# Single Image Super Resolution based on a Modified U-net with Mixed Gradient Loss

Zhengyang Lu, Ying Chen

*Key Laboratory of Advanced Process Control for Light Industry (Ministry of Education), Jiangnan University*

*Wuxi 214122, Jiangsu, People's Republic of China*

**Abstract**

Single image super-resolution (SISR) is the task of inferring a high-resolution image from a single low-resolution image. Recent research on super-resolution has achieved great progress due to the development of deep convolutional neural networks in the field of computer vision. Existing super-resolution reconstruction methods have high performances in the criterion of Mean Square Error (MSE) but most methods fail to reconstruct an image with shape edges. To solve this problem, the mixed gradient error, which is composed by MSE and a weighted mean gradient error, is proposed in this work and applied to a modified U-net network as the loss function. The modified U-net removes all batch normalization layers and one of the convolution layers in each block. The operation reduces the number of parameters, and therefore accelerates the reconstruction. Compared with the existing image super-resolution algorithms, the proposed reconstruction method has better performance and time consumption. The experiments demonstrate that modified U-net network architecture with mixed gradient loss yields high-level results on three image datasets: SET14, BSD300, ICDAR2003. Code is available online[1]

*Keywords:* Image super-resolution, Network achitecture, Gradient loss.

## 1. Introduction

Significant progress of neural networks for computer vision has been given to the field of Single Image Super-Resolution (SISR). SISR aims to reconstruct a super-resolution image $I^{SR}$ from a single low-resolution image $I^{LR}$. This task, refered as super-resolution, finds direct applications in numerous areas such as medical image processing [1, 2], HDTV [3], face recognition [4, 5], satellite image

---

[1]The project is coded by PyTorch and is proposed on `https://github.com/MnisterLu/simplifiedUnetSR`

processing [6, 7, 8] and surveillance [9, 10]. Furthermore, the super-resolution task can also be expressed as a one-to-many mapping from low-resolution to high-resolution space for each pixel. Super-resolution is an inference problem and the main solution to this problem relies on the statistical model in recent researches [11].

Recently, deep learning based SISR has attracted wide attention. He [12] introduced the first deep learning method for single image super-resolution called Super-Resolution Convolutional Neural Network (**SRCNN**). Then, Fast Super-Resolution Convolutional Neural Network (**FSRCNN**) [13] built a lightweight framework to solve real-time problems of super-resolution. Super-resolution Generative Adversarial Networks (**SRGAN**) [14] was the first Generative Adversarial Networks (**GAN**) [15] for super-resolution task and considered the human subjective evaluation of reconstructed images. Enhanced Deep Residual Networks for Single Image Super-Resolution (**EDSR**) [16] was the winner of the NTIRE2017 Super-Resolution Challenge Competition and achieved an outstanding reconstruction performance. Deep Back-Projection Networks (**DBPN**) [17] exploited iterative up-sampling and down-sampling and providing an error feedback mechanism.

One of the main disadvantages of previous works was that it is difficult to reconstruct clear boundaries and high gradient components. The other was that the architecture of the previous network lacked of skip connection between the convolution layers at the same depth, which caused the information loss. To solve the problem, an improved U-net with a mixed gradient loss is proposed for SISR. The contribution of the new network can be summarized as follows:

1. The improved **U-net** network architecture is proposed for super-resolution on a single image. The **U-net**, which has been approved effective for image segmentation, is modified for SISR task by removing all batch normalization layers and one convolution layer in each block. The operation reduces the computation cost and meanwhile keeps the reconstruction performance. The input image is up-scaled for the larger size and build a new convolution block on the large scale, which has a skip-connection with the output block on the same scale. The direct up-scaled images avoid the errors caused by redundant calculations.

2. A mixed gradient error (MixGE) is proposed and applied to SISR loss function. As a combination of mean square error (MSE) [18] and mean gradient error (MGE), not only pixel error but also gradient error is considered and the improved loss is called MixGE. In the paper, we use the classic gradient calculation method which was proposed by Sobel[19]. The method for MGE computation are presented and analyzed.

3. The modified **U-net** for common scenes task super-resolution yields the outstanding performance over existing methods on SET14 and BSD300 dataset. Moreover, the proposed network successfully outperforms other state-of-art methods on texture tasks such as ICDAR2003 dataset.

## 2. Related Work

The goal of super-resolution methods is to recover a high-resolution image from a low-resolution image [20]. Recently most popular SISR methods were implemented by deep learning network instead of traditional mathematical models. This section will analyze and summarize most effective existing super-resolution reconstruction methods which realized by deep learning methods.

### 2.1. Super-resolution network

For super resolution, the methods can be divided into video super-reslution reconstruction, such as **VESPCN** [21], **VSRCNN** [22], and single image super-resolution, which includes **SRCNN** [12], **FSRCNN** [13], **VDSR** [23], **DRCN** [24], **SRGAN** [14], **ESPCN** [25], **EDSR** [16] and **DBPN** [17].

To better solve the distortion problem of high-resolution image reconstruction, He [12] presented a network architecture named **SRCNN** which was the earliest proposal to solve super-resolution problem through deep learning network. He demonstrated that a convolution network can be used to learn an efficacious mapping from low-resolution to high-resolution in an end-to-end architecture.

Although the **SRCNN** introduced a new way to solve the super-resolution reconstruction problem by deep learning, this method had 2 obvious limitations. First, the observation field was too narrow that can not get enough corresponding information. Second, the training of **SRCNN** was hard to converge and easy to be overfitted.

To solve these two limitations, **FSRCNN** [13] was proposed, which was different from **SRCNN** in 3 aspects. First, **FSRCNN** adopted the original low-resolution image as input instead of the **bicubic** [26] interpolate image in **SRCNN**. Second, The de-convolution layer was added at the end of the network for upsampling. The third and the most efficient part was to adapt smaller filter kernels and a deeper network into super-resolution task.

The **VDSR** [23] came up with a much deeper network. To get a large observation field, **VDSR** chose a convolution kernel of 3×3 in the deep network. Otherwise, very deep networks converged too slowly because of the big number of parameters. In the **VDSR** method, residual learning and gradient clipping were chosen to be the solution to the training problem.

Taking an interpolated image as the input, **DRCN** [24] was composed of three modules, namely embedding network for feature extraction, inference network for nonlinear feature mapping, and reconstruction network for SR image generation. The Inference network was a recursive network, that is, data looped through the layer multiple times. Expanding this loop was equivalent to multiple concatenated convolution layers using the same set of parameters.

Generative adversarial networks (**GAN**) [15], which was proposed by Goodfellow in 2014, was set up for estimating generative models via adversarial process. First **GAN** proposed for super-resolution task was **SRGAN** [14]. The experimental results in this work showed that the generator network trained on MSE loss function could output SR images with high Peak signal-to-noise ratio (PSNR) but over-smoothed. The output images of **SRGAN** had a better visual

effect than other methods without adversarial process. The work proposed a **GAN** that apply a deep residual network (**ResNet**) [27] with skip-connection. For improving the visual effect of the super-resolution reconstructed results, they put forward the perceptual loss instead of the MSE loss function.

To solve the high computational complexity of deep network, **ESPCN** [25] was proposed for super-resolution task with a much higher processing speed than previous methods. The core concept of **ESPCN** was the sub-pixel convolutional layer. The input to the network was the original low-resolution image. After passing through two convolutional layers, the result feature image was the same size as the input image, and the feature channel was $r^2$ where $r$ was the magnification scale. Feature images of size $r^2 \times H \times W$ were re-arranged into high-resolution images of size $1 \times rH \times rW$. The operation of convolution layer re-arrangement greatly improved the efficiency of de-convolution and the **ESPCN** had a better performance on super-resolution tasks.

**EDSR** [16] was the winner of the NTIRE2017 Super-Resolution Challenge Competition. As stated in the paper, the most significant performance improvement of **EDSR** was to remove the batch normalization layers of **SRResNet**, which can expand the size of the model to improve the quality of the results.

Haris [17] proposed Deep Back-Projection Networks (**DBPN**), applying iterative up-sampling and down-sampling and exploiting an error feedback mechanism for projection errors in every block. This work constructed mutually-connected up-sampling and down-sampling blocks, which represented different parts of low-resolution and high-resolution components.
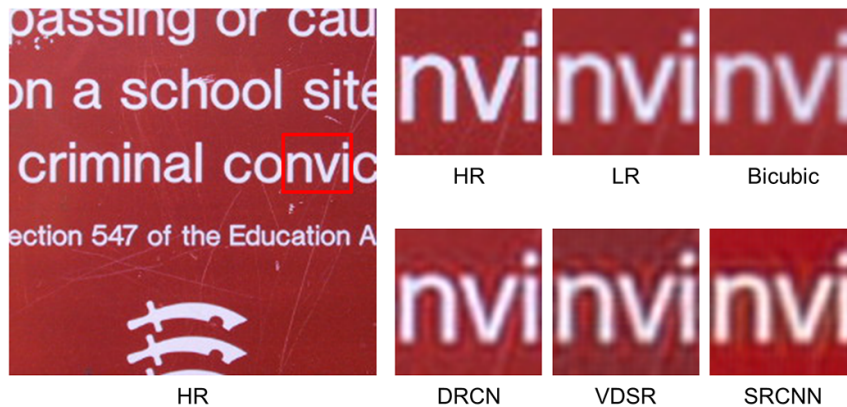
*2.2. Problem analysis*



Figure 1: Super-resolution results by previous works on ICDAR2003 dataset($\times$4).

Previous works on SISR reconstruction task achieved a good performance, but there were still many shortcomings and many parts that lacked reasonable explanations. First, previous studies of SISR had not dealt with the problem of blurred edge, shown as reconstructed results by methods of **SRCNN** and **VDSR** in Figure 1, because these researches were limited to MSE loss function. Though the performance of **SRCNN** on MSE loss function was outstanding, subjective observation of the reconstructed images was blurred. Therefore, gradient components should be considered into the SISR task.

Second, most convolutional layers of SISR deep networks were directly connected which led to a part of data loss in low-dimensional feature layers. High-resolution images which only reconstructed by high-level semantic information lost most basic fine-grained texture details. It is necessary to build a skip-connection between same depth layers to combine high-level semantic information with low-level fine-grained texture details.

Third, according to the impressive analysis of **VDSR** [23] which concluded that better performance could be achieved with deeper network, the number of parameters of the most SISR network were very large due to very deep networks were used to obtain high performance. The networks suffered from large computation cost which makes it difficult to run in real-time.
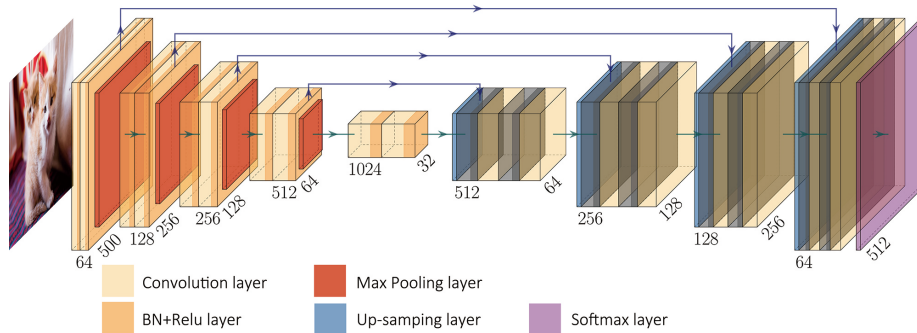
## 3. Method

The task of the single image super-resolution was to reconstruct a super-resolution image $I^{SR}$ from a low-resolution image $I^{LR}$. To solve this single image super-resolution problem, network architecture is constructed where input image $I^{LR}$ is represented as real-valued tensors of size $H \times W \times C$ and output image $I^{SR}$ represented as a tensor of size $rH \times rW \times C$.

This work proposes an improved network called modified U-net and a mixed gradient loss function which combines with Mean Square Error and mean gradient error.
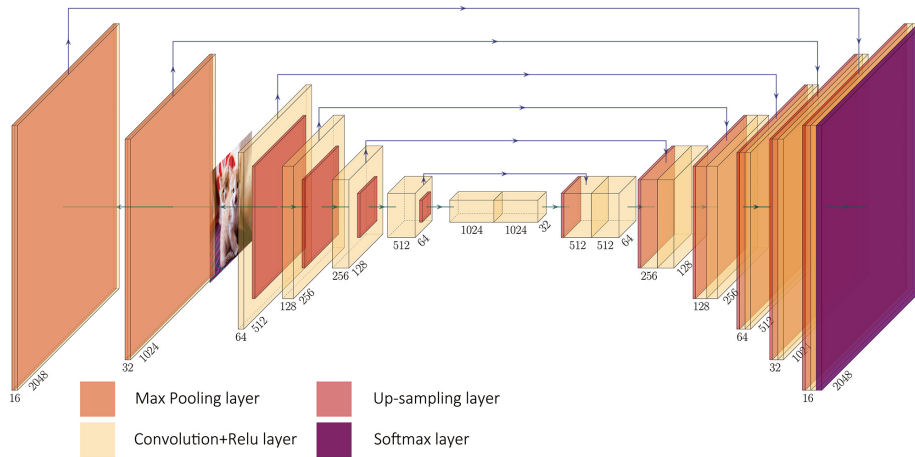
### 3.1. Modified U-net

**U-net** [28] was first proposed by Ronneberger and widely used in the field of semantic segmentation. Ronneberger showed that the **U-net** can be trained as a end-to-end system from extremely few images. Meanwhile, this network architecture surpassed the previous state-of-the-art method on the ISBI challenge for semantic segmentation of neuronal structures.

The **U-net** architecture is illustrated in Figure 2(a). The left side is the contracting path which is used for extracting features and the right side is called expansive path for decoding. The contracting path contains the continuous part of two 3×3 kernels and each part is followed by a Rectified Linear Unit (ReLU) layer and then a 2×2 max-pooling operation with stride 2 for down-sampling. So downscale of each contracting block is 2. Each step in the expansive path includes an up-sampling of the feature map, which followed by a 2×2 convolution that halves the number of feature channels, and two 3×3 convolution kernels, each followed by a ReLU layer.

(a) U-net network architecture.



(b) modified U-net network architecture.

Figure 2: Original U-net and modified U-net network architecture

The most important contribution of **U-net** is to put forward the skip-connection. The **U-net** network architecture is proposed after the Fully Convolutional Network (**FCN**)[29] and modified in a creative method that it yields the state-of-the-art segmentation in medical image processing at that time. Compared with **FCN**, the first main difference is **U-net** is symmetric and the second one is the skip connections between the down-sampling path and the up-sampling path employ a concatenation operation instead of a direct sum. Skip connections in the network aim to provide the original local information to the global information when up-sampling. Because of the symmetry of **U-net** network architecture, the **U-net** has numerous feature maps in up-sampling blocks that allows the efficacious transmission of information.

As illustrated in Figure 2(b), the modified **U-net** designed in this work for SISR has three differences comparing with the original **U-net**. First, it removes all batch normalization layers and one of the convolution layers in each block.

The reason for this is that super-resolution reconstruction is a pixel-level task because the solution to interpolation problem is mainly consider pixels in a certain area. Directly magnified images avoid the errors caused by redundant calculations, which can get closer to real results.

Second, the input image is up-scaled for the larger size and build a new convolution block on the larger scale, which has a skip-connection with the output block on the same scale. The direct up-scaled images avoid the errors caused by redundant calculations, which can get much closer to ground truth. Each addition of upscale layer means expand the size into twice. In other words, the network should have 3 upscale layers if upscale size is 8 and similarly 2 upscale layers for upscale size is 4.

Third, the depth of the original **U-net** is fixed to 4, which means there are 4 downscale blocks and corresponding 4 upscale blocks. Previous work [28] has shown that the deeper network leads to the better performance and higher computation cost in most situations. In our work, the depth of the modified **U-net** for SISR is discussed in section 4, showing a trade off between reconstruction accuracy and computation cost.

*3.2. Mixed Gradient Error*

The aim of SISR is to learn a mapping function $f$ for generating an high-resolution image $\hat{Y} = f(X)$ from a low-resolution image $X$ that is close to the ground truth image $Y$. MSE [30] is widely used as most loss function, which measures the average of the squares of each pixel errors in super-resolution reconstruction. MSE can be shows as follows:

$$MSE = \frac{1}{n}\frac{1}{m}\sum_{i=1}^{n}\sum_{j=1}^{m}(Y(i,j) - \hat{Y}(i,j))^2 \tag{1}$$

where $n$ is the number of horizontal pixels and $m$ is the number of vertical pixels.

To solve the gradient error measurement problem, we introduce classic gradients to the SISR loss function. In our work, Sobel operator[19] is used for gradient calculation, that is,

$$G_x = Y * \begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{bmatrix} \tag{2}$$

where $*$ is the convolution operation.

The gradient map $G$ in $y$ direction of the ground truth image $Y$ shows below:

$$G_y = Y * \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix} \tag{3}$$

Then we combine the gradient value of $x$ and $y$ direction as follows:

7

$$G(i,j) = \sqrt{G_x^2(i,j) + G_y^2(i,j)} \tag{4}$$

Meanwhile, the gradient map $\hat{G}$ can be calculated in the same way. The aim of measure the gradient error is to learn a sharp edge which is close to the groud truth edge. The Mean Gradient Error (MGE) shows as follows:

$$MGE = \frac{1}{n}\frac{1}{m}\sum_{i=1}^{n}\sum_{j=1}^{m}(G(i,j) - \hat{G}(i,j))^2 \tag{5}$$

After achieving the MGE, it should be emphasized in the mixed gradient error. As the main component of mixed gradient error, Mean Square Error forms a Mixed Gradient Error (MixGE) by adding Mean Gradient Error with a weight of $\lambda_G$.

$$MixGE(Y,\hat{Y}) = MSE + \lambda_G MGE \tag{6}$$

## 4. Experiment

### 4.1. Dataset

ICDAR2003 [31] is a dataset of the ICDAR Robust Reading Competition for text detection and recognition tasks. For the single image super-resolution task, there is not a commonly used dataset. The ICDAR2003 dataset consists of 258 training images and 249 testing images, which contains texts in most of the complex circumstances in common life. Because of the resolution of images varies from 422×102 to 640×480, we resize them into 224×224 with **bicubic** interpolation. This network is also compared with other existing methods on standard benchmark dataset: SET14 [32], BSD300 [33].

### 4.2. Evaluation method

Two commonly used performance metric are employed for evaluation and comparsion: PSNR [34] and Structural Similarity Index (SSIM) [35]. The super-resolution results are evaluated with the criteria of PSNR [34] and SSIM [35] on three channels in RGB colour space. The criterion of PSNR is based on the error between each corresponding pixels. Because this criterion only consider the numerical error, most high PSNR results do not have good visual performance. Therefore, the criterion of SSIM [35] observes the distortion of the image by comparing the changes in the image structure information, thereby obtaining an objective quality evaluation. The criteria of PSNR and SSIM are all based on luminance. The higher the value of these criteria, the better the performance of image reconstruction.

Two criteria are described as follows. Let $Y$ donate the ground truth and $\hat{Y}$ donate the reconstructed high-resolution images respectively.

8

$$MSE = \frac{1}{MN} \sum_{i=1}^{M} \sum_{j=1}^{N} (\hat{Y}(i,j) - Y(i,j))^2 \tag{7}$$

$$PSNR(Y, \hat{Y}) = 10 log_{10} \frac{255^2}{MSE} \tag{8}$$

The criterion of SSIM between patches $P_{\hat{Y}}$ and $P_Y$ at the same location on ground truth images $\hat{Y}$ and reconstructed high-resolution image $Y$ is defined as

$$SSIM(P_{\hat{Y}}, P_Y) = \frac{(2\mu_{P_{\hat{Y}}}\mu_{P_Y} + c_1)(2\sigma_{P_{\hat{Y}}}\sigma_{P_Y} + c_2)}{(\mu_{P_{\hat{Y}}}^2 + \mu_{P_Y}^2 + c_1)(\sigma_{P_{\hat{Y}}}^2 + \sigma_{P_Y}^2 + c_2)} \tag{9}$$

where $\mu_{P_{\hat{Y}}}$ and $\mu_{P_Y}$ are the mean of patch $P_{\hat{Y}}$ and $P_Y$ respectively. Meanwhile, $\sigma_{P_{\hat{Y}}}$ and $\sigma_{P_Y}$ are the deviation of patch $P_{\hat{Y}}$ and $P_Y$. $c_1$ and $c_2$ are small constants. Then, The criterion of $SSIM(\hat{Y}, Y)$ is the average of patch-based SSIM over the image.

### 4.3. Implement Details

These existed SISR methods are evaluated on 3 common used dataset: SET14, BSD300, ICDAR2003. SET14 and BSD300 consist of natural scenes and ICDAR2003 contain various types of texts in a robust common scene. The ground truth high-resolution images are downscaled by **bicubic** interpolation to generate low-resolution and high-resolution image pairs for training and testing on Table.1. We convert all images into RGB colour space and processing the data into three channels.

Table 1: Image size of different scales.

| Scale | Image size | |
|---|---|---|
| | LR | HR |
| ×2 | 112×112 | 224×224 |
| ×4 | 56×56 | 224×224 |
| ×8 | 28×28 | 224×224 |

In the training parameter set, the batch of data is set to 1. Our model is trained by Adam optimizer [36] with $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$. The learning rate is set to $10^{-3}$ initially and decreases to half every 25 epoch. The PyTorch implement our models with one RTX 2080 GPU and code is presented online.

*4.4. Results*

*4.4.1. Network analysis*

Before comparison with other state-of-the-art methods, some parameters of the modified **U-net** and MixGE loss function need to be determined. The first part is the depth of modified **U-net**, which influences the performance of super-resolution results and computational complexity. The second one is the MGE weight for MixGE loss function. 50 randomly selected images from training set are taken as validation data for parameter determination..

**Depth analysis.** Figure 3 shows the PSNR with increasing depth of U-net, in which the numbers of the networks parameter are shown above each bar. It can be seen from the figure that the performance has a significant improvement from the depth of 2 to 5, while has a slight improvement from the depth of 5 to 8 with large increase in model complexity. Therefore, the depth of 5 is a proper choice that keep a trade-off between reconstruction accuracy and computation cost.
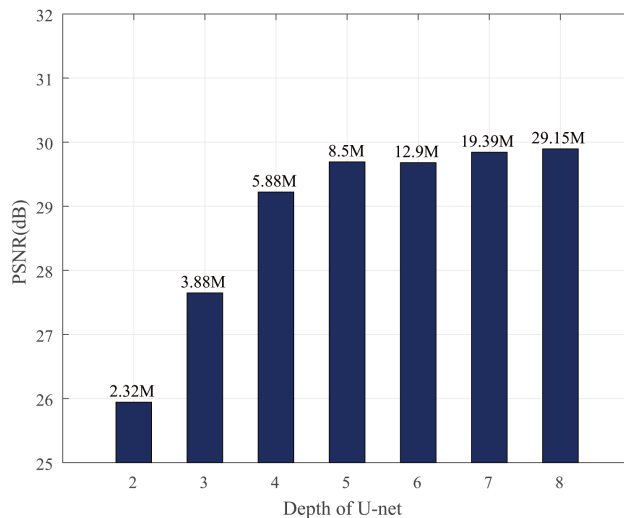


Figure 3: PSNR comparison with different depths of U-net on BSD300(2×) dataset.

**Loss analysis.** For the comparison of MixGE and MSE loss, Figure 4 shows that the results have a great improvement on each dataset if considering the Mean Gradient Error. The method with MixGE loss performs better than that with MSE loss.

In this experiment, the independent variable, $\lambda_G$, was set as $10^{-4}$, $10^{-3}$, $10^{-2}$, $10^{-1}$ and 1. It can be clearly seen from the Figure 4 that the performance is getting better with the increasing value of $\lambda_G$ from $10^{-4}$ to $10^{-1}$ and achieves the peak of performance on $10^{-1}$. This shows that it will greatly improve the performance when choosing the MSE as main component and the MGE as an

important auxiliary component. The best performance reaches the PSNR of 29.4dB when the weight equals 0.1.

It can be seen that MSE loss function is still an irreplaceable loss component in MixGE. MGE becomes an auxiliary component in MixGE to support deep networks to build sharp-edged, gradient-accurate reconstructed images.
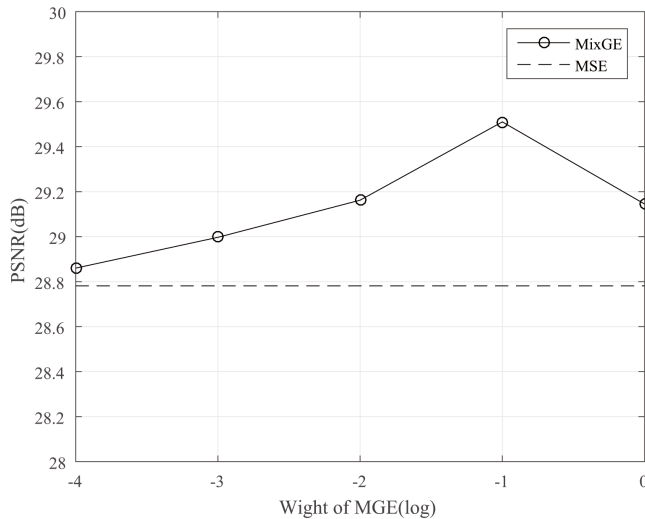
[th]



Figure 4: PSNR comparison between MSE loss and MixGE loss with different weights on BSD300(2×) dataset.

*4.4.2. Comparison with the state-of-the-arts*

To assess the performance of our method, we compare the modified **U-net** architecture with existing super-resolution deep learning network, including **SRCNN** [12], **FSRCNN** [13], **VDSR** [23], **DRCN** [24], **SRGAN** [14], **ES-PCN** [25], **EDSR** [16] and **DBPN** [17]. In this experiment, the performance of **bicubic** [26] interpolation is chosen to be the baseline method for evaluating performance of deep learning method.

Table.2 shows the super-resolution performance of different methods on SET14, BSD300, ICDAR2003 dataset, where **UnetSR** denotes the modified **U-net** with MSE and **UnetSR+** denotes the modified **U-net** with MixGE. In the Table.2, the best performance is marked in red and the second best performance is marked in blue in each row. It can seen clearly from the table that **UnetSR+** has the highest PSNR, which means the best performance, at upscale size 2 and 8 on ICDAR2003 dataset. Meanwhile, **UnetSR+** has the second best performance when upscale size is 4.

11

Table 2: Number of parameters comparison between different network architectures.

| Method | Scale | SET14 | | BSD300 | | ICDAR2003 | |
|--------|-------|-------|------|--------|------|-----------|------|
| | | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| Bicubic [26] | ×2 | 24.4523 | 0.8482 | 26.6538 | 0.7924 | 32.9327 | 0.9028 |
| ESPCN [25] | ×2 | 26.7606 | 0.8999 | 28.9832 | 0.8732 | 35.6041 | 0.9243 |
| SRCNN [12] | ×2 | 25.9711 | 0.8681 | 28.6943 | 0.8671 | 35.2711 | 0.9234 |
| VDSR [23] | ×2 | 28.6617 | 0.9269 | 29.3889 | 0.8785 | 36.2323 | 0.9375 |
| EDSR [16] | ×2 | 24.0624 | 0.8383 | 28.3119 | 0.8621 | 34.5047 | 0.9258 |
| FSRCNN [13] | ×2 | 23.1284 | 0.8123 | 28.7534 | 0.8681 | 35.0533 | 0.9355 |
| DRCN [24] | ×2 | 24.4234 | 0.8458 | 27.5089 | 0.8088 | 33.7849 | 0.9185 |
| SRGAN [14] | ×2 | 23.9553 | 0.8195 | 28.7072 | 0.8633 | 33.2834 | 0.9135 |
| DBPN [17] | ×2 | 28.4092 | 0.9202 | 29.8675 | 0.8834 | 36.2344 | 0.9401 |
| UnetSR | ×2 | 26.7241 | 0.8735 | 29.4241 | 0.8813 | 35.7147 | 0.9388 |
| UnetSR+ | ×2 | 28.3965 | 0.9198 | 29.8403 | 0.8816 | 37.3673 | 0.9675 |
| Bicubic [26] | × 4 | 19.7167 | 0.6089 | 23.5053 | 0.6157 | 28.1135 | 0.7875 |
| ESPCN [25] | × 4 | 20.6292 | 0.6333 | 24.4899 | 0.6641 | 29.4861 | 0.8214 |
| SRCNN [12] | × 4 | 20.5825 | 0.6288 | 24.2232 | 0.6597 | 28.1906 | 0.7661 |
| VDSR [23] | × 4 | 21.4763 | 0.6991 | 24.7077 | 0.6816 | 30.5267 | 0.8321 |
| EDSR [16] | × 4 | 19.9784 | 0.6269 | 23.9192 | 0.6513 | 27.9723 | 0.7101 |
| FSRCNN [13] | × 4 | 19.3255 | 0.5941 | 24.2499 | 0.6599 | 28.0231 | 0.7652 |
| DRCN [24] | × 4 | 19.7077 | 0.6078 | 23.3462 | 0.6132 | 27.7174 | 0.7764 |
| SRGAN [14] | × 4 | 19.3877 | 0.5976 | 24.1675 | 0.6485 | 27.5605 | 0.7654 |
| DBPN [17] | × 4 | 21.7657 | 0.7171 | 25.0644 | 0.6967 | 29.8832 | 0.8224 |
| UnetSR | × 4 | 20.8891 | 0.6693 | 24.8332 | 0.6843 | 29.3374 | 0.8202 |
| UnetSR+ | × 4 | 21.6825 | 0.7112 | 24.9522 | 0.6901 | 31.8966 | 0.8898 |
| Bicubic [26] | × 8 | 16.1132 | 0.3673 | 21.3115 | 0.4933 | 24.3856 | 0.6831 |
| ESPCN [25] | × 8 | 16.3441 | 0.3628 | 21.6447 | 0.5064 | 25.1132 | 0.6964 |
| SRCNN [12] | × 8 | 16.3853 | 0.3614 | 21.8101 | 0.5075 | 22.6281 | 0.6103 |
| VDSR [23] | × 8 | 16.7994 | 0.4095 | 21.9697 | 0.5181 | 25.6303 | 0.7104 |
| EDSR [16] | × 8 | 15.7257 | 0.3209 | 21.6573 | 0.5067 | 23.5578 | 0.5987 |
| FSRCNN [13] | × 8 | 14.5788 | 0.2541 | 21.3311 | 0.5011 | 22.5721 | 0.6155 |
| DRCN [24] | × 8 | 16.1497 | 0.3685 | 21.2771 | 0.4934 | 24.2561 | 0.6725 |
| SRGAN [14] | × 8 | 15.7133 | 0.3221 | 21.8766 | 0.5121 | 23.5621 | 0.6425 |
| DBPN [17] | × 8 | 16.7398 | 0.4122 | 22.0577 | 0.5229 | 26.3482 | 0.7196 |
| UnetSR | × 8 | 16.7001 | 0.4093 | 21.9865 | 0.5231 | 25.7734 | 0.7106 |
| UnetSR+ | × 8 | 17.8289 | 0.4103 | 22.0368 | 0.5235 | 28.2512 | 0.8101 |

We average values of PSNR and SSIM of different datasets and upscale sizes for each method. The results are shown in Table 3, together with parameter numbers of each network. It is apparent from the Table.3 that the proposed **UnetSR+** achieves the best reconstruction accuracy on average. Commonly, the ISIR performance improves with more network size, while the **UnetSR+**

only has 36% parameter numbers of **DBPN**, but the experimental performance exceeds **DBPN** 2.25% on PSNR and 2.47% on SSIM. In addition, it can be seen from the table that the proposed **UnetSR+** shows a distinguished trade-off between reconstruction accuracy and model complexity. For example, comparing with **SRGAN**, **DBPN** increased model size by more than 2 times with PSNR increased 8.3% and SSIM increased 9.04%, while the **UnetSR+** improves 10.76% on PSNR and 11.73% on SSIM with parameters increased only 30%.

Table 3: Average results and number of parameters comparison between different network architectures.

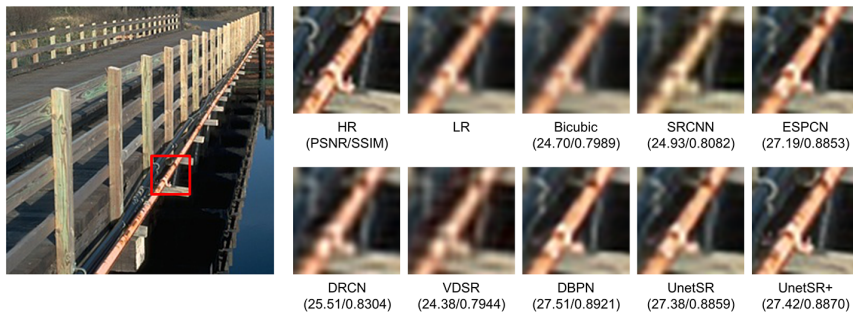| Method | Number of parameters (M) | PSNR (dB) | SSIM |
|--------|--------------------------|-----------|------|
| Bicubic | None | 24.1316 | 0.6777 |
| ESPCN | 0.08 | 25.4506 | 0.7091 |
| SRCNN | 0.17 | 24.8618 | 0.6880 |
| VDSR | 0.22 | 26.1548 | 0.7326 |
| EDSR | 0.78 | 24.4100 | 0.6712 |
| FSRCNN | 0.03 | 24.1128 | 0.6673 |
| DRCN | 0.11 | 24.2413 | 0.6783 |
| SRGAN | 6.54 | 24.2460 | 0.6761 |
| DBPN | 23.21 | 26.2633 | 0.7372 |
| UnetSR | 8.50 | 25.7092 | 0.7234 |
| UnetSR+ | 8.50 | 26.8546 | 0.7554 |



Figure 5: Super-resolution results on BSD300 dataset(×2).

For subjective visual evaluation, we compare reconstruction examples with different deep learning methods of SISR for common scenes task on SET14 and BSD300 dataset in Figures 5 and 6 and examples for texture task on ICDAR2003 dataset in Figures 7 and 8. In Figure 5, super-resolution results
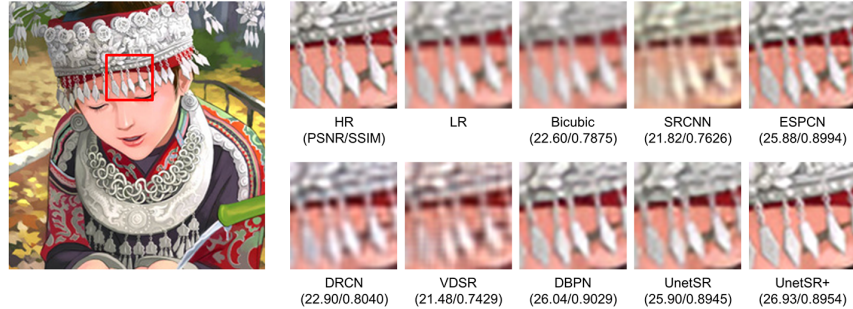
Figure 6: Super-resolution results on SET14 dataset(×2).

on BSD300 dataset for common scenes tasks with the scale of 2 are shown with up-scaled details on the right. As illustrated in the figure, the iron pin can be figured out in four methods, which is **ESPCN**, **DBPN**, **UnetSR** and **UnetSR+**. However, the pip strap, a much smaller detail, only reconstructed perfectly by the method of **UnetSR+**. It is also clear from the Figure 6 that though methods of **ESPCN**, **DBPN**, **UnetSR** and **UnetSR+** reconstruct four distinguishable pendants, only the method of **UnetSR+** provides clear and proper edges for pedants.
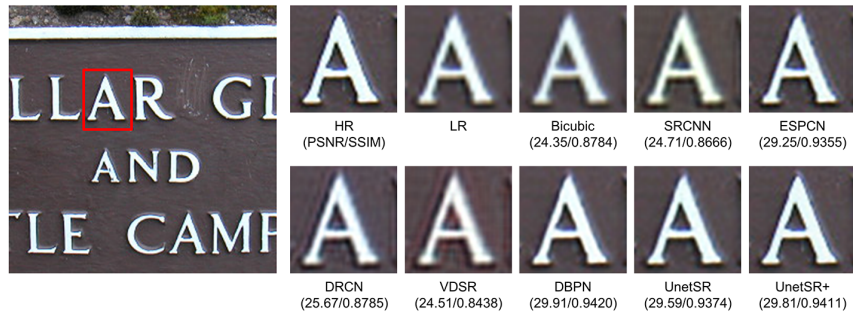


Figure 7: Super-resolution results on ICDAR2003 dataset(×4).

For texture task, super-resolution results on ICDAR2003 dataset on the scale of 4 are shown in Figure 7. The left original image has clear convex white characters on a uniform dark brown background. Under this classic super-resolution scene, the reconstructed images of **Bicubic**, **SRCNN**, **ESPCN**, **DRCN** and **VDSR** produced a serious overlap phenomenon. Contrary to
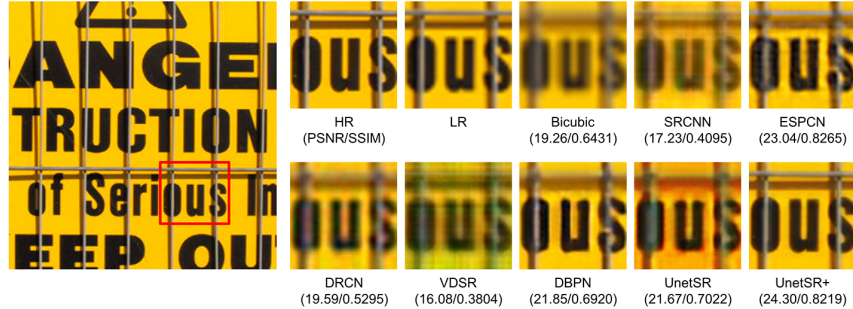
14

Figure 8: Super-resolution results on ICDAR2003 dataset($\times$8).

these results, each image reconstructed by methods of **DBPN**, **UnetSR** and **UnetSR+** produced a visible border. The method of **UnetSR+** not only restores the sharp border, but also even restores the convexity of the original character. In Figure 8, it shows super-resolution results on ICDAR2003 dataset with the scale of 8 under a more complex scene. The methods of **ESPCN**, **DBPN** and **UnetSR+** are still able to provide distinguishable reconstruction results, in which the **UnetSR+** restores almost all the details and boundaries with sharper edges.

## 5. Conclusion

**U-net** for image segmentation is modified in the paper and applied to SISR. First, we remove the superfluous part which includes batch normalization layers and one of the convolution layers. Second, the input image is up-scaled for the larger size and the direct up-scaled images avoid the errors caused by redundant calculations. The modified **U-net** network achieve a high value of PNSR with the trade-off the depth of 5. Furthermore, a Mixed Gradient Loss, which is combined with Mean Square Error and Mean Gradient Error, is proposed for sharp edge reconstruction. Experiments show that the proposed network successfully outperforms other state-of-art methods on SET14, BSD300 and ICDAR2003 datasets. In view of the outstanding performance of the work on texture reconstruction task, it is readily applicable to low-resolution texture detection and recognition in future work.

## Acknowledgments

# References

[1] H. Greenspan, G. Oz, N. Kiryati, and S. Peled, "Mri inter-slice reconstruction using super-resolution," *Magnetic resonance imaging*, vol. 20, no. 5, pp. 437–446, 2002.

[2] S. Peled and Y. Yeshurun, "Superresolution in mri: application to human white matter fiber tract visualization by diffusion tensor imaging," *Magnetic Resonance in Medicine: An Official Journal of the International Society for Magnetic Resonance in Medicine*, vol. 45, no. 1, pp. 29–35, 2001.

[3] T. Goto, T. Fukuoka, F. Nagashima, S. Hirano, and M. Sakurai, "Super-resolution system for 4k-hdtv," in *2014 22nd International Conference on Pattern Recognition*, pp. 4453–4458, IEEE, 2014.

[4] B. K. Gunturk, A. U. Batur, Y. Altunbasak, M. H. Hayes, and R. M. Mersereau, "Eigenface-domain super-resolution for face recognition," *IEEE transactions on image processing*, vol. 12, no. 5, pp. 597–606, 2003.

[5] F. W. Wheeler, X. Liu, and P. H. Tu, "Multi-frame super-resolution for face recognition," in *2007 First IEEE International Conference on Biometrics: Theory, Applications, and Systems*, pp. 1–6, IEEE, 2007.

[6] M. W. Thornton, P. M. Atkinson, and D. Holland, "Sub-pixel mapping of rural land cover objects from fine spatial resolution satellite sensor imagery using super-resolution pixel-swapping," *International Journal of Remote Sensing*, vol. 27, no. 3, pp. 473–491, 2006.

[7] A. J. Tatem, H. G. Lewis, P. M. Atkinson, and M. S. Nixon, "Super-resolution target identification from remotely sensed images using a hopfield neural network," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 39, no. 4, pp. 781–796, 2001.

[8] A. J. Tatem, H. G. Lewis, P. M. Atkinson, and M. S. Nixon, "Super-resolution land cover pattern prediction using a hopfield neural network," *Remote Sensing of Environment*, vol. 79, no. 1, pp. 1–14, 2002.

[9] L. Zhang, H. Zhang, H. Shen, and P. Li, "A super-resolution reconstruction algorithm for surveillance images," *Signal Processing*, vol. 90, no. 3, pp. 848–859, 2010.

[10] F. C. Lin, C. B. Fookes, V. Chandran, and S. Sridharan, "Investigation into optical flow super-resolution for surveillance applications," 2005.

[11] Z. Wang, J. Chen, and S. C. Hoi, "Deep learning for image super-resolution: A survey," *arXiv preprint arXiv:1902.06068*, 2019.

[12] 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, 2015.

[13] C. Dong, C. C. Loy, and X. Tang, "Accelerating the super-resolution convolutional neural network," in *European conference on computer vision*, pp. 391–407, Springer, 2016.

[14] 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*, pp. 4681–4690, 2017.

[15] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, pp. 2672–2680, 2014.

[16] 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*, pp. 136–144, 2017.

[17] 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*, pp. 1664–1673, 2018.

[18] Y. Ephraim and D. Malah, "Speech enhancement using a minimum-mean square error short-time spectral amplitude estimator," *IEEE Transactions on acoustics, speech, and signal processing*, vol. 32, no. 6, pp. 1109–1121, 1984.

[19] N. Kanopoulos, N. Vasanthavada, and R. L. Baker, "Design of an image edge detection filter using the sobel operator," *IEEE Journal of solid-state circuits*, vol. 23, no. 2, pp. 358–367, 1988.

[20] D. G. S. B. M. Irani, "Super-resolution from a single image," in *Proceedings of the IEEE International Conference on Computer Vision, Kyoto, Japan*, pp. 349–356, 2009.

[21] J. Caballero, C. Ledig, A. Aitken, A. Acosta, J. Totz, Z. Wang, and W. Shi, "Real-time video super-resolution with spatio-temporal networks and motion compensation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4778–4787, 2017.

[22] A. Kappeler, S. Yoo, Q. Dai, and A. K. Katsaggelos, "Video super-resolution with convolutional neural networks," *IEEE Transactions on Computational Imaging*, vol. 2, no. 2, pp. 109–122, 2016.

[23] 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*, pp. 1646–1654, 2016.

[24] 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*, pp. 1637–1645, 2016.

[25] 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*, pp. 1874–1883, 2016.

[26] C. De Boor, "Bicubic spline interpolation," *Journal of mathematics and physics*, vol. 41, no. 1-4, pp. 212–218, 1962.

[27] 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*, pp. 770–778, 2016.

[28] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*, pp. 234–241, Springer, 2015.

[29] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3431–3440, 2015.

[30] D. M. Allen, "Mean square error of prediction as a criterion for selecting variables," *Technometrics*, vol. 13, no. 3, pp. 469–475, 1971.

[31] D. Karatzas, F. Shafait, S. Uchida, M. Iwamura, L. G. i Bigorda, S. R. Mestre, J. Mas, D. F. Mota, J. A. Almazan, and L. P. De Las Heras, "Icdar 2013 robust reading competition," in *2013 12th International Conference on Document Analysis and Recognition*, pp. 1484–1493, IEEE, 2013.

[32] M. Bevilacqua, A. Roumy, C. Guillemot, and M. L. Alberi-Morel, "Low-complexity single-image super-resolution based on nonnegative neighbor embedding," 2012.

[33] D. Martin, C. Fowlkes, D. Tal, J. Malik, *et al.*, "A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics," Iccv Vancouver:, 2001.

[34] Q. Huynh-Thu and M. Ghanbari, "Scope of validity of psnr in image/video quality assessment," *Electronics letters*, vol. 44, no. 13, pp. 800–801, 2008.

[35] A. Hore and D. Ziou, "Image quality metrics: Psnr vs. ssim," in *2010 20th International Conference on Pattern Recognition*, pp. 2366–2369, IEEE, 2010.

[36] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.