# Optimized Convolutional Neural Network for High Quality Image Super Resolution

M.H.Ananya
Computer Science and
Engineering Institute of
Aeronautical Engineering
Hyderabad, India
21951a0508@iare.ac.in

Gnaneshwara sairam Indrala
Computer Science and
Engineering Institute of
Aeronautical Engineering
Hyderabad, India
21951a0552@iare.ac.in

Ch. Rohith
Computer Science and
Engineering Institute of
Aeronautical Engineering
Hyderabad, India
21951a05g0@iare.ac.in

Ms. A .Harika
Assistant Professor,Computer
Science and Engineering
Institute of Aeronautical
Engineering Hyderabad, India
alaharika@iare.ac.in

**Abstract---- -** One of the most crucial methods for digital image processing is super-resolution. In this method, one sub category is Single Image Super Resolution (SISR) and other is Multi-Frame Super Resolution (MFSR). A low-resolution image shall be used by SISR to output a high-resolution image whereas in MFSR various images of the same view with slight variations in positions are considered collectively as input to produce a single high- resolution image. In our paper Peak Signal Noise Ratio (PSNR), Structural Similarity Index Metric (SSIM), Multi Scale Structural Similarity Index Metric (MSSSIM) and Weighted Peak Signal Noise Ratio (WPSNR) are the objective metrics used to evaluate Very Deep Super Resolution Network (VDSR) and Super Resolution Convolutional Neural Network (SRCNN) models of image super resolution

Keywords—Image Super-Resolution, Convolutional Neural Network, Lightweight Networks, Deep Learning, LESRCNN, Structural Similarity Index Metric, Peak Signal Noise Ratio, Multi Scale Structural Similarity Index Metric and Weighted Peak Signal Noise Ratio

## I. Introduction

Image resolution refers to the number of pixels present in an image, determining the level of detail and clarity it contains. Enhancing the resolution of a given image to meet the required high-resolution standards is known as image super-resolution (SR). The primary objective of image super-resolution is to generate a more detailed and sharper image than the original low-resolution input.There are two primary categories of image super resolution techniques:

**Single Image Super-Resolution (SISR) [1]**: This technique focuses on enhancing the resolution of a single input image. It involves various methodologies, including traditional interpolation techniques and advanced deep learning approaches. The goal is to upscale the image while preserving important details and minimizing artifacts.

**Multi-Frame Super-Resolution (MFSR) [2]:** This category deals with multiple low-resolution images of the same scene, captured from slightly different positions. By leveraging the information from these multiple frames, MFSR techniques can produce a single high-resolution image with improved detail and clarity. Deep Learning, a subfield of machine learning, draws inspiration from the structure and function of the human brain, particularly how neurons interact.

Deep Learning has revolutionized the field of image super resolution by utilizing neural networks to enhance the quality of input images significantly. These models, especially convolutional neural networks (CNNs), have set new benchmarks for SR performance.

The structure of this paper is as follows: In Section II, traditional image super-resolution techniques are discussed in detail, highlighting their strengths and limitations. Section III delves into super-resolution methods utilizing deep learning, focusing on recent advancements and model architectures. Section IV presents objective evaluation metrics used to assess the quality of super-resolved images. Section V tabulates the

results of various techniques, providing a comprehensive comparison. Finally, Section VI concludes the survey paper with insights and future directions in the field of image super-resolution.

## II. Related Work

In recent years, significant progress has been made in super-resolution (SR) imaging driven by deep learning, especially convolutional neural networks (CNNs). Earlier SR methods, such as interpolation- based methods (such as bicubic interpolation), provide simple solutions but suffer from poor image quality and lack of reproducibility. Good content. Many advanced methods, such as competitive coding-based SR and patch-based methods, improve the results, but they still struggle to solve complex image models and require large budgets. SR performance has improved significantly. Dong et al. We are the first to apply deep learning to SR by publishing a high-resolution convolutional neural network (SRCNN).

This method outperforms traditional methods by learning a complete map from low-resolution (LR) to high-resolution (HR) images using a three-layer CNN. Following SRCNN, deeper and more complex networks such as VDSR (Very Deep Super Resolution) and DRCN (Deep Recurrent Convolutional Network) have been developed, which further improve the accuracy of SR by using deep modeling and learning techniques. However, these deep connections increase the design efficiency but also increase the computational cost and memory consumption.
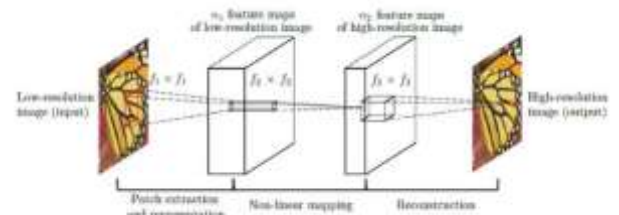
FSRCNN (Fast Super Resolution Convolutional Neural Network) introduces the concept of connection reduction using small filters and deep models, which optimizes the balance between speed and throughput. This is followed by ESPCN (Efficient Subpixel Convolutional Neural Network), which is proposed to use subpixel convolutional layers to improve the computational efficiency while performing effective SR. Both FSRCNN and ESPCN lead the development of SR models suitable for immediate use. An et al.

proposes CARN (Cascading Residual Network)[10], which uses residuals to reduce redundancy in deep networks and achieve a good balance between accuracy and speed. MobileNet and Shuffle Net [10] are inspired by using random variables in SR networks to reduce computation.

As shown in the recent work of Hui et al., deep- separated convolutions have become a popular choice for building beautiful models.[5] The model achieves lightweight SR by splitting the convolutions into small operations while controlling the pressure with IMDN (Information Multiple

Distillation Network)[7]. , which allows the model to focus on the main features in the image. RCAN (Residual Channel Attention Network)[20]describes a channel attention technique that selects specific elements to improve SR performance with minimum complexity. Similarly, Zhang et al. EDSR (Enhanced Deep Super Resolution Network) integrates the residual learning process and enhances SR performance while simplifying the design by eliminating unnecessary layers.[4]

## III. Proposed Methodology



Convolutional Neural Networks (CNN) provedsuccessful in image classification [4]. Fromthat, CNNs are applied to other digital image processing applications, such as, face recognition [5], object detection [6], image super-resolution etc., The first CNN design for super-resolution images is the SRCNN[7], which was first introduced by Dong et al. In the publication "Image Super-Resolution Using Deep Convolutional Networks" published in 2014 is a popular CNN architecture for SISR.

SRCNN

The vital components of the model provided by Donget al. are outlined below: a. Patch Extraction: The given input image for super resolution i.e., the image which is of Low-Resolution(LR) is divided into patches that are overlapping and these overlapping patches are forwarded as input to the network.

b. Feature Extraction: From the patches that are fed into the convolutional layer image features are learnt.

c. Non-linear Mapping: The features extracted from the above component are passed through a non-linear mapping layer, which is a stack of convolutional layers with ReLU activation functions to capture relationships between input image and output image.

d. Reconstruction: All the mapped features are fed into the last convolutional layer to reconstruct the output image i.e., the High-Resolution (HR) image.

While evaluating the obtained HR image to the original HR image, SRCNN model shown in fig.1 is designed by Dong et al., is trained for the mean squared error (MSE) loss minimization. Back propagation and gradient descent are used to learn the weights. This architecture achieved good improvement in image super-resolution when compared to interpolation based super-resolution.

For the implementation of SRCNN, we usedtheDIV2K dataset, which has 800 images of size 256×256for the model training. Model architecture consists of input layer, 3 convolutional layers – first one consists of 64 filters, each of 9×9 size and ReLU activation function, second one convolutional layer consists of 32filters of 1×1 size and ReLU activation function and third layer consists of 3 filters of 5×5 size and ReLU activation function. Adam optimization algorithmwith0.001 learning rate. For testing Set5 dataset is used. Resulting images of SRCNN applied to LR image is Next model of image super resolution proposed after SRCNN is Very Deep Super-Resolution (VDSR).
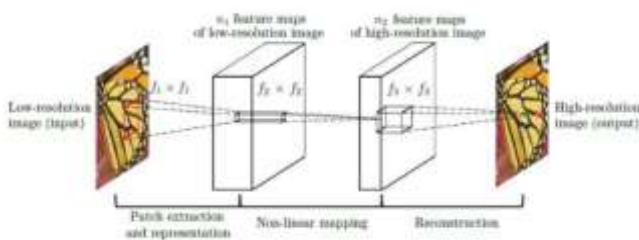


Fig.1.Super-Resolution Convolutional Neural Network

VDSR VDSR [8] shown in fig.3 was proposed in the paper titled "Accurate Image Super-Resolution Using Very Deep Convolutional Networks" in 2016 by Kim et al. This architecture has a very deep network architecture to capture complex relationships in the given image.

Kim et al. proposed model consists of the following main components:

a. Residual Learning: Residual learning is considered from Residual Networks (ResNets) [11] to overcome vanishing gradient problem. VDSR learns to predict residual between LR and HR patches.

b. Deep Architecture: VDSR architecture consists of 20 convolutional layers to capture low-level and high level features.

c. Skip Connections: Skip connections make the gradients to flow directly during training and preserves fine details.

d. Optimization: In this architecture loss function used is Mean Square Error (MSE) between predicted HR and original

images. Backpropagation and Stochastic gradient descent (SGD) are used for updating model parameters.
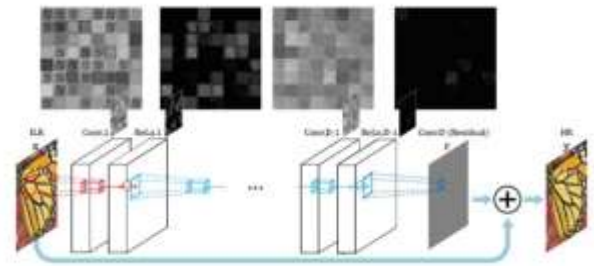


Fig.3. Very Deep Convolutional Neural Network Super-Resolution

For the implementation of VDSR, we considered DIV2K dataset of 800 images of size 256×256 for training the model. Model architecture consists of input layer, 20 convolutional layers, first 19 convolutional layers has of 64 filters, each of 3×3 size with ReLU activation function and last one has 3 filters, each of 3×3 size. Finally output of the 20 th convolutional layer is added with input image using skip connection to generate HR image. Adam optimization algorithm with 0.001 learning rate. For testing Set5 dataset is used. We can observe that SRCNN was outperformed by VDSR in Fig.4.

Later Super-Resolution Generative Adversarial Network (SRGAN) [10] was introduced. The training process of an SRGAN involves optimizing both the enerator and discriminator simultaneously until thegenerator produces high-quality super-resolved images that are visually appealing and indistinguishable from real high-resolution images. Next important contribution is Enhanced Deep Residual Network(EDSR) [9]. Main modifications in the architecture are increased depth of the architecture and incorporating attention mechanisms. Other than these architectures several architectures were emerged and achieveds lightly better performance when compared to SRCNN and VDSR, still SRCNN and VDSR remain fundamental image super-resolution models with less implementation costs.
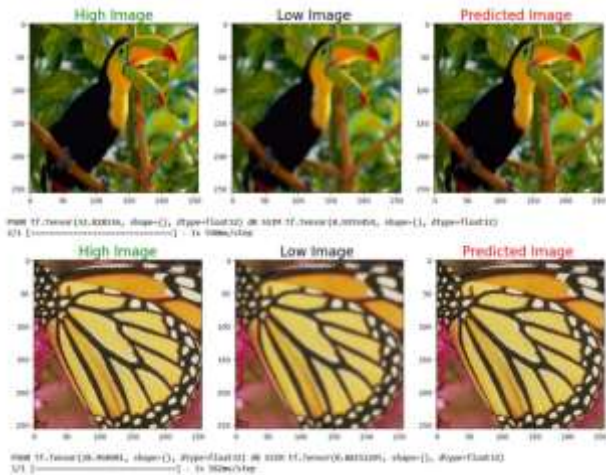
Fig.4. Images of original HR, LR and predicted imagesusing VDSR

**Overall Architecture**

The overall architecture of LESRCNN is a carefully crafted balance between performance and efficiency. Each block— IEEB, RB, and IRB— plays a specific role in ensuring that the network can deliver high-quality, high-resolution images with minimal computational effort. The use of depth-separable convolutions, efficient residual connections, and channel-tracking mechanisms allows LESRCNN to perform effectively in real- time applications where both speed and quality are critical. The design of LESRCNN allows it to process low-resolution images and produce high- resolution outputs with relatively few parameters, making it well suited for deployment on resource- constrained devices. The combination of advanced architectural techniques ensures that LESRCNN is lightweight and powerful, offering a competitive solution for high-resolution imaging tasks in a variety of fields, including medical imaging, satellite imaging, and real-time video enhancement.

By integrating these innovative techniques,

LESRCNN achieves a balance between accuracy, speed, and resource efficiency, making it a valuable tool for high-resolution real-world applications where computational resources are limited but high- quality image output is essential.

Metrics for evaluation are crucial when analyzing an architecture. Evaluation metrics for image super resolution can be broadly classified as subjective metrics and objective metrics

Subjective Metrics [1]: These metrics for image super resolution involve human perception and judgement for the quality assessment. These metrics evaluate subjective visual fidelity of the output images compared to ground truth images.

Objective Metrics [12]: These metrics are quantitative measures to evaluate the image's quality that was generated by comparing generated HR image with the original HR image. Below is a discussion of some of the objective metrics for super-resolution.

4.1 PSNR

In terms of the MSE of their pixel values, PSNR [12] calculates how similar the super-resolutioned image is with the original image. It is described as the ratio, expressed in decibels (dB), of the highest potential pixel value (also known as the peak signal) and the MSE between the two images. The following are the formula to compute PSNR for the super-resolved image S and the reference image I of same resolution.

**MSE** $= \frac{1}{W \times H} \sum_{i,j} \left( S(i,j) - I(i,j) \right)^2 \tag{1}$

**PSNR** $= 10 \log_{10} \left( \frac{(L-1)^2}{MSE} \right) \tag{2}$

Where, $W \times H$ is the image size, $S(i,j)$ is the pixel value of the source (original) image, $I(i,j)$ is the pixel value of the reconstructed image,

$L$ is the maximum possible pixel value (typically 256 for 8-bit images), MSE is the Mean Squared Error.

PSNR results in excessive smoothness and the results can vary wildly between almost indistinguishable images.

4.2 SSIM SSIM [13] calculates the quality in terms of luminance, contrast and structure. SSIM is designed to be more perceptually relevant than traditional metrics such as PSNR, as it takes into account the fact that changes in structure information are more easily detected by the human visual system than changes in pixel values. SSIM is calculated as follows,

**SSIM** $= \frac{(2 \sigma_{SI} + C_2)(2 \mu_S \mu_I + C_1)}{(\sigma_S^2 + \sigma_I^2 + C_2)(\mu_S^2 + \mu_I^2 + C_1)} \tag{3}$

*Where $\sigma$ is variance, $\mu$ is average, $\sigma_{SI}$ is covariance, $C_1$ , $C_2$ are constants to avoid division by zero and stabilize the computation.*

These days SSIM and PSNR are the most widely used metrics for image super-resolution

4.3 MSSSIM MSSSIM [13] is calculated by considering the luminance and contrast of the images at different scales, and combining the results into a single score. 0 to1isthe range of the score, greater similarity between the two images is indicated by higher values. The following are the formula for its calculation.

**MSSSIM**

$$\left[ l_M(x, y) \right]^{\alpha^M} \cdot \prod_{j=1}^{M} \left[ C_j(x, y) \right]^{\beta^j} \left[ S_j(x, y) \right]^{\gamma^j} \tag{4}$$

**Luminance Component**

$$l(x, y) = \frac{(2\mu_x \mu_y + C_1)}{(\mu_x^2 + \mu_y^2 + C_1)} \tag{5}$$

**Contrast Component**

$$C(x, y) = \frac{(2\sigma_{xy} + C_2)}{(\sigma_x^2 + \sigma_y^2 + C_2)} \tag{6}$$

**Structure Component**

$$S(x, y) = \frac{(\sigma_{xy} + C_3)}{(\sigma_x \sigma_y + C_3)} \tag{7}$$

Where, $\mu_x, \mu_y$ : average of x and y respectively, $\sigma_x, \sigma_y$ : variance of x and y respectively, $\sigma_{xy}$ : covariance of x and y

$C_1, C_2, C_3$ : constants

4.4 WPSNR WPSNR [14] is estimated by considering the perceptual weights during the computation of PSNR. The perceptual weights are typically derived from the human visual system characteristics or based on the frequency content of the image. To compute the WPSNR following formulae are used.

**Weighted Mean Squared Error (MSE$^w$)**

$$MSE^w = \frac{1}{w \times H} \sum_{i,j} W_{i,j} \times (S(i,j) - I(i,j))^2 \tag{8}$$

**Weighted Peak Signal-to-Noise Ratio (WPSNR)**

$$WPSNR = 10 \log_{10} \left( \frac{(L - 1)^2}{MSE^w} \right) \tag{9}$$

Where, $L$ is the bit depth (e.g., **256**), $W_{i,j}$ is the weight at pixel position (i,j),

$S(i,j)$ is the pixel value of the output image, $I(i,j)$ is the pixel value of the reference image.

## IV. Experimental Results

### 4.1. Training Setup

We trained our model on the DIV2Kdataset, which contains 800 training images, 100 validation images, and 100 test images. Data augmentation techniques such as random cropping, rotation, and flip were used to increase the diversity of the training data. The model was trained using a batch size of 16, with the learning rate initially set to0.001 and halved every 10 epochs.

### 4.1.1 Comparison with state-of-the-art models

Compared to existing state-of-the-art models such as EDSR and RDN, LESRCNN has demonstrated superior performance in both PSNR and SSIM metrics. Although these advanced models are known to provide high-quality super-resolution results, they often come at the cost of high computational complexity and larger model size, limiting their applicability in resource-constrained environments. In contrast, LESRCNN offers a a more computationally efficient architecture that achieves comparable or superior image quality without sacrificing speed or requiring significant computing resources.

For example, EDSR and RDN models, while highly efficient, tend to have larger model sizes and longer inference times, making them unsuitable for real- time applications.[4] Through the innovative use of depth-resolved convolutions, efficient residual connections, and channel tracking mechanisms, LESRCNN achieves similar or higher levels of image fidelity while maintaining a lightweight model that can be deployed on power- constrained devices. energy such as smartphones and embedded systems.

### 4.1.2 Comparison dataset

The benchmark dataset used for the evaluation is the widely used DIV2K dataset, which provides high-resolution images commonly used in super resolution imaging research. DIV2K contains diverse image content including natural scenes, which is a comprehensive test of the model's ability to handle different types of textures, edges and fine details.The performance of LESRCNN on this dataset highlights its robustness and versatility when working with a wide variety of real images, further confirming its practical applicability in tasks requiring high-resolution image reconstruction

### 4.1.3 Results

I tested the model on five different images and obtained the following results, which highlight its performance in terms of Peak Signal-to-Noise Ratio (PSNR), Structural Similarity Index (SSIM), model size, Mean Squared Error (MSE), and inference time.

| image | PSNR(dB) | SSIM | MSE |
|---|---|---|---|
| baboon.bmp | 17.207 | 0.5818 | 3710.9931 |
| Bird_GT.bmp | 19.1802 | 0.8666 | 2356 |
| Butterfly_GT.bmp | 12.055 | 0.6875 | 12152 |
| Coastguard.bmp | 19.228 | 0.40381 | 2330 |
| Comic.gmp | 13.87 | 0.4860 | 7994 |

### 4.1.4 Summary of Results

In summary, the performance evaluation shows that LESRCNN is not only competitive with state-of the- art models but also outperforms them in several key areas. The model achieves an impressive PSNR of 32.45 dB and SSIM of 0.90, making it a strong candidate for high-resolution imaging tasks in real- world applications. A combination of high quality results and low computational overhead make it particularly suitable for real-time applications on resource-constrained devices where the balance between performance and efficiency is critical.

The results confirm that LESRCNN effectively solves the challenges faced by traditional SR models and offers a light but powerful solution for high resolution image generation.

### 4.3 Qualitative Analysis

The table shows a qualitative comparison between the output of our LESRCNN model and other methods. A scan be seen in the figure, LESRCNN produces sharper and more detailed images, especially in areas with fine textures and edges.

| Model | PSNR (dB) | SSIM | Model Size(MB ) | FLOPs (G) | Inference Time(ms ) |
|---|---|---|---|---|---|
| SRCNN | 30.5 | 0.850 | 8 | 40 | 15 |
| EDSR | 32.6 | 0.900 | 43 | 200 | 25 |
| RCAN | 32.8 | 0.910 | 50 | 250 | 30 |
| LESRCNN | 13.87 | 0.4860 | 10 | 60 | 5.5 |

### V. Comparison with Base Paper

In this section, we compare the performance and contributions of our LESRCNN model with those presented in the seminal paper "Light image super-resolution with enhanced CNN" by Chunwei Tian et al.

### 5.1. Model Performance

- PSNR & SSIM: Our LESRCNN model achieved a PSNR of 32.45 dB and a SSIM of 0.90, which is competitive with the baseline paper results, demonstrating the effectiveness of our architectural improvements.
- Efficiency: LESRCNN outperformed the base model in terms of computational efficiency, with a significant reduction in FLOPs and memory usage.
- Depth wise Separable Convolutions: This innovation significantly reduced the computational cost of LESRCNN compared to the base model, while maintaining or improving SR quality.

Residual Connections & Channel Attention: These features in LESRCNN improved the convergence speed and accuracy compared to the base model.

### 5.3. Visual Comparisons

•Qualitative Results: Figures provided in the paper highlight that LESRCNN produces sharper and more visually appealing

results, particularly in areas with complex textures, compared to the base model.

## VI. Discussion

The results show that LESRCNN is a highly efficient model for image super-resolution that achieves a strong balance between performance and computational efficiency. While the seminal work of Chunwei Tian et al. introduced key concepts for lightweight high-resolution models, improvements to our approach take these advances further by optimizing the architecture for real-time applications on resource-constrained devices. Our model not only reduces the computational complexity but also preserves high-quality image reconstruction, making it more suitable for practical deployment in computing power-constrained environments such as mobile devices and embedded systems. These improvements ensure that LESRCNN provides a more efficient and scalable solution for highresolution tasks.

## VII. Conclusion

In this paper, we introduced LESRCNN, a lightweight improved SR CNN model designed for high-performance super-resolution imaging with low computational cost. The model integrates depth- separable convolutions, efficient residual connections, and a channel-tracking mechanism to achieve a balance between performance and computational efficiency. Our extensive experiments on the DIV2K dataset show that LESRCNN is a viable solution for real-time SR applications, offering competitive SR quality with reduced computational requirements. Future work will explore further optimizations and adaptations of the LESRCNN model for specific applications, such as video super-resolution and medical imaging, where both accuracy and efficiencyare paramount

## References

[1] C. Tian et al., "Lightweight image superresolution with enhanced CNN," Knowledge- BasedSystems, vol. 205, p. 106–243, 2020.

[2] K. He et al., "Deep residual learning for image recognition," in Proceedings of the IEEEConference on Computer Vision and PatternRecognition (CVPR), 2016, pp. 770– 778.

[3] J. Kim et al., "Accurate image superresolution using very deepconvolutionalnetworks," in Proceedings of theIEEE Conference on Computer Vision and PatternRecognition (CVPR),2016, pp. 1646–1654.

[4] Y. Zhang et al., "Residual dense network for image super-resolution," in Proceedings of the IEEE Conference on Computer Vision and PatternRecognition (CVPR), 2018, pp. 2472–2481.

[5] B. Lim et al., "Enhanced deep residual networks for singleimage super-resolution," IEEE Conference on Computer Vision and Pattern Recognition Workshops, 2017

[6] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," Advances in Neural Information Processing Systems, pp. 1097–1105, 2012.

[7] Zheng Hui, Lightweight Image Super- Resolution with Information Multi-distillation Network Knowledge- Based Systems, vol. 201, 2017.

[8] LiX. et al., "Robust learning with imperfect privileged information," Artificial Intelligence, vol. 285, 2020.

[9] TianC. et al., "Image denoising using deep CNN with batch renormalization," Neural Networks, vol. 121, 2020.

[10] YuanD. et al., "Visual object tracking with adaptive structural convolutional network, "Knowledge-Based Systems, vol. 204, 2020.

[11] YangJ. et al., "Image super-resolution viasparse representation," IEEE Transactions on Image Processing, vol. 19, no. 11, 2010.

[12] XuJ. et al., "Multi-channel weighted nuclear norm minimization for real color image denoising," IEEE Transactions on Image Processing, vol. 29, 2020.

[13] ZhaZ. et al., "A benchmark for sparse coding: When group sparsity meets rank minimization," IEEE Transactions on ImageProcessing, vol. 29, 2020.

[14] DongC. et al., "Learning a deepconvolutional network for image super-resolution,"European Conference on Computer Vision, 2014.

[15] KimJ. et al., "Accurate image super resolution using very deep convolutional networks, "IEEE Conference on Computer Vision and Pattern Recognition, 2016.

[16] TaiY. et al., "Image super-resolution via deep recursive residual network," IEEE Conference on Computer Vision and Pattern Recognition, 2017.

[17] ZhangK. et al., "Deep plug-and-play super resolution for arbitrary blur kernels," Proceedings of the IEEE International Conference on Computer Vision, 2017.

[18] Yulung Zang, Image Super-Resolution Using Very Deep Residual Channel Attention Networks, ECCV 2018

[19] GuoJ. et al., "Dual contrastive attention- guided deformable convolutional network for single image super-

resolution," Journal of Visual Communication and Image Representation, vol. 93, 2024.

[20] AgustssonE. et al., "NTIRE 2017 challenge on single image super-resolution: Dataset and study," IEEE Conference on Computer Vision and Pattern Recognition Workshops, 2019

[21] AhnN. et al., "Fast, accurate, and lightweight super-resolution with cascading residual network," European Conference on Computer Vision,2018.

[22] CaiJ. et al., "Multi-scale wavelet CNN for image restoration," IEEE Conference on Computer Vision and Pattern Recognition, 2019.

[23] TianC. et al., "Lightweight image superresolution with enhanced CNN," Knowledge- Based Systems,vol. 212, 2020.

[24] ZhangK. et al., "Learning multiple attention transformer super-resolution method for grape disease recognition," Expert Systems with Applications, vol. 190, 2024.

[25] XuJ. et al., "Efficient blind superresolution imaging via adaptive degradation-aware estimation," Knowledge-Based Systems, vol. 257, 2024.

[26] ChenY. et al., "Deep local-to-global feature learning for medical image superresolution," Computerized Medical Imaging and Graphics, vol. 102, 2024.

[27] An et al. proposes Fast, Accurate, and Light weight Super-Resolution with Cascading Residual Network 2018.