang, N. Ahuja, and M.-H. Yang. Deep laplacian pyramid networks for fast and accurate super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 624–632, 2017.
- [20] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4681–4690, 2017.
- [21] J. Li, F. Fang, K. Mei, and G. Zhang. Multi-scale residual network for image super-resolution. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 517–532, 2018.
- [22] Z. Li, J. Yang, Z. Liu, X. Yang, G. Jeon, and W. Wu. Feed-back network for image super-resolution. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3867–3876, 2019.
- [23] 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*, pages 136–144, 2017.
- [24] J. Liu, W. Zhang, Y. Tang, J. Tang, and G. Wu. Residual feature aggregation network for image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2359–2368, 2020.
- [25] A. Muqet, J. Hwang, S. Yang, J. H. Kang, Y. Kim, and S.-H. Bae. Ultra lightweight image super-resolution with multi-attention layers. *arXiv preprint arXiv:2008.12912*, 2020.
- [26] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer. Automatic differentiation in pytorch. 2017.

- [27] T. Peleg and M. Elad. A statistical prediction model based on sparse representations for single image super-resolution. *IEEE transactions on image processing*, 23(6):2569–2582, 2014.
- [28] W. Shi, J. Caballero, F. Huszár, J. Totz, A. P. Aitken, R. Bishop, D. Rueckert, and Z. Wang. Real-time single image and video super-resolution using an efficient subpixel convolutional neural network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1874–1883, 2016.
- [29] R. K. Srivastava, K. Greff, and J. Schmidhuber. Highway networks. *arXiv preprint arXiv:1505.00387*, 2015.
- [30] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- [31] Y. Tai, J. Yang, and X. Liu. Image super-resolution via deep recursive residual network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3147–3155, 2017.
- [32] Y. Tai, J. Yang, X. Liu, and C. Xu. Memnet: A persistent memory network for image restoration. In *Proceedings of the IEEE international conference on computer vision*, pages 4539–4547, 2017.
- [33] 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*, pages 4799–4807, 2017.
- [34] X. Wang, K. Yu, S. Wu, J. Gu, Y. Liu, C. Dong, Y. Qiao, and C. Change Loy. Esrgan: Enhanced super-resolution generative adversarial networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 0–0, 2018.
- [35] Z. Wang, J. Chen, and S. C. Hoi. Deep learning for image super-resolution: A survey. *arXiv preprint arXiv:1902.06068*, 2019.
- [36] J. Yu, Y. Fan, J. Yang, N. Xu, Z. Wang, X. Wang, and T. Huang. Wide activation for efficient and accurate image super-resolution. *arXiv preprint arXiv:1808.08718*, 2018.
- [37] R. Zeyde, M. Elad, and M. Protter. On single image scale-up using sparse-representations. In *International conference on curves and surfaces*, pages 711–730. Springer, 2010.
- [38] K. Zhang, W. Zuo, S. Gu, and L. Zhang. Learning deep cnn denoiser prior for image restoration. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3929–3938, 2017.
- [39] K. Zhang, W. Zuo, and L. Zhang. Learning a single convolutional super-resolution network for multiple degradations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3262–3271, 2018.
- [40] Y. Zhang, K. Li, K. Li, L. Wang, B. Zhong, and Y. Fu. Image super-resolution using very deep residual channel attention networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 286–301, 2018.
- [41] 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*, pages 2472–2481, 2018.