# Channel Attention and Multi-level Features Fusion for Single Image Super-Resolution

1st Yue Lu
*Beijing University of*
*Posts and Telecommunications*
Beijing, China
lu_yue@bupt.edu.cn

2nd Yun Zhou
*Academy of Broadcasting Science*
Beijing, China
zhouyung@abs.ac.cn

3rd Zhuqing Jiang
*Beijing University of*
*Posts and Telecommunications*
Beijing, China
jiangzhuqing@bupt.edu.cn

4th Xiaoqiang Guo
*Academy of Broadcasting Science*
Beijing, China
guoxiaoqiang@abs.ac.cn

5th Zixuan Yang
*Beijing University of*
*Posts and Telecommunications*
Beijing, China
hsuanr@qq.com

*Abstract*—**Convolutional neural networks (CNNs) have demonstrated superior performance in super-resolution (SR). However, most CNN-based SR methods neglect the different importance among feature channels or fail to take full advantage of the hierarchical features. To address these issues, this paper presents a novel recursive unit. Firstly, at the beginning of each unit, we adopt a compact channel attention mechanism to adaptively recalibrate the channel importance of input features. Then, the multi-level features, rather than only deep-level features, are extracted and fused. Additionally, we find that it will force our model to learn more details by using the learnable upsampling method (i.e., transposed convolution) only on residual branch (instead of using it both on residual branch and identity branch) while using the bicubic interpolation on the other branch. Analytic experiments show that our method achieves competitive results compared with the state-of-the-art methods and maintains faster speed as well.**

*Index Terms*—**Super-Resolution, Convolutional Neural Networks, Recursive Unit, Channel Attention, Multi-level Features Fusion**

## I. INTRODUCTION

The aim of Single Image Super-Resolution (SISR) [1] is to recover a high resolution (HR) image from its corresponding low resolution (LR) input image. SISR is widely used in computer vision applications such as security and surveillance imaging, satellite imaging and medical imaging.

Since deep learning has fundamentally changed how computers learn features, the field of SISR makes impressive improvements by using Convolutional Neural Networks in recent years. Dong et al. [2] firstly proposed a fully convolutional neural network, termed SRCNN, which applies SR after the bicubic interpolation operation. To improve the representational power of SR network, Kim et al. [3] successfully trained a 20 layers model (VDSR) to learn the global residual instead of the actual whole image, achieving vast improvements over SRCNN. Kim et al. [4] further proposed a deeply-recursive convolutional network (DRCN) to efficiently reuse weight parameters while exploiting a large image context. Lai et al.

[5] proposed LapSRN to progressively predict residual image in a coarse-to-fine manner with Charbonnier loss, striking a balance between reconstruction accuracy and execution time. In DRRN, Tai et al. [6] built a very deep network (up to 52 layers) with deep recursions for SR.

While these methods have significantly improved the performance of SISR, there remain three issues to be noticed:

First, a convolutional feature channel often corresponds to a specific functionality like texture extraction or intensity detection. Therefore, in certain recursion some feature channels are more significant than others [7]. However, most of existing CNN-based SR methods treat the channel relationship equally without considering different importance or utilizing channel attention to flexibly adjust features.

Second, when features flow in a network, they can be roughly divided into original, shallow-level and deep-level features. These hierarchical features rely on different receptive fields and therefore carry diverse information. Most of existing methods, however, do not fully exploit the hierarchical features or simply combine these features in a element-wise summation manner. Despite some methods [4]–[6] adopt local skip connections, their basic recursion is still a plain network (i.e., cascade of convolutional layers). The same shortage is also observed in many other computer vision tasks [8], [9].

Third, LapSRN [5] claims that using interpolation to upscale input images to the desired spatial resolution increases unnecessary computational cost and often results in visible reconstruction artifacts. Therefore, it employs transposed convolution both on the residual branch and identity branch, making two branches become learnable. However, we experimentally find that this strategy not only fails to explicitly learn sophisticated high-frequency residuals which contain more image details, but also fluctuates the training process.

To address these issues, we introduce a novel recursive unit into our model as shown in Fig. 1. Our main contribution is threefold:
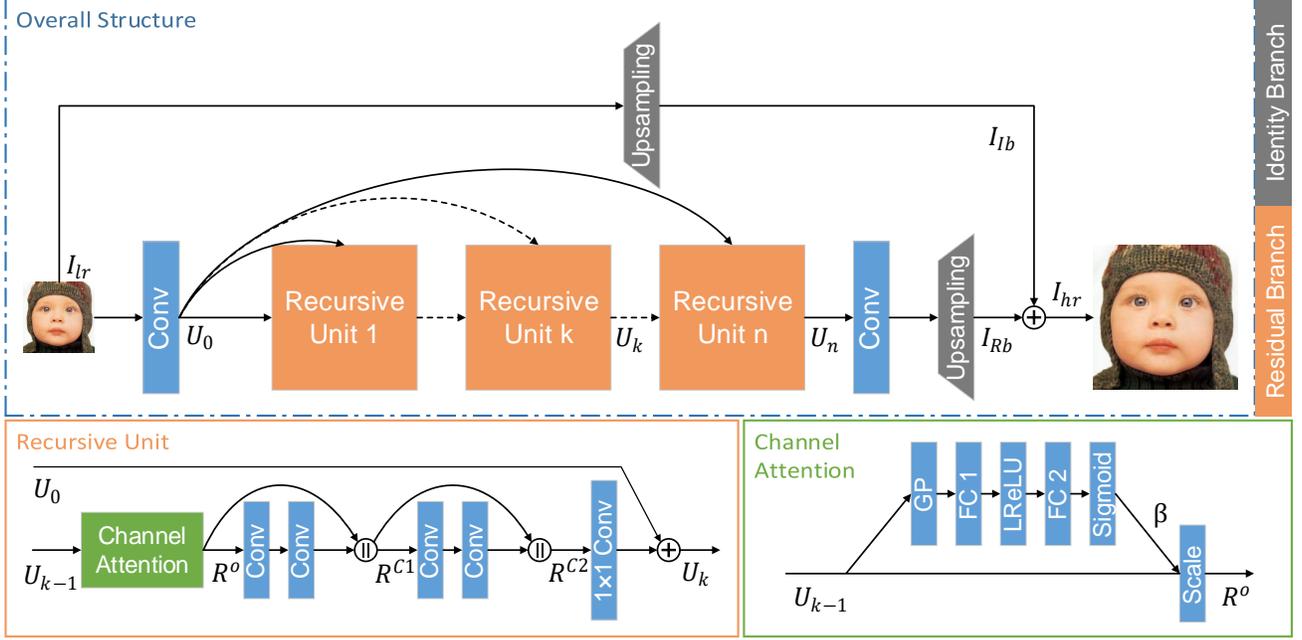
Fig. 1. The architecture of our proposed network (for 2× SR), recursive unit and channel attention. The symbol ∥ indicates concatenation and $GP$ denotes the global average pooling operation. Each convolution layer has a LReLU activation before it, and for concision, we omit it.

- *Channel Attention:* At the beginning of the recursive unit, we utilize a channel attention mechanism to adaptively recalibrate the channel importance of input features.
- *Multi-level Features Fusion:* We extract the deep-level features and aggregate multi-level features simultaneously. Moreover, the LR feature map, extracted from first Conv layer, is added at the end of each recursive unit.
- *Focusing more on details learning:* We use transposed convolution to upscale the residual branch and bicubic interpolation to upscale the identity branch. Experiments show that using a fixed method on one branch and modify the other branch is better than mutual adjustment.

## II. PROPOSED METHOD

In this section, we will formulate the proposed method, including the recursive unit and overall structure.

### A. Recursive Unit

The recursive unit is built upon the channel attention and multi-level features fusion. Let $\mathbf{U}_{k-1}$, $\mathbf{U}_k \in \mathbb{R}^{H \times W \times C}$ denote the input and output of $k$-th recursive unit, where $C$ represents the number of feature map channels.

*Channel Attention:* Our goal is to obtain a discriminative input feature of a certain recursion. Inspired by [10], a lightweight channel attention mechanism is introduced, which allows for global information to selectively emphasise informative features and restrain less useful ones via a one-dimensional vector $\boldsymbol{\beta} = [\beta_1, ..., \beta_i, ...\beta_C]$. Each scalar $\beta_i$ represents the calibration factor of $i$-th channel. As illustrated in Fig. 1, we first adopt a global average pooling across

spatial dimensions $H \times W$ to extract the global information $\boldsymbol{\alpha} \in \mathbb{R}^{1 \times 1 \times C}$ from $\mathbf{U}_{k-1}$. Then, it is followed by a dimension reduction layer with reduction ratio 4, a LReLu activation, a dimension increase layer and a sigmoid activation to generate $\boldsymbol{\beta}$. The two computable layers are implemented by fully connected layers (i.e., $1 \times 1$ convolution layers). The final output of the recalibration (denoted as $\mathbf{R}^o \in \mathbb{R}^{H \times W \times C}$) is acquired by rescaling the input features $\mathbf{U}_{k-1}$ with $\boldsymbol{\beta}$:

$$\mathbf{R}^o = \boldsymbol{\beta} \odot \mathbf{U}_{k-1}, \tag{1}$$

where $\odot$ refers to channel-wise multiplication between the calibration factor $\beta_i$ and the feature channel $\mathbf{u}_{k-1,i} \in \mathbb{R}^{H \times W}$, $i = 1, 2, ...C$.

*Multi-level Features Fusion:* In order to make full use of hierarchical features and further improve the information flow between layers, unlike LapSRN [5] and DRRN [6], we fuse multi-level features in a recursive unit rather than only using deep-level features. As shown in Fig. 1, we denote $\mathbf{R}^{c1} \in \mathbb{R}^{H \times W \times 2C}$, $\mathbf{R}^{c2} \in \mathbb{R}^{H \times W \times 4C}$ as the outputs of the two concatenation operations:

$$\mathbf{R}^{c1} = \mathbf{R}^o \parallel H_1(\mathbf{R}^o), \tag{2}$$

$$\mathbf{R}^{c2} = \mathbf{R}^{c1} \parallel H_2(\mathbf{R}^{c1}) = \mathbf{R}^o \parallel H_1(\mathbf{R}^o) \parallel H_2(\mathbf{R}^{c1}), \tag{3}$$

where the symbol ∥ denotes concatenation and $H_1$, $H_2$ represent the first two and second two convolution operations respectively. In Eq. (3), $\mathbf{R}^o$, $H_1(\mathbf{R}^o)$ and $H_2(\mathbf{R}^{c1})$ can be considered as original, shallow-level and deep-level features which are concatenated in the channel dimension. Next, we utilize a $1 \times 1$ convolution layer to compress this dimension
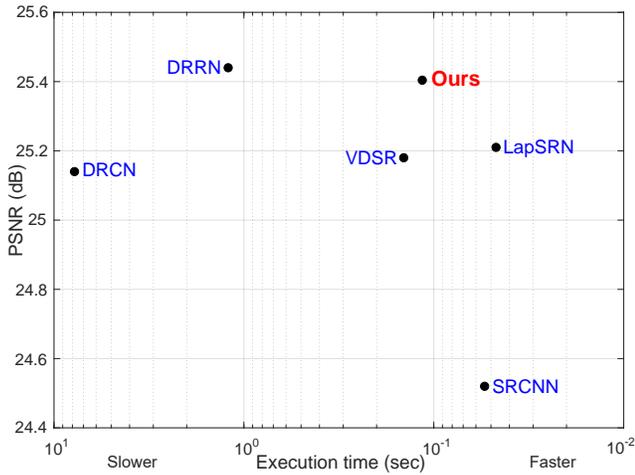
Fig. 2. Runtime and performance trade-off. The results are evaluated on Urban100 [11] with the scale factor $4\times$.



Fig. 3. Reconstruction details for $4\times$ SR. (a) Ground Truth. (b) Ours (using two learnable branches). (c) Ours (using one learnable branch)

to $C$. To solve the gradient vanishing problem, the compressed feature map is finally added to the LR feature map $\mathbf{U}_0$:

$$\mathbf{U}_k = H_3(\mathbf{R}^{c2}) + \mathbf{U}_0, \tag{4}$$

where $H_3$ refers to the function of $1 \times 1$ convolution layer.

### B. Overall Structure

We adopt the pyramid structure to recover HR images, which is introduced by [5]. Fig. 1 only illustrates the $2\times$ SR model which consists of the residual branch and the identity branch. On the residual branch, we use one convolution layer with LReLU to extract features $\mathbf{U}_0$ directly from the LR input. Then, several recursive units are stacked. Supposing there are n recursive units, the output $\mathbf{U}_n$ can be obtained by

$$\mathbf{U}_n = F_n(...(F_2(F_1(\mathbf{U}_0)))...), \tag{5}$$

where $F_n$ denotes the function of the $n$-th recursive unit. To avoid artifacts, we adopt a transposed convolution layer to up-sample the global residual image (denoted as $\mathbf{I}_{Rb}$). On the identity branch, the LR input is up-sampled by the bicubic interpolation instead of a learnable method. In contrast to LapSRN that performs transposed convolution on the identity branch, our method stabilizes the training process and focuses more on image details. The up-sampled image ($\mathbf{I}_{Ib}$) is finally combined with $\mathbf{I}_{Rb}$ to estimate the HR image ($\mathbf{I}_{hr}$) via an element-wise summation, which can be formulated as:

$$\mathbf{I}_{hr} = F_{up1}(\mathbf{U}_n) + F_{up2}(\mathbf{I}_{lr}) = \mathbf{I}_{Rb} + \mathbf{I}_{Ib}, \tag{6}$$

where $F_{up1}$ and $F_{up2}$ refer to the upsampling operations on the residual branch and identity branch respectively.

## III. EXPERIMENTS

### A. Implementation Details

We train all of the models mentioned later with 291 images, where 91 images are from Yang et al. (T91) [12] and other 200 images are from Berkeley Segmentation Dataset (BSDS200)

[13]. By following [5], we augment the training dataset via randomly scaling, rotating and flipping. For testing, we use three widely used benchmark datasets: Set5 [14], Set14 [15] and BSDS100 [13]. Our model is evaluated with PSNR and SSIM. We convert all images into YCbCr color space and only use the Y-channel to process.

All recursive units share same parameters, and we use 64 convolutional filters for the first two layers and 128 filters for the second two layers. All convolutional filters (except the $1 \times 1$ convolution layer) have the same kernel size ($3 \times 3$) and are initialized by the method of He et al. [16]. The learning rate is set to $10^{-5}$ and decreased by a factor of 2 for every 80 epochs. We set patch size to $128 \times 128$ and batch size to 64. Our method is implemented with MatConvNet toolbox [17] and a NVIDIA GTX 1080ti GPU.

### B. Comparison with the State-of-the-Art

Taking the trade-off between runtime and performance into account, we use 6 recursive units in our method. We compare the proposed method with 6 SR methods: Bicubic, SRCNN [2], VDSR [3], DRCN [4], LapSRN [5] and DRRN [6]. The quantitative results of exiting methods are quoted from their papers and shown in Table I. Our method significantly outperforms the prior methods (except DRRN) in all testing datasets and scale factors. Compared with DRRN, our method performs slightly worse on PSNR while better on SSIM in most instances, noting that SSIM focuses on measuring structural and detail similarities. Since both LapSRN and our method do not train $3\times$ SR model exclusively, we achieve relatively poor performance of $3\times$ SR.

As for execution time, we use the original codes of the compared methods to evaluate the runtime on a same machine with 3.6GHz Intel i7 CPU (32G RAM) and NVIDIA GTX 1080ti GPU (11G Memory). Fig. 2 shows the trade-off between the execution time and performance on Urban100 [11] dataset for $4\times$ SR. Our method is 0.04dB lower than DRRN on PSNR, but approximately 10 times faster.

TABLE I
BENCHMARK RESULTS. AVERAGE PSNR/SSIM FOR SCALE FACTOR ×2, ×3 AND ×4. RED COLOR INDICATES THE BEST PERFORMANCE AND BLUE COLOR INDICATES THE SECOND BEST PERFORMANCE.

| Dataset | Scale | Bicubic | SRCNN [2] | VDSR [3] | DRCN [4] | LapSRN [5] | DRRN [6] | Ours |
|---|---|---|---|---|---|---|---|---|
| Set5 | ×2 | 33.66/0.9299 | 36.66/0.9542 | 37.53/0.9587 | 37.63/0.9588 | 37.52/0.9591 | 37.74/0.9591 | 37.69/0.9598 |
| | ×3 | 30.39/0.8682 | 32.75/0.9090 | 33.66/0.9213 | 33.82/0.9226 | 33.82/0.9227 | 34.03/0.9244 | 33.94/0.9233 |
| | ×4 | 28.42/0.8104 | 30.48/0.8628 | 31.35/0.8838 | 31.53/0.8854 | 31.54/0.8855 | 31.68/0.8888 | 31.67/0.8888 |
| Set14 | ×2 | 30.24/0.8688 | 32.45/0.9067 | 33.03/0.9124 | 33.04/0.9118 | 33.08/0.9130 | 33.23/0.9136 | 33.20/0.9139 |
| | ×3 | 27.55/0.7742 | 29.30/0.8215 | 29.77/0.8314 | 29.76/0.8311 | 29.79/0.8320 | 29.96/0.8349 | 29.89/0.8340 |
| | ×4 | 26.00/0.7027 | 27.50/0.7513 | 28.01/0.7674 | 28.02/0.7670 | 28.19/0.7720 | 28.21/0.7720 | 28.21/0.7726 |
| BSDS100 | ×2 | 29.56/0.8431 | 31.36/0.8879 | 31.90/0.8960 | 31.85/0.8942 | 31.80/0.8950 | 32.05/0.8973 | 32.00/0.8971 |
| | ×3 | 27.21/0.7385 | 28.41/0.7863 | 28.82/0.7976 | 28.80/0.7963 | 28.82/0.7973 | 28.95/0.8004 | 28.88/0.8000 |
| | ×4 | 25.96/0.6675 | 26.90/0.7101 | 27.29/0.7251 | 27.23/0.7233 | 27.32/0.7280 | 27.38/0.7284 | 27.37/0.7294 |

TABLE II
ABLATION EXPERIMENTS OF OUR MODEL ON SET5 FOR 4× SR. REMOVING EACH COMPONENT WILL DEGRADE PERFORMANCE.

| Fusion | Attention | Learnable Branch | Set5 |
|---|---|---|---|
| √ | | One | 31.62 |
| | √ | One | 31.58 |
| √ | √ | Two | 31.66 |
| √ | √ | One | **31.67** |

## C. Network Analysis

To demonstrate the effect of each component, we carry out four ablation experiments of channel attention, multi-level features fusion and the number of learnable branches (i.e., using bicubic interpolation or transposed convolution in the identity branch). By removing the multi-level features fusion, our model falls back to a network similar to LapSRN [5] but with the channel attention. The results confirm that making full use of multi-level features will significantly improve performance. One possible reason is that fusing hierarchical features improves the information flow and eases the difficulty of training. We can conclude from Table II that the model with full components achieves the best performance. Fig. 3 shows the different reconstruction details between the two learnable branches and one learnable branch. Although using one learnable branch improves only 0.01dB on PSNR, it gets much better visual effects in image details.

## IV. CONCLUSION

In this paper, a novel recursive unit is proposed for boosting single image super-resolution. The proposed unit first performs channel recalibration on input features, and then the multi-level features are extracted and fused. Moreover, we use transposed convolution on the residual branch and bicubic interpolation on the identity branch, which reconstructs more details in terms of visual perception. Experiments show that our network achieves competitive reconstruction performance and maintains faster execution speed.

## ACKNOWLEDGMENT

## REFERENCES

[1] D. Glasner, S. Bagon, and M. Irani, "Super-resolution from a single image," in *Computer Vision, 2009 IEEE 12th International Conference on*. IEEE, 2009, pp. 349–356.

[2] 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.

[3] 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*, 2016, pp. 1646–1654.

[4] ——, "Deeply-recursive convolutional network for image super-resolution," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1637–1645.

[5] W.-S. Lai, J.-B. Huang, N. Ahuja, and M.-H. Yang, "Deep laplacian pyramid networks for fast and accurate super-resolution," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 624–632.

[6] Y. Tai, J. Yang, and X. Liu, "Image super-resolution via deep recursive residual network," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, vol. 1, no. 4, 2017.

[7] Q. Wang, Z. Teng, J. Xing, J. Gao, W. Hu, and S. Maybank, "Learning attentions: residual attentional siamese network for high performance online visual tracking," in *CVPR*. IEEE Computer Society, 2018.

[8] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *CVPR*, vol. 1, no. 2, 2017, p. 4.

[9] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 4, pp. 834–848, 2018.

[10] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," *arXiv preprint arXiv:1709.01507*, 2017.

[11] J.-B. Huang, A. Singh, and N. Ahuja, "Single image super-resolution from transformed self-exemplars," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 5197–5206.

[12] J. Yang, J. Wright, T. S. Huang, and Y. Ma, "Image super-resolution via sparse representation," *IEEE transactions on image processing*, vol. 19, no. 11, pp. 2861–2873, 2010.

[13] P. Arbelaez, M. Maire, C. Fowlkes, and J. Malik, "Contour detection and hierarchical image segmentation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 33, no. 5, pp. 898–916, 2011.

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

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

[16] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1026–1034.

[17] A. Vedaldi and K. Lenc, "Matconvnet: Convolutional neural networks for matlab," in *Proceedings of the 23rd ACM international conference on Multimedia*. ACM, 2015, pp. 689–692.