

Single Image Super-resolution via a Lightweight Residual Convolutional Neural Network

Yudong Liang, Ze Yang, Kai Zhang, Yihui He, Jinjun Wang, *Senior Member, IEEE*,
and Nanning Zheng, *Fellow, IEEE*

Abstract—Recent years have witnessed great success of convolutional neural network (CNN) for various problems both in low and high level visions. Especially noteworthy is the residual network which was originally proposed to handle high-level vision problems and enjoys several merits. This paper aims to extend the merits of residual network, such as skip connection induced fast training, for a typical low-level vision problem, i.e., single image super-resolution. In general, the two main challenges of existing deep CNN for super-resolution lie in the gradient exploding/vanishing problem and large numbers of parameters or computational cost as CNN goes deeper. Correspondingly, the skip connections or identity mapping shortcuts are utilized to avoid gradient exploding/vanishing problem. In addition, the skip connections have naturally centered the activation which led to better performance. To tackle with the second problem, a lightweight CNN architecture which has carefully designed width, depth and skip connections was proposed. In particular, a strategy of gradually varying the shape of network has been proposed for residual network. Different residual architectures for image super-resolution have also been compared. Experimental results have demonstrated that the proposed CNN model can not only achieve state-of-the-art PSNR and SSIM results for single image super-resolution but also produce visually pleasant results. This paper has extended the mmm 2017 oral conference paper with a considerable new analyses and more experiments especially from the perspective of centering activations and ensemble behaviors of residual network.

Index Terms—super-resolution, deep residual convolutional neural network, skip connections, parameter numbers

I. INTRODUCTION

SINGLE image super-resolution (SISR) [1], [2], [3] aims to recover a high-resolution (HR) image from the corresponding low-resolution (LR) image. It is a very practical technique due to its high value in various fields, such as producing high-definition images from low-cost image sensors, medical imaging and satellite imaging. Restoring the HR image from the single LR input is also a very difficult problem of high theoretical values which arouses more and more interests from the academic communities [4], [5], [6], [7] and large companies [8], [9], [10]. Typically, it is very challenging to restore the missing pixels from an LR observation since the number of pixels to be estimated in the HR image is usually

much larger than that in the given LR input. The ill-posed nature of single image super-resolution problem makes restoring HR images an arena to evaluate inference and regression techniques. Generally, SISR techniques can be roughly divided into three categories: the interpolation methods, the reconstruction methods [11] and the example based methods [12], [13].

Most of the recent SISR methods fall into the example based methods which learn prior knowledge from LR and HR pairs, thus alleviating the ill-posedness of SISR. Representative methods mainly rely on different learning techniques to incorporate image priors for super-resolution process, including neighbor embedding regression [14], [15], [16], sparse coding [13], [17], tree based regressions [18], [19] and deep convolutional neural network (CNN) [20], [5], [21], [22].

Among the above techniques, deep learning techniques especially deep CNN have largely promoted the state-of-the-art performances in SISR area. Dong *et al.*[20] firstly proposed a deep convolutional neural network termed SRCNN with three convolutional layers for image super-resolution which gave the best practise at that time. Later, Dong *et al.*[5] extended SRCNN with larger filter sizes and filter numbers while kept the depth of CNN fixed to further improve the performance. They found that deeper models were hard to train and [5] failed to boost the performance by increasing the depth. Such findings indicate that deeper models are not suitable for image super-resolution, which is counter-intuitive as deeper models have been proved more effective in many tasks [23], [24], [25]. Instead of directly predicting the HR output, Kim *et al.*[22] proposed a very deep CNN (VDSR) of depth up to 20 by a large skip connection to predict the residual image, i.e., the high frequency of the HR image. VDSR surpasses SRCNN with a large margin which mainly benefits from two aspects: deeper architecture and predicting high frequency of images only which is called residual learning by [22].

As demonstrated in [22], the SR results have been improved as the VDSR network goes deeper to a certain depth (20). Although VDSR has achieved impressive results, the plain structure of VDSR which simply stacks layers hampers the convergence of deeper architectures due to the gradient exploding/vanishing problem. It would bring little improvement as the network goes deeper. Fortunately, the residual network [24], [25] has successfully addressed this issue. As a result, different from VDSR, this paper has designed a novel very deep residual convolutional neural network shown in Fig. 1. As LR image and target HR image are highly correlated, predicting high frequency of the image only is a kind of residual learning which largely lowers the price for training. A

Yudong Liang is with the Key Laboratory of Computational Intelligence and Chinese Information Processing of Ministry of Education, and the School of Computer and Information Technology, Shanxi University, 92 Wucheng Road, Taiyuan, Shanxi Province, 030006, China. A large part of the work is done when he was with the Institute of Artificial Intelligence and Robotics, Xi'an Jiaotong University, China. e-mail: liangyudong006@163.com.

Ze Yang, Yihui He, Jinjun Wang, and Nanning Zheng are with the Institute of Artificial Intelligence and Robotics, Xi'an Jiaotong University, China.

Kai Zhang is with Harbin Institute of Technology, China

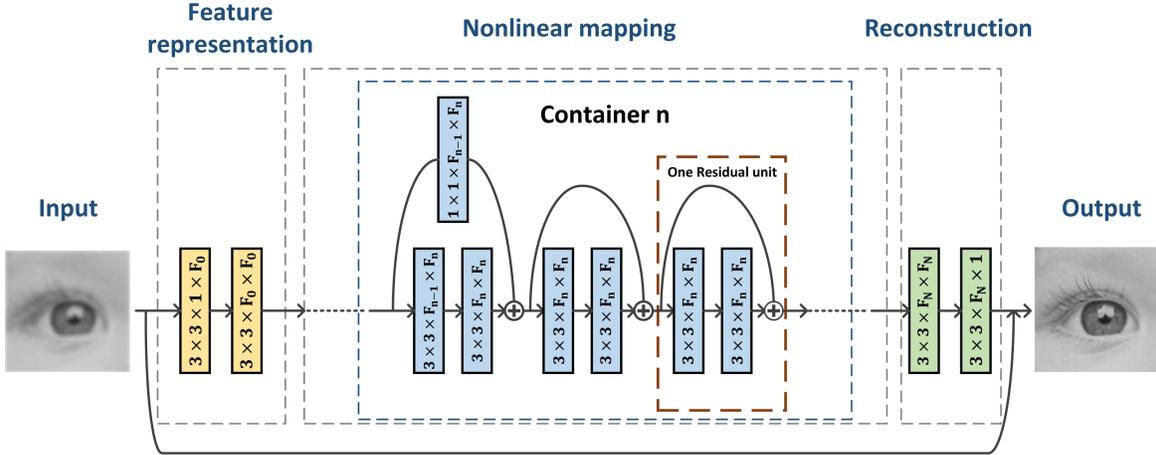


Fig. 1. The architecture of our residual model.

totally deep residual CNN will be expected to fully take advantage of the correlations between LR and HR images. Moreover, skip connections or identity mapping shortcuts in deep residual CNN would alleviate gradient vanishing/exploding problem when the network becomes increasingly deeper.

For neural network it is known that mean shifts toward zero or centering the activations speeds up learning [26], [27] by bringing the normal gradient closer to the unit natural gradient [28]. LeCun *et al.* [26] justified the benefits of centering activation functions and advised to center the distribution of supervision information. The mean value of distribution for high frequency in the images is around zeros, while the mean of raw HR image pixels biases above zero. Thus, it is easy to understand that predicting residual images (high frequency) instead of HR images has largely improve the convergence of the network [22]. Batch Normalization (BN) [29] also aimed to center activations by reducing the internal covariate shift with extra moving average computations. Raiko *et al.* [30] proved that shortcut connections made the Fisher information matrix closer to a diagonal matrix and standard gradient closer to the natural gradient which contributes to centering the outputs and slopes of the hidden layers in multi-layer perceptron network. Thus, identity mapping shortcuts naturally help the residual network center the activations.

The Batch Normalization (BN) layers in the conventional residual branches [24], [25] are abandoned in our proposed architecture as skip connections and predicting high frequency (with a zero mean) have ensured centering the activations if the network is not too deep. Our designed residual architectures which we refer to as “SRResNetNB” have consumed less computational resources and achieved better performances empirically.

While very deep CNN model would increase the model capacity, on the other hand, it would introduce a huge number of parameters which is sometimes unacceptable for applications. Thus, when hardware resources are limited, a lightweight architecture using less parameters is essential for real word applications. In this paper, the ‘shape’ of deep CNN has been investigated to largely reduce the parameter numbers. The ‘shape’ of deep CNN refers to depth, all the filter sizes

and numbers of each layer which decide sizes and numbers of feature maps in each layer to form a global shape. With a residual architecture and lightweight ‘shape’ design, the proposed model can not only achieve state-of-the-art PSNR and SSIM results for single image super-resolution but also produce visually pleasant results.

A preliminary version of this work was presented earlier [31]. The present work adds to the conference oral version in significant ways: first, different deep architectures of residual branches are explored to further conclude a principle of designing a deep network for image super-resolution. Second, a considerable new analysis from the perspective of centering activations and ensemble behaviors of residual networks has been represented and intuitive explanations are supplied to the result. In particular, a strategy of gradually varying the ‘shape’ of the residual network has been clarified in constructing a lightweight structure, based on the assumption that the residual network has been seen as an ensemble of relatively shallow networks with a large capacity [32]. Third, more detailed experiments are represented to design the structures and retrench the parameters of the residual model.

II. RELATED WORKS

In the pioneer work by Freeman *et al.* [12], the co-occurrence priors were proposed that similar LR local structures often relate to similar HR local information. From LR and corresponding HR images, LR and HR examples (patches or sub images) could be extracted to form training databases. The mappings from LR to HR examples call for accurate regression methods to be applied. In fact, the learning based regression methods especially deep learning based methods have dominated the example based methods.

Since the work of SRCNN [20], deep CNNs have refreshed the state-of-the-art performances in super-resolution area. Simply elaborating the filter sizes and filter numbers for SRCNN [5] had further improved the performance. Wang *et al.* [33] designed the CNN architecture to incorporate the sparse coding prior based on the learned iterative shrinkage and thresholding algorithm (LISTA). With sparsity prior mod-

eling, the performance boosted even with a model of smaller size compared with SRCNN.

Kim *et al.* [22] made a breakthrough for image super-resolution by predicting residual images and using a much deeper CNN termed VDSR up to 20 layers, which largely accelerated the speed of training and outperformed SRCNN presented by Dong *et al.* [5]. To ensure the fast convergence of deep CNN and avoid gradient vanishing or exploding, a much larger learning rate for training was cooperated with adjustable gradient clipping in VDSR training. VDSR is inspired by the merits of VGG net which attempts to train a thin deep network. However, this kind of plain networks are not easy to be optimized when they go even deeper as demonstrated by He *et al.* [24], [25].

The difficulties of training deeper plain networks were carefully analyzed by He *et al.* [24], [25]. The degradation problem [24] has been observed that the testing accuracy even the training accuracy becomes saturated then degrades rapidly as plain networks go deeper. This degradation is caused by the difficulties of training other than overfitting. It has been demonstrated that learning a residual function is much more easier than learning the original prediction function with very deep CNN. Residual networks with a surprising depth were designed for image classification problems with skip connections or identity mapping shortcuts. Later, a detailed analysis on the mechanisms of identity mapping in deep residual networks and a new residual unit design have been represented in [25].

Residual network has also been applied in conjunction with perceptual loss [34] to generate style transferred image and produce visual more pleasing HR images in large magnification factors. SRResNet [9], another famous concurrent work with us has also designed a residual network with skip connections for image super-resolution which serves as the generator network in a generative adversarial framework, termed SRGAN. To produce photo-realistic results, SRGAN [9] exploited an adversarial loss to replace the traditional mean squared reconstruction error. This adversarial framework recovered images with better perceptual quality and especially favored large upscaling factors (*e.g.*, 4). The success of these work and our previous version [31] has indicated the importance of skip-connections for image super-resolution. Later, Tai *et al.* [35] introduced skip connections of multiple paths and shared the weights of residual units in one recursive block. Most of the residual networks for image super-resolution designed the residual branches as a combination of convolution, nonlinear activation (such as ReLU or PReLU) and Batch Normalization (BN) layers, which are the same as the residual branches for image classification task.

The idea of shortcuts has been related with centering the activations at zero for multi-layer perceptron network [30]. Raiko *et al.* [30] proposed to transform the outputs of each hidden layers in multi-layer perceptron network to have zero output and zero slope on average and use separate shortcut connections to model the linear dependencies. It is known that centering the activations accelerates learning [36], [30], [26], [27]. LeCun *et al.* [26] analyzed the eigenvalues of Hessian matrix during the gradient descend process and give a

theoretical justification for the benefits of centering activation functions. The applied skip-connections have already centered the activations at zero within certain depth and the mean of distributions for high frequencies in images is close to zero, which indicate the BN layers could be eliminated in our residual units.

Our design has been further supported by the very recent work [37], which wins the first prize in Ntire 2017 challenge on single image super-resolution [38]. Liang *et al.* [39] further extended the identity skip connections to projection skip connections and explored the power of internal priors for deep architectures.

After largely easing the difficulties of training much deeper CNN with residual functions by shortcuts or skip connections, the huge number of parameters in deep architecture is still a big problem for computational resources and storages. The evolvement of Inception models [40], [29], [41], [42] has demonstrated that carefully designed topologies enable compelling performances with less parameters. He *et al.* [24], [25] attempt to alleviate the problem by bottleneck architectures. The bottleneck architectures first utilize 1×1 convolutions to reduce the dimensions, then after some operations, 1×1 convolutions are applied again to increase the dimensions. With such a residual unit design, the number of parameters could be largely reduced. Thus, the ‘shape’ of CNN could be potentially explored to reduce the parameters while maintain the performances. The bottleneck architectures decrease the parameter numbers at the expense of increasing the depth of the network to mountain the performances. In the meanwhile, contextual information is very important for image super-resolution [22], [5], such 1×1 convolutions design may give a negative effort to the SR results. A study on the influences of the ‘shape’ (the filter sizes, depth and numbers of convolutions in each layer) on the performance of image super-resolution has been represented in the following sessions. With a carefully design and exploration of the ‘shape’ of the network, novel residual deep models are proposed for image super-resolution task in this paper.

III. A LIGHTWEIGHT RESIDUAL DEEP MODEL FOR IMAGE SUPER-RESOLUTION

Following the example based methods, HR examples I^h and LR examples I^l are extracted from HR images I^H and LR images I^L respectively. The degeneration process of LR images I^L from the corresponding HR images I^H could be considered as the following blurring process related with blur kernel G and downsampling process \downarrow_s with a scale factor s

$$I^L = (I^H \otimes G) \downarrow_s. \quad (1)$$

In the experiments, this process is simulated by a ‘bicubic’ downscale interpolation.

In the next part, residual deep models for image super-resolution will be designed from the perspective of centering the activations to speed up learning.

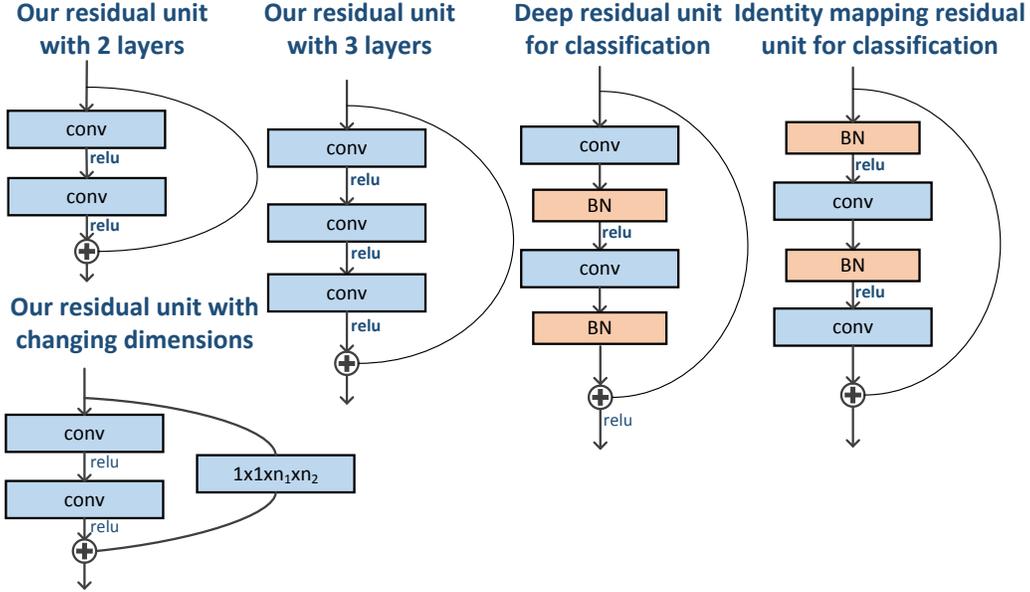


Fig. 2. The architectures of different residual units.

A. Deep Residual Models

The gradient exploding/vanishing problems are largely alleviated by skip connections in the deep residual models [24], [25]. The architectures of our deep residual models especially the residual branches will be further designed from the perspective of centering activations.

Simply stacking the convolutional layers and rectified linear units as VDSR fashions [22] will have a mean activation larger than zero [36]. The non-zero mean activation acts as bias for the next layer. The more the layers are correlated, the higher their bias shift.

For a multi-layer perceptron network, Raiko *et al.* [30] proved that the transformation by shortcuts centered the activations which made the Fisher information matrix closer to a diagonal matrix, and thus standard gradient closer to the natural gradient. The transformations can be as

$$x^{k+1} = \mathbf{A} \cdot \mathbf{T}(\mathbf{B} \cdot x^k) + \mathbf{C} \cdot x^k, \quad (2)$$

where \mathbf{A} , \mathbf{B} , \mathbf{C} is the weight matrices, \mathbf{T} is a nonlinearity activation, x^k is the output of the neurons of the k th layer.

Similarly, for convolutional neural network, transformations can be as

$$x^{k+1} = f(\theta^k, x^k) + \mathbf{C} \cdot x^k, \quad (3)$$

where f is a function composed by convolutions, nonlinearity activation, and Batch Normalization (BN). When the weight matrix \mathbf{C} becomes identity matrix, function $f(\theta^k, x^k)$ will become our residual branches. Thus, our residual networks with skip connections can naturally centering activations and speed up learning.

For image super-resolution problems, super-resolution is only applied on the luminance channel (Y channel in YCbCr color space) in most of previous study [15], [20], [22]. It is obvious that the distribution of values on the luminance channel in the output HR images doesn't center at zero, while the residual images (high frequency of the images) have means

towards zero. To center the activations, our deep residual CNN applies a large skip-connection as [22] which makes the network predict the residual images (the high frequency of the images). Predicting the residual images has largely improved the training speed and convergency results.

Our deep residual CNN for image super-resolution is an end-to-end mapping model which can be roughly divided into three sub-networks to perform three steps: feature representation, nonlinear mapping, and reconstruction.

The feature representation sub-network extracts LR discriminative features from the LR input images, while nonlinear mapping part maps the LR feature representations into HR feature representations. Reconstruction part restores the HR images from HR feature representations. Feature representation sub-network applies plain network stacking convolutional and ReLU layers as shown in Fig. 1 and reconstruction sub-network only uses convolutional layers as [39]. The main body of our model, nonlinear mapping part consists of residual units which center the activations with shortcuts and ease the difficulties of training.

Typical units of our deep residual CNN are shown in Fig. 2. Empirically, residual unit with 2 or 3 convolutional layers works well for image super-resolution problem, those two kinds of units are applied in the experiments. When featuremap dimensions change, the identity shortcut becomes a projection to change feature dimensions. The second right and rightmost are one unit of residual net for image classification problems proposed by He *et al.* [24], [25] respectively. Compared with them, the architectures of our residual functions are composed of convolutional, ReLU layers and shortcuts, which are quite different. Batch Normalization units are discarded and deployments are different. Batch Normalization [29] reduces the distribution variations of layers (internal covariate) by normalizing the input of the layers. With an input x , the output

of BN layer is given by

$$BN_{\gamma,\beta} = \gamma \left(\frac{x - \mu}{\sqrt{\sigma^2 + \varepsilon}} \right) + \beta, \quad (4)$$

where γ and β are learnable parameters, μ and σ are the mean and variance of activations in the mini-batch, respectively, ε is a small constant for numerical stability. Obviously, the activation after Batch Normalization operation has also been centered. As skip connections (Eq. (3)) have naturally corrected the bias shift, thus if the residual network is not that deep¹, the BN layers can be abandoned as it needs extra learning and inference computations which take much more computational resources.

Shortcuts or skip connections which are identity mappings are realized by element-wise additions. As this element-wise addition increases very little computations, our feed-forward deep residual CNN has a similar computational complexity with VDSR [22] fashions network. Similar with VDSR [22], small convolutional filter of size 3×3 has been applied. Assuming the input of k -th residual unit as x^k , the residual functions have the following form

$$x^{k+1} = x^k + f(\theta^k, x^k), \quad (5)$$

where θ^k are the parameters of k -th residual unit.

A simple Euclidean loss function is adopted to make predictions approximate the high frequencies of examples

$$\mathcal{L} = \frac{1}{2n} \sum_{i=1}^n \|\mathcal{F}(\theta, I_i^l) - (I_i^h - I_i^l)\|^2 \quad (6)$$

where n is the number of patch pairs (I^l, I^h) , $\mathcal{F}(\theta, I^l)$ denotes the predictions of our deep residual CNN with parameter θ . Our deep residual CNN is composed of several **Containers** which have certain number of residual units. For succinctness, the filter numbers keep the same in each single container. The architectures of our deep residual CNN will be described as a sequence of the filter numbers $(N1_{k1}, N2_{k2}, \dots)$ in containers. If subscript k exists for N_k , it means there are k residual units with each having a filter number of N in this container.

Stochastic gradient descent (SGD) with the standard back-propagation [43] is applied to train our deep residual CNN. In particular, the parameter is updated as Eq. (7), where m denotes the momentum parameter with a value of 0.9 and η is the learning rate.

$$\Delta_{i+1} = m \cdot \Delta_i + \eta \cdot \frac{\partial \text{loss}}{\partial \theta_i}, \quad \theta_{i+1} = \theta_i + \Delta_{i+1} \quad (7)$$

High learning rates are expected to boost training with faster and better convergency. Adjustable gradient clipping [22] is utilized to keep learning rates high while at the same time to prevent the net from gradient exploding problems. Gradients $\frac{\partial \text{Loss}}{\partial \theta_i}$ are clipped into the range of $[-\frac{\tau}{\eta}, \frac{\tau}{\eta}]$, where τ is a constant value.

¹The bias from zero will accumulate as the network goes deeper.

B. Lightweight Design for the Proposed Model

In this section, the ‘shape’ of deep CNN has been explored to achieve better performances but with less number of parameters. The ‘shape’ of deep CNN is determined by all the sizes and numbers of filters in each layer besides the depth of the network. Thin but small filter size works well with padding which leads to larger receptive field as network goes deep, in specific, 3×3 filter size has been applied. It is general that deeper and wider network will have larger model capacity and better feature representational ability. However, the number of parameters is restricted by the hardware or computational resources. Using less parameters to achieve better performances is essential for applications. Next, filter numbers and the combinations of filter numbers will be discussed to retrench parameters for a better performance.

1) *Exploring the Shape of the Architecture*: Inspired by the evolvement of Inception models [40], [29], [41], [42] and the bottle-neck architecture [25], it is supposed that changing the shape of the architecture may maintain the performance while largely reduce the computational parameters. Instead of applying 1×1 convolutions as bottle-neck architecture, the 3×3 convolutions are applied as image SR process largely depends on the contextual information in local neighbor areas.

The filter numbers of VDSR are kept the same. There seems to be few principles to decide filter numbers and the combinations of filter numbers in a network. Instead of using a same number of filters in a network, the filter numbers can be varied to potentially reduce parameters which could enable a deeper or wider network.

Residual networks can be interpreted as an ensemble of many paths of differing depth [32] and residual networks enable very deep networks by leveraging only the short paths during training [32]. According to this assumption, if the models of short paths in the residual network have been less disturbed, the performance of residual network which is an ensemble could keep stable.

A strategy of gradually varying the ‘shapes’ of residual models is proposed by us to reduce parameters. Gradually varying the shape of network means the filter numbers of the adjacent layers should increase or decrease gradually. This has been illustrated as Fig. 3. In Fig. 3, different residual branches and corresponding skip connections are denoted by different colors. The residual networks can be unfolded as a summations of models from different paths of residual networks. Considering a residual network with three units or sub residual network, if the filter numbers of the adjacent layers change gradually, e.g., only the filter numbers of R3 changes (e.g., decreases), a lot of paths are unaffected. Thus, the residual networks are more robust to the shape varying and our strategy can be applied to achieve better performances with less parameters.

The impacts of the feature map numbers in each layer on performance are carefully explored in the following fashions as Fig. 4: gradually, the numbers of feature maps

- increase monotonically up to N.
- decrease monotonically from N.
- increase up to N then decrease.

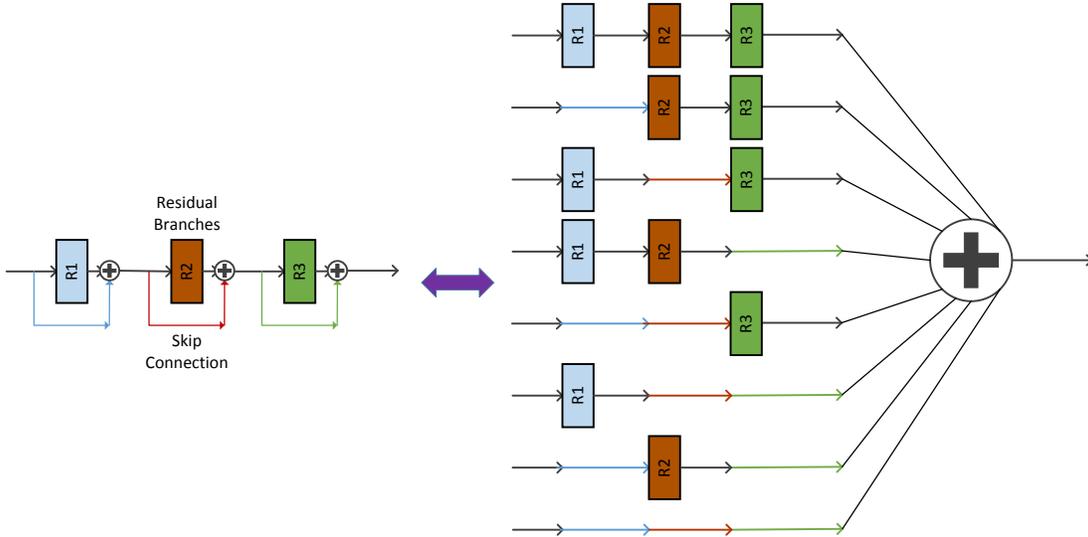


Fig. 3. Residual network behaves like an ensemble of networks

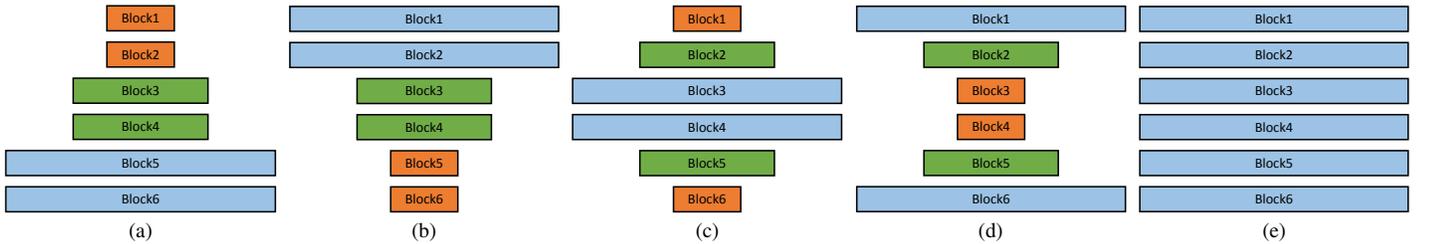


Fig. 4. Different ‘shapes’ of networks which gradually vary the feature map numbers. The width of the block correlates to the number of feature maps in the layer.

- decrease from N then increase.
- keep the same as N (baseline).

In Fig. 4, the width of the square block correlates to the numbers of the feature map in the layer. The larger width of square block indicates there are more feature maps in that layer. Compared with the baseline way that the feature map numbers keep the same, applying gradually varying the shape strategy has largely reduced the parameters.

The experiments demonstrate that different lightweight designs have achieved comparable performances with less parameters. This will be further discussed in the experiments part. In comparison with our residual CNN, the performances of VDSR with different shapes fluctuate heavily. This proves our residual architectures are more robust to the shape varying of CNN and our strategy of gradually varying the ‘shape’ of residual network could be applied to achieve better performances with less parameters.

2) *Training with Multiple Upscaling Factors to Retrench Parameters:* It has been pointed out that it is feasible to train a deep CNN for different upscaling factors [22]. Training datasets for different specified upscaling factors are combined together to enable our deep residual CNN to handle multiple upscaling factors, as images across different scales share some common structures and textures. Parameters are shared across different predefined upscaling factors which further dispenses with the trouble of retaining different models for different upscaling factors. It will retrench parameters when multiple

upsampling factors are required.

IV. EXPERIMENTS

In this section, we conducted a series of experiments to explore the empirical principles to design a deep architecture for image super-resolution problem. The performances of the proposed method against the state-of-the-art SISR methods are compared which clearly demonstrate better or comparable subjective scores and more visual pleasing results.

The same 291 training images applied by VDSR were utilized for training, including 91 images proposed in Yang *et al.*[13] and 200 natural images from Berkeley Segmentation Dataset (BSD). For testing, four datasets were investigated: ‘Set5’ and ‘Set14’ [15], [20], ‘Urban100’ [7] and ‘BSD100’ [15], [6].

The size of example was set as 41×41 and the batch size was chosen as 64. Momentum and weight decay parameters were fixed as 0.9 and 0.0001 respectively. Multi-scale training was applied in all of the following experiments. Weight initialization methods [24], [25] were applied with small modulations. Learning rate was initially set to 0.1 and then decreased by a factor of 10 every 30 epochs. All these settings ensure us to make a fair comparison with the competing approaches including VDSR method.

TABLE I
COMPARISON IN DIFFERENT DATASETS AND WITH DIFFERENT SCALES.

Dataset	Scale	Bicubic PSNR/SSIM	A+[16] PSNR/SSIM	RFL[18] PSNR/SSIM	SelfEx[7] PSNR/SSIM	SRCNN[20] PSNR/SSIM	VDSR[22] PSNR/SSIM	SRResNetNB PSNR/SSIM	R-basic PSNR/SSIM
Set5	×2	33.66/0.9299	36.54/0.9544	36.54/0.9537	36.49/0.9537	36.66/0.9542	37.53/0.9587	37.51/0.9587	37.27/0.9577
	×3	30.39/0.8682	32.58/0.9088	32.43/0.9057	32.58/0.9093	32.75/0.9090	33.66/0.9213	33.72/0.9215	33.43/0.9190
	×4	28.42/0.8104	30.28/0.8603	30.14/0.8548	30.31/0.8619	30.48/0.8628	31.35/0.8838	31.37/0.8838	31.15/0.8796
Set14	×2	30.24/0.8688	32.28/0.9056	32.26/0.9040	32.22/0.9034	32.42/0.9063	33.03/0.9124	33.10/0.9131	32.86/0.9113
	×3	27.55/0.7742	29.13/0.8188	29.05/0.8164	29.16/0.8196	29.28/0.8209	29.77/0.8314	29.80/0.8317	29.67/0.8297
	×4	26.00/0.7027	27.32/0.7491	27.24/0.7451	27.40/0.7518	27.49/0.7503	28.01/0.7674	28.06/0.7681	27.90/0.7648
BSD100	×2	29.56/0.8431	31.21/0.8863	31.16/0.8840	31.18/0.8855	31.36/0.8879	31.90/0.8960	31.91/0.8961	31.76/0.8940
	×3	27.21/0.7385	28.29/0.7835	28.22/0.7806	28.29/0.7840	28.41/0.7863	28.82/0.7976	28.83/0.7980	28.73/0.7954
	×4	25.96/0.6675	26.82/0.7087	26.75/0.7054	26.84/0.7106	26.90/0.7101	27.29/0.7251	27.29/0.7248	27.19/0.7221
Urban100	×2	26.88/0.8403	29.20/0.8938	29.11/0.8904	29.54/0.8967	29.50/0.8946	30.76/0.9140	30.88/0.9150	30.47/0.9100
	×3	24.46/0.7349	26.03/0.7973	25.86/0.7900	26.44/0.8088	26.24/0.7989	27.14/0.8279	27.17/0.8283	26.92/0.8208
	×4	23.14/0.6577	24.32/0.7183	24.19/0.7096	24.79/0.7374	24.52/0.7221	25.18/0.7524	25.22/0.7537	25.02/0.7452

A. Comparisons with the State-of-the-art Methods

Table I shows the quantitative comparisons with A+ [16], RFL [18], SelfEx [7], SRCNN [20] and VDSR [22]. Visual results are also represented to give intuitive assessment. In Table I, two models of our deep residual CNN with different depth have been investigated, denoted as **R-basic** and **SRResNetNB** respectively. The residual unit in R-basic and deeper and larger model SRResNetNB has two convolutional layers. R-basic ($16_3, 32_3, 64_3$) has 22 layers, while SRResNetNB ($16_3, 32_3, 64_3, 128_3, 256_3$) has 34 layers. SRResNetNB has achieved the best performances compared with other methods in most cases and comparable results in other situations.

R-basic outperforms the other methods except VDSR. However, the performances of VDSR (20 layers) have not been obtained by our reimplementation. For example, the average PSNR of VDSR by our reimplementation for Set5 and Set14 are 37.32dB and 32.89dB respectively, with a gap of more than 0.1db from the reported results. Assisted with the missing tricks, the performance of our model is expected to be further boosted. In Fig. 5, the PSNR against training epochs has been compared among R-basic, SRResNetNB, and VDSR trained by us. Deeper and larger model SRResNetNB outperformed VDSR at very beginning with a large margin. Although R-basic contains much less parameters, R-basic model has obtained comparable performances with VDSR.

In Fig. 8, all the compared results are obtained by the released code of the authors or from the reported ones in the paper. Visually pleasing results have been achieved by our model. Restorations of our method contain more authentic texture and more clear details compared with the results by other methods such as the texture of the zebra head. Our method has provided less artifacts, *e.g.*, all the other methods except ours have restored obvious artifacts at the location of book. Shaper edges have appeared in our restorations which have represented visually more pleasing results.

B. Number of Parameters

For R-basic model, there are 22 convolutional layers and 0.3M(322721) parameters accumulated by the numbers of corresponding weights and bias. For SRResNetNB model, 34 convolutional layers and 5M(4975905) parameters are applied. The compared VDSR in Table I is 20 layers and has

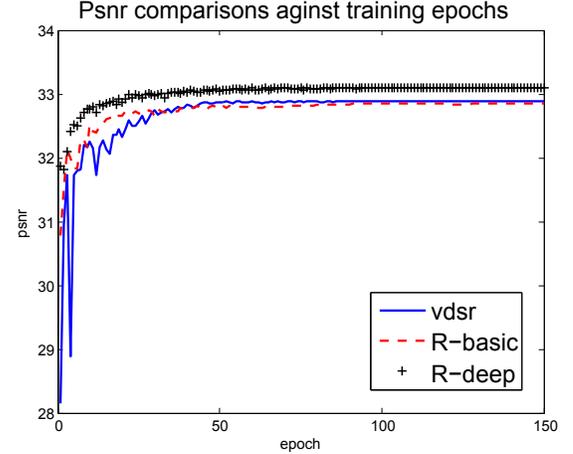


Fig. 5. Comparisons of test psnr on Set 14 against training epochs among SRResNetNB(denoted as R-deep), R-basic and VDSR.

0.7M(664704) parameters. Although SRResNetNB has more parameters, our SRResNetNB model is still acceptable which can be efficiently trained with a single GPU.

C. The Position of ReLU

In the residual branches, convolutional and ReLU layers are applied. The performances compared with the positions of ReLU layers (ReLU before/after conv) as in Fig. 6 are represented in Table II on Set14. The compared network has a same depth and corresponding convolutional layers among these networks have the same parameter numbers.

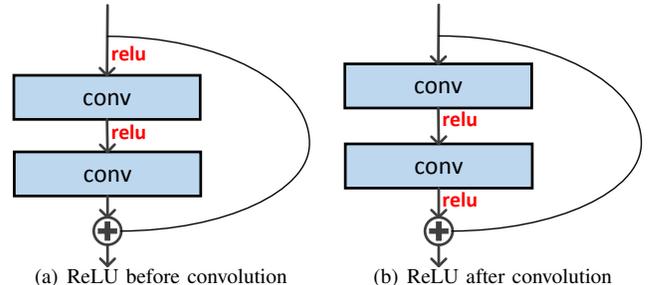


Fig. 6. The positions of ReLU in residual branches.

TABLE II
ABLATION COMPARISONS FOR RESIDUAL NETWORK WITH DIFFERENT
ORDERS OF CONVOLUTION AND RELU LAYERS IN TERMS OF AVERAGE
PSNR (DB) ON SET14.

scale	identity+ ReLU after conv	identity+ ReLU before conv
× 2	32.97	33.01
× 3	29.75	29.77
× 4	28.02	28.02

From the results in Table II, we conclude that the positions of ReLU in the residual branches make small differences.

D. Impacts of Batch Normalization on SISR

In Fig. 7, test PSNR of Set 14 against training epochs by our R-basic with and without BN are compared to demonstrate the impacts of Batch Normalization on SISR problems. In Fig. 7, the compared structure with BN layers is the same as the structure applied for image classifications [25], showed in the rightmost column in Fig. 2.

It seems adding BN operations has hampered further improvement when more epochs have been performed. Normalizing input distribution of mini-batch to suppress data shifting has been proved powerful and largely accelerated the training convergency speed. It also enables deeper architectures and larger learning rates to be utilized in other tasks. However, whiten input and output of the intermediate layer may not be suitable for image super-resolution task which needs precise output. Another suspect may be regularization effects of BN have not been fully exploited as the training set of Fig. 7 is still limited in contrast with ImageNet. As larger learning rates were enabled by gradient clipping methods, the benefits of BN for leaning rates are alleviated.

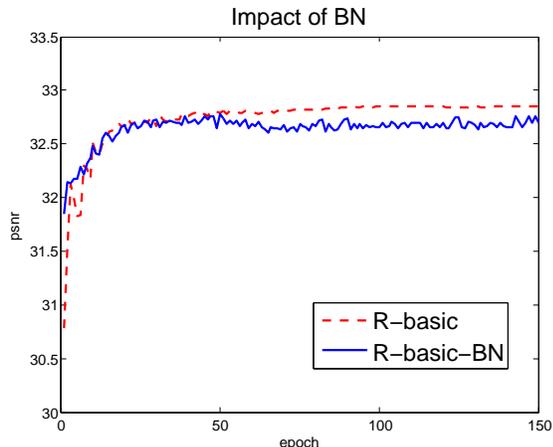


Fig. 7. The impacts of BN: test psnr of Set 14 against training epochs by our R-basic with and without BN.

From the perspective of centering activations, the skip connection itself has the benefits of centering the activations which partially reduces the necessities of BN operations when the network is not too deep to correct the mean bias. Moreover, the BN operation takes extra computations during learning and inference. Without BN operation, provided with certain

computational resources, larger and wider deep architectures can be enabled to get better performances.

The impacts of Batch Normalization on SISR are still an open issue for the future study.

E. The Deeper the Better, the Wider the Better

SRResNetNB performs much better than R-basic and VDSR model with deeper and wider network. Next, ablations of our system would be evaluated to unpack this performance gain. The skip connections and two factors, width (related to filter numbers) and depth of our model would be analyzed in the following steps.

TABLE III
PSNR COMPARISON BETWEEN OUR RESIDUAL CNN AND VDSR TRAINED
BY US

	Set5	Set14	BSD100	Urban100
$R(64_8)$	37.28	32.91	31.72	30.45
VDSR	37.32	32.89	31.77	30.51

First, 20-layer VDSR has been added with 8 identity shortcuts to form a residual network, denoted $R(64_8)$. Each residual unit has two convolutional and Relu layers. The performance of $R(64_8)$ is roughly the same as VDSR in Table III. The shortcuts have very little impacts on the descriptive power. From the perspective of the centering the activations, to predict the high frequencies of image has pushed the final output activation centered. Within certain depth (*e.g.*, 20), the difficulties of learning has been alleviated, thus the shortcuts of the residual network have less impacts on the descriptive power.

If the network goes even deeper, the mean bias accumulate and difficulties of training increase. Then the benefits of skip connection will dominate that it alleviates gradient vanishing/exploding problems and helps centering the activations in the layers of the net, which enable a deeper network and greatly improve the performance.

Second, fixing the depth of the model, simply broadening the width will improve the performance as showed in Table IV, *e.g.*, $R(16_3, 32_3, 64_3)$ vs $R(32_3, 64_3, 128_3)$, $R(4_3, 8_3, 16_3, 32_3, 64_3)$ vs $R(16_3, 32_3, 64_3, 128_3, 256_3)$. Increasing the filter numbers would enlarge the model capacity which enables modeling more complex nonlinear mappings from LR examples to HR examples.

Third, the deeper the architecture, the better the performance. Adding more residual units, *e.g.*, $R(16_3, 32_3, 64_3, 128_3)$ vs $R(32_3, 64_3, 128_3)$ will improve the performance. Certainly, the depth should be no more than certain limit to avoid the overfitting problem and computational resource limitations. Within this limit, the deeper the better. Our residual unit eases the training difficulties which enables a deeper CNN architecture to improve the situation. On the other side, when model goes deeper as our residual SRResNetNB, plain deep CNN like VDSR fashions can not converge well and the restorations deteriorate. Another attempt to facilitate deeper net is the lightweight design which aims to solve the problem of too many parameters. It will be discussed next.

TABLE IV
PSNR BY THE RESIDUAL MODEL OF DIFFERENT DEPTH AND WIDTH WITH A MAGNIFICATION FACTORS 2 IN SET14.

	$R(16_3, 32_3, 64_3)$	$R(32_3, 64_3, 128_3)$	$R(16_3, 32_3, 64_3, 128_3)$	$R(16_3, 32_3, 64_3, 128_3, 256_3)$	$R(4_3, 8_3, 16_3, 32_3, 64_3)$
PSNR(dB)	32.85	32.96	33.00	33.10	32.91

F. Lightweight Design

In this part, the proposed strategy of gradually varying the ‘shape’ of residual network has been investigated. The performances of different architectures with different shapes have been investigated for our residual net in Table V and VDSR fashions in Table VI counterpart.

The number of featuremap has been gradually varied. To be specific, there are 28 layers as 6 **containers** stack, each **container** contains 2 residual units (2 convolutional layers in each residual unit). The depth can be calculated as $28 = 2 + 6 \times 2 \times 2 + 2$, where feature representation sub-network and reconstruction sub-network each have 2 convolutional layers. For models of VDSR fashions, 12-layer VDSR have been explored. For residual architectures, networks of different ‘shapes’ have achieved comparable results. On the contrary, the performances of VDSR structures have largely fluctuated when the shapes of the networks vary.

Residual networks can be interpreted as an ensemble of models which are the paths of differing depth in the residual network [32]. When the ‘shapes’ of residual models are gradually changing, some short paths in the residual network have been less disturbed as Fig. 3. Thus, the performances of the ensembles are nearly unchanged. On the contrary, the single path VDSR network are more disturbed by the variations of the shape. Instead of keeping the filter number fixed, less parameters can be applied for the residual network with our strategy to achieve comparable performances.

G. Training with Multiple vs Single Upscaling Factors

In this section, we compare the performances of networks handling multiple with respect to single upscaling factors as [22] in Table VII. The training examples from different upscaling factors were mixed together to enable the model handling multiple upscaling factors. It seems mixing samples augmentations strategy [22] from different upscaling factors has slightly boosted the performances, especially for large upscaling factors.

V. CONCLUSION

In this paper, from the perspective of centering activations and ensemble behaviors of residual network, a novel residual deep CNN which takes advantage of skip connections or identity mapping shortcuts in avoiding gradient exploding/vanishing problem was proposed for single image super-resolution. In particular, the ‘shape’ of CNN has been carefully designed such that a very deep convolutional neural network with much fewer parameters can produce even better performance. Based on the investigations into the influences of the network ‘shape’ on the performances, a strategy of gradually varying the ‘shape’ of the network has been proposed to construct this lightweight model. Experimental results have

demonstrated that the proposed method can not only achieve state-of-the-art PSNR and SSIM results for single image super-resolution but also produce visually pleasant results.

ACKNOWLEDGMENT

This work is partially supported by National Science Foundation of China under Grant NO. 61473219.

REFERENCES

- [1] K. Zhang, X. Zhou, H. Zhang, and W. Zuo, “Revisiting single image super-resolution under internet environment: Blur kernels and reconstruction algorithms,” in *Pacific Rim Conference on Multimedia*. Springer, 2015, pp. 677–687. 1
- [2] K. Zhang, B. Wang, W. Zuo, H. Zhang, and L. Zhang, “Joint learning of multiple regressors for single image super-resolution,” *IEEE Signal Processing Letters*, vol. 23, no. 1, pp. 102–106, 2016. 1
- [3] Y. Liang, J. Wang, S. Zhang, and Y. Gong, “Incorporating image degeneration modeling with multitask learning for image super-resolution,” in *Image Processing (ICIP), 2015 IEEE International Conference on*. IEEE, 2015, pp. 2110–2114. 1
- [4] R. Timofte, R. Rothe, and L. Van Gool, “Seven ways to improve example-based single image super resolution,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1865–1873. 1
- [5] C. Dong, C. C. Loy, K. He, and X. Tang, “Image super-resolution using deep convolutional networks,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 38, no. 2, pp. 295–307, 2016. 1, 2, 3
- [6] C.-Y. Yang, C. Ma, and M.-H. Yang, “Single-image super-resolution: A benchmark,” in *Computer Vision—ECCV 2014*. Springer, 2014, pp. 372–386. 1, 6
- [7] J.-B. Huang, A. Singh, and N. Ahuja, “Single image super-resolution from transformed self-exemplars,” in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2015, pp. 5197–5206. 1, 6, 7
- [8] Y. Romano, J. Isidoro, and P. Milanfar, “Rairs: rapid and accurate image super resolution,” *IEEE Transactions on Computational Imaging*, vol. 3, no. 1, pp. 110–125, 2017. 1
- [9] 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*, 2017. 1, 3
- [10] 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 sub-pixel convolutional neural network,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1874–1883. 1
- [11] M. Irani and S. Peleg, “Motion analysis for image enhancement: Resolution, occlusion, and transparency,” *Journal of Visual Communication and Image Representation*, vol. 4, no. 4, pp. 324–335, 1993. 1
- [12] W. T. Freeman, E. C. Pasztor, and O. T. Carmichael, “Learning low-level vision,” *International journal of computer vision*, vol. 40, no. 1, pp. 25–47, 2000. 1, 2
- [13] J. Yang, J. Wright, T. Huang, and Y. Ma, “Image super-resolution as sparse representation of raw image patches,” in *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on*. IEEE, 2008, pp. 1–8. 1, 6
- [14] H. Chang, D.-Y. Yeung, and Y. Xiong, “Super-resolution through neighbor embedding,” in *Computer Vision and Pattern Recognition, 2004. CVPR 2004. Proceedings of the 2004 IEEE Computer Society Conference on*, vol. 1. IEEE, 2004, pp. 1–275. 1
- [15] R. Timofte, V. De, and L. V. Gool, “Anchored neighborhood regression for fast example-based super-resolution,” in *Computer Vision (ICCV), 2013 IEEE International Conference on*. IEEE, 2013, pp. 1920–1927. 1, 4, 6

TABLE V
PERFORMANCE BY DIFFERENT RESIDUAL MODELS WHICH HAVE DIFFERENT SHAPES WITH A MAGNIFICATION FACTOR 2 IN SET14.

residual	$R(16_4, 32_4, 64_4)$	$R(64_4, 32_4, 16_4)$	$R(16_2, 32_2, 64_2, 64_2, 32_2, 16_2)$	$R(64_2, 32_2, 16_2, 16_2, 32_2, 64_2)$
PSNR(dB)	32.91	32.85	32.94	32.89

TABLE VI
PERFORMANCE BY VDSR MODELS WHICH HAVE DIFFERENT SHAPES WITH A MAGNIFICATION FACTOR 2 IN SET14.

VDSR	$(8_2, 16_2, 64_2)$	$(64_2, 16_2, 8_2)$	$(8, 16, 64, 64, 16, 8)$	$(64, 16, 8, 8, 16, 64)$	(64_{16})
PSNR(dB)	32.68	32.59	32.66	32.50	32.85

TABLE VII
PSNR COMPARISONS BETWEEN MODELS HANDLING MULTIPLE VS SINGLE UPSCALING FACTORS, DENOTED AS ‘MULTISCALE’ AND ‘SINGLE SCALE’

		Set5	Set14	BSD100	Urban100
		PSNR/SSIM	PSNR/SSIM	PSNR/SSIM	PSNR/SSIM
× 2	Multiscale	37.51/0.9587	33.10/0.9131	31.91/0.8961	30.88/0.9150
	single scale	37.52/0.9589	33.03/0.9129	31.90/0.8958	30.84/0.9143
× 3	Multiscale	33.72/0.9215	29.80/0.8317	28.83/0.7980	27.17/0.8283
	single scale	33.6/0.9212	29.75/0.8313	28.79/0.7967	27.08/0.8255
× 4	Multiscale	31.37/0.8838	28.06/0.7681	27.29/0.7251	25.22/0.7537
	single scale	31.30/0.8824	27.99/0.7668	27.24/0.7237	25.14/0.7051

- [16] R. Timofte, V. De Smet, and L. Van Gool, “A+: Adjusted anchored neighborhood regression for fast super-resolution,” in *Computer Vision—ACCV 2014*. Springer, 2014, pp. 111–126. **1, 6, 7**
- [17] J. Yang, J. Wright, T. S. Huang, and Y. Ma, “Image super-resolution via sparse representation,” *Image Processing, IEEE Transactions on*, vol. 19, no. 11, pp. 2861–2873, 2010. **1**
- [18] S. Schuler, C. Leistner, and H. Bischof, “Fast and accurate image upscaling with super-resolution forests,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3791–3799. **1, 6, 7, 11**
- [19] J. Salvador and E. Pérez-Pellitero, “Naive bayes super-resolution forest,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 325–333. **1**
- [20] C. Dong, C. C. Loy, K. He, and X. Tang, “Learning a deep convolutional network for image super-resolution,” in *Computer Vision—ECCV 2014*. Springer, 2014, pp. 184–199. **1, 2, 4, 6, 7, 11**
- [21] Y. Liang, J. Wang, S. Zhou, Y. Gong, and N. Zheng, “Incorporating image priors with deep convolutional neural networks for image super-resolution,” *Neurocomputing*, 2016. **1**
- [22] J. Kim, J. K. Lee, and K. M. Lee, “Accurate image super-resolution using very deep convolutional networks,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016. **1, 2, 3, 4, 5, 6, 7, 9, 11**
- [23] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014. **1**
- [24] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” *arXiv preprint arXiv:1512.03385*, 2015. **1, 2, 3, 4, 6**
- [25] —, “Identity mappings in deep residual networks,” *arXiv preprint arXiv:1603.05027*, 2016. **1, 2, 3, 4, 5, 6, 8**
- [26] Y. Le Cun, I. Kanter, and S. A. Solla, “Eigenvalues of covariance matrices: Application to neural-network learning,” *Physical Review Letters*, vol. 66, no. 18, p. 2396, 1991. **2, 3**
- [27] G. B. Orr and K.-R. Müller, *Neural networks: tricks of the trade*. Springer, 2003. **2, 3**
- [28] S.-I. Amari, “Natural gradient works efficiently in learning,” *Neural computation*, vol. 10, no. 2, pp. 251–276, 1998. **2**
- [29] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” in *International Conference on Machine Learning*, 2015, pp. 448–456. **2, 3, 4, 5**
- [30] T. Raiko, H. Valpola, and Y. LeCun, “Deep learning made easier by linear transformations in perceptrons,” in *Artificial Intelligence and Statistics*, 2012, pp. 924–932. **2, 3, 4**
- [31] Z. Yang, K. Zhang, Y. Liang, and J. Wang, “Single image super-resolution with a parameter economic residual-like convolutional neural network,” in *International Conference on Multimedia Modeling*. Springer, 2017, pp. 353–364. **2, 3**
- [32] A. Veit, M. Wilber, and S. Belongie, “Residual networks behave like ensembles of relatively shallow networks,” *Advances in Neural Information Processing Systems*, pp. pages 550–558, 2016. **2, 5, 9**
- [33] Z. Wang, D. Liu, J. Yang, W. Han, and T. Huang, “Deep networks for image super-resolution with sparse prior,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 370–378. **2**
- [34] J. Johnson, A. Alahi, and L. Fei-Fei, “Perceptual losses for real-time style transfer and super-resolution,” in *European Conference on Computer Vision*. Springer, 2016, pp. 694–711. **3**
- [35] 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*, 2017. **3**
- [36] D.-A. Clevert, T. Unterthiner, and S. Hochreiter, “Fast and accurate deep network learning by exponential linear units (elus),” *arXiv preprint arXiv:1511.07289*, 2015. **3, 4**
- [37] B. Lim, S. Son, H. Kim, S. Nah, and K. M. Lee, “Enhanced deep residual networks for single image super-resolution,” in *IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*. IEEE, 2017, pp. 1132–1140. **3**
- [38] R. Timofte, E. Agustsson, L. Van Gool, M.-H. Yang, L. Zhang, B. Lim, S. Son, H. Kim, S. Nah, K. M. Lee *et al.*, “Ntire 2017 challenge on single image super-resolution: Methods and results,” in *Computer Vision and Pattern Recognition Workshops (CVPRW), 2017 IEEE Conference on*. IEEE, 2017, pp. 1110–1121. **3**
- [39] Y. Liang, R. Timofte, J. Wang, Y. Gong, and N. Zheng, “Single image super resolution-when model adaptation matters,” *arXiv preprint arXiv:1703.10889*, 2017. **3, 4**
- [40] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, “Going deeper with convolutions,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9. **3, 5**
- [41] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, “Rethinking the inception architecture for computer vision,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2818–2826. **3, 5**
- [42] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi, “Inception-v4, inception-resnet and the impact of residual connections on learning,” in *AAAI*, 2017, pp. 4278–4284. **3, 5**
- [43] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in neural information processing systems*, 2012, pp. 1097–1105. **5**

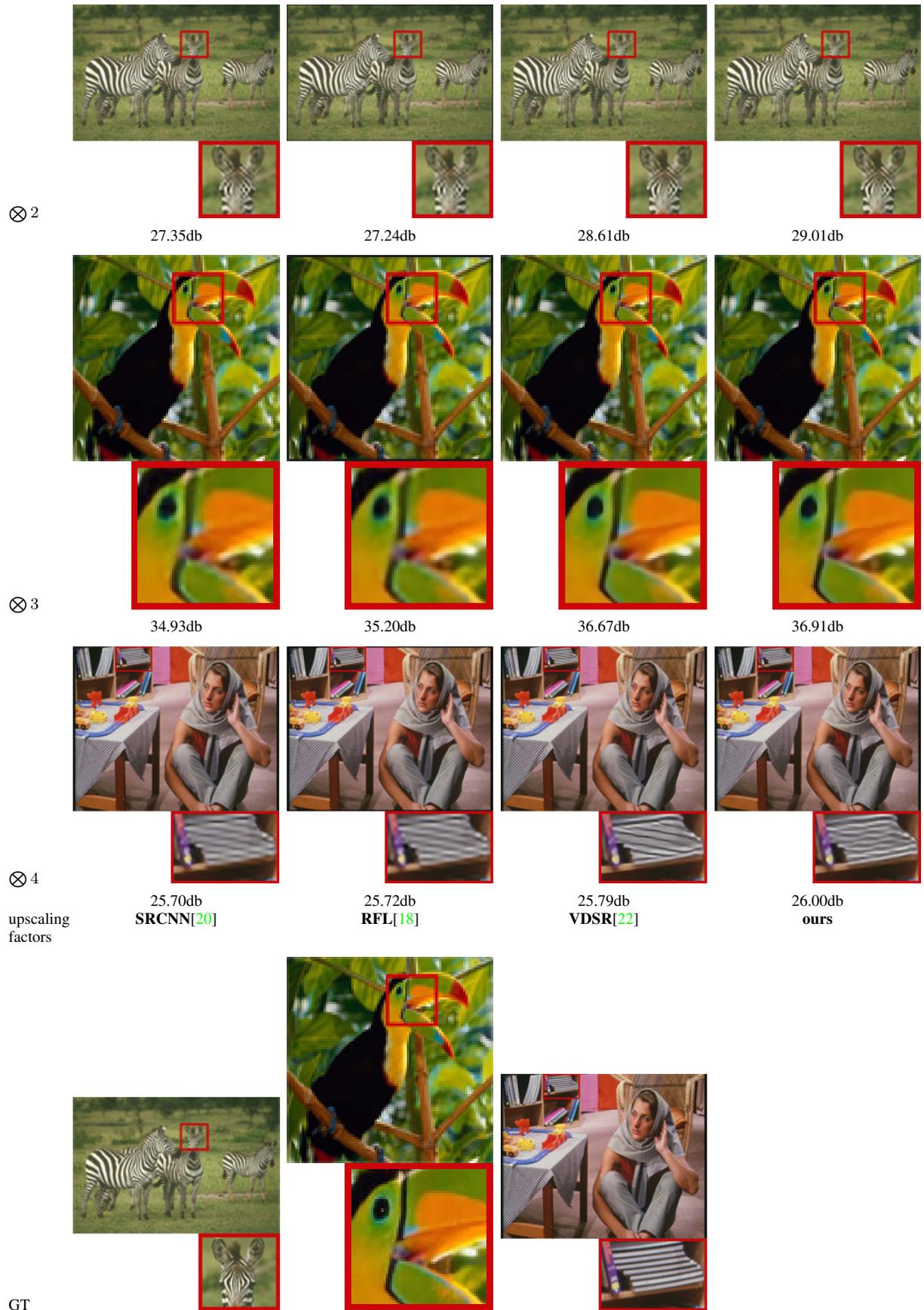


Fig. 8. Comparisons of image SR results with different methods in different upscaling factors