

Image Super-resolution using Efficient Subpixel CNN with Residual Dense Blocks

Mihika Dhariwal, Kavyasri R

Abstract: Image super-resolution is integral to various applications, such as medical imaging and remote sensing, where the enhancement of low-resolution images is crucial for accurate analysis and interpretation. Our method introduces an approach to image super-resolution, integrating an Efficient Subpixel Convolutional Neural Network (CNN) with Residual Dense Blocks for enhanced performance. In comparison to traditional methods like Single Convolutional Neural Network (SRCNN), bicubic interpolation, and the standalone ESPCN model, the incorporation of Residual Dense Blocks adds depth and improves the model's capacity to capture intricate features within image datasets. Feature maps are extracted in its Low-Resolution space as opposed to earlier techniques where this was done in the High-Resolution space. In addition to the sub-pixel convolution layer, which learns an array of filters to upscale the final LR feature maps into the HR output, the Residual Dense Blocks (RDBs) heavily focuses on keeping the information extracted in previous layers alive, while computing outputs for the present layer.

Keywords: Image super-resolution, low-resolution images, efficient subpixel CNN, residual dense blocks, single convolutional neural network (SRCNN), bicubic interpolation, ESPCN model, sub-pixel convolution layer.

1 Introduction

Images play a crucial role in conveying information across diverse sectors, including industrial production, biomedical applications, and robot vision. The intricacies of image processing, involving various devices, often lead to the acquisition of poor-quality images characterized by issues like low resolution, blur, distortion, noise, and loss of detail. Addressing these challenges, image super-resolution (SR) reconstruction technology has been introduced, utilizing multiple low-quality and low-resolution (LR) images with complementary information to generate a single high-quality and high-resolution (HR) image. This technology finds extensive applications in fields such as remote sensing [1], medical imaging [2], and monitoring [3]. Over the past decades, numerous SR methods have emerged, categorized into four groups based on interpolation, reconstruction, and machine learning or deep learning methods.

The most straightforward method for enhancing the resolution of an image involves a basic upsampling through interpolation. Interpolation achieves this by directly augmenting the image's dimensions, adding new pixels or data points to the low-resolution image. There are three widely used interpolation techniques:

(1) Nearest Neighbour interpolation, (2) Bilinear interpolation, and (3) Bicubic interpolation. [4]

In light of recent AI advancements, there is a notable surge in the exploration of deep learning-based super-resolution models. Numerous open-source models have demonstrated impressive performance in benchmark studies. Among these, the Super Resolution Convolutional Neural Network (SRCNN) stands out as one of the pioneering deep learning methods that surpassed earlier traditional super-resolution techniques. Leveraging a Convolutional Neural Network (CNN) architecture, SRCNN enables the AI model to learn a direct mapping between low-resolution and high-resolution images [5]. However, this method exhibits certain limitations. Firstly, it necessitates the use of interpolation techniques, such as bicubic interpolation, for the upsampling of low-resolution (LR) images. Secondly, it elevates the resolution either before or at the initial layer of the network. In essence, the CNN approach directly applies the convolutional neural network to the upsampled LR image, leading to heightened computational complexity and increased memory costs.

To address these issues, a novel approach called Efficient Subpixel Convolution Neural Network (ESPCN) has been introduced. This method introduces an efficient sub-pixel convolutional layer to the CNN network, aiming to mitigate the drawbacks associated with traditional CNN-based super-resolution techniques. [6]

Our approach is the ESPCN model integrated with Residual Dense Blocks (RDB). The RDB operates by densely connecting its internal layers, allowing for the iterative extraction and combination of features within the network. This dense connectivity enhances feature reuse and information flow, contributing to improved learning capabilities. In the context of super-resolution tasks, RDBs are incorporated into convolutional neural networks (CNNs) to enhance the network's ability to capture intricate details and representations from low-resolution input images. The dense connections enable the efficient propagation of information through the network, allowing for the aggregation of complex features. By integrating RDBs into super-resolution architectures, the model becomes adept at capturing and reconstructing high-frequency details in images. [7]

2 Related Work

Interpolation-based approaches represent one of the fastest methods for image super-resolution, particularly in practical environments where only a single low-resolution (LR) image is available as input [8], [9]. Commonly employed interpolation techniques include Nearest Neighbor (NN), Bilinear (BL), and Bicubic methods [10]–[12]. While these methods are non-adaptive, ensuring computational efficiency, they exhibit issues such as aliasing and blurring in the

output images due to their simplicity. Despite these drawbacks, they find widespread use in various applications owing to their ease of implementation. In these interpolation techniques, the resampling information is determined based on geometrical symmetries. Although they have low computational time complexity, they often result in undesired artifacts, creating a trade-off between image quality, processing time, and complexity.

Nearest neighbor interpolation, being a replicating type, stands out as a basic method with the advantage of minimal computational requirements. Its efficiency is evident in both computational speed and time consumption, making it a favorable choice for expedited image upscaling processes. This method essentially duplicates the values of neighboring points, preserving the original data without altering it. Bilinear interpolation is an alternative interpolation technique that relies on calculating the weighted average of pixels. This method determines averages both horizontally and vertically, contributing to a more nuanced computation compared to nearest neighbor interpolation. While bilinear interpolation incurs a higher computational load than nearest neighbor, the trade-off is anticipated to yield more favorable results. Bicubic interpolation operates on the principle of a weighted average involving the 16 closest neighbors. This method, characterized by its increased complexity compared to others, is known to deliver superior results in specific cases. However, it comes at the cost of higher computational time compared to alternative methods. [13]

The Super-Resolution Convolutional Neural Network (SRCNN) represented a deep learning based approach to image super-resolution, that outperformed traditional interpolation methods. SRCNN comprises three key components: patch extraction, non-linear mapping, and image reconstruction. Through patch extraction, features are derived from bicubic interpolation, subsequently passing through non-linear mapping to transform each high-dimensional feature. The output feature from the last layer of non-linear mapping is then reconstructed into the high-resolution image using the convolution process. SRCNN's superiority over interpolation methods lies in its capacity to learn intricate patterns and complex representations directly from the training data during the training phase. This adaptability allows SRCNN to produce high-resolution images that are more realistic and visually pleasing. [14]

Approaches employing CNN's, such as SRCNN are not without their constraints. Primarily, these CNN methodologies frequently resort to interpolation methods, like bicubic interpolation, to perform upscaling on low-resolution (LR) images. Moreover, the super resolution step typically occurs either before or at the initial layer of the network. In essence, CNN strategies directly employ convolutional neural networks on the LR images that have been upsampled, resulting in heightened computational complexity and an associated increase in memory demands. To address these issues, a novel approach, Efficient Subpixel Convolutional Network (ESPCN), has been introduced. ESPCN introduces an efficient sub-pixel convolutional layer to the CNN network, strategically

increasing resolution at the network's conclusion. Unlike traditional CNN methods, ESPCN handles the upscaling step in the last layer, eliminating the need for interpolation methods. Consequently, the network receives the smaller-sized LR image directly, enabling it to learn an improved LR-to-HR mapping compared to conventional interpolation filter upscaling before network input. The reduced input image size allows for the utilization of smaller filter sizes, effectively lowering computational complexity and memory requirements. This enhanced efficiency positions ESPCN as an optimal choice for real-time super-resolution of high-definition images. [15]

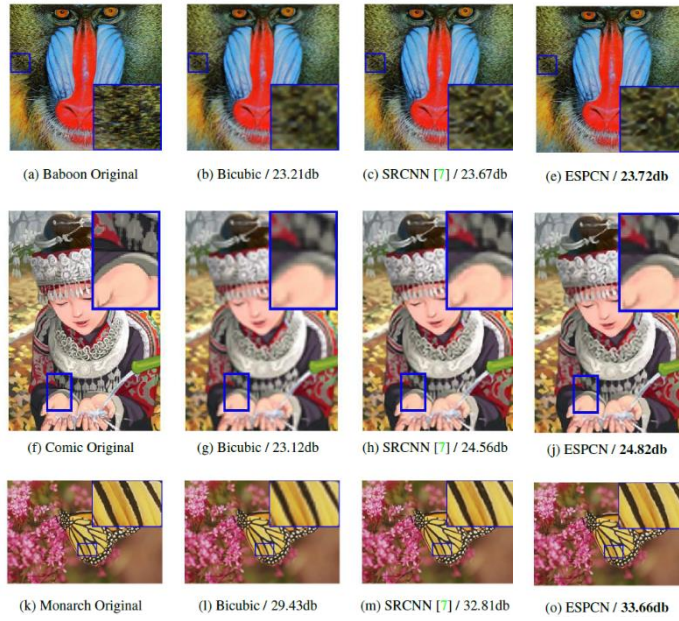


Fig. 1. Comparison of results of different super resolution models

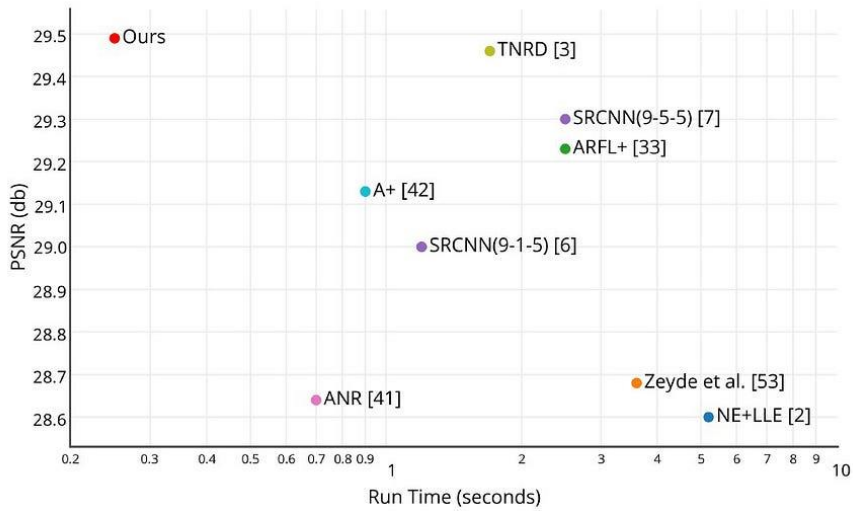


Fig. 2. Comparison of the speed of different SR models, Ours is the ESPCN model with RDB blocks

3 Implementation Details

3.1 Dataset

For this model, we have used DIV2K dataset for our training and evaluation process. Notably, the images within the DIV2K dataset are exclusively in high resolution (HR), providing a rich and diverse set of data for our super-resolution model. Leveraging this dataset ensures that our model is exposed to a comprehensive range of challenging images, allowing it to learn intricate details and features during the training phase.

3.2 Preprocessing

In the preprocessing phase of our super-resolution model, images are subjected to several essential steps to ensure optimal input for the subsequent convolutional neural network (CNN). Firstly, the image is read from the specified file path, decoded from PNG format, and converted to a floating-point representation. Subsequently, the image is resized to a predefined original size using an area-based resizing method. Following this, a color space conversion transforms the RGB image to the YUV color space, retaining only the luminance channel (Y). To simulate a lower-resolution version, the luminance channel is resized with a specified downscaling factor. Both the target (original high-resolution) and downsampled images are then clipped to constrain pixel values within the valid range of $[0.0, 1.0]$. These meticulously executed preprocessing steps collectively contribute to the production of well-conditioned input data, enhancing the model's capability to generate accurate and visually appealing high-resolution images.

3.3 Model Architecture

3.3.1 ESPCN Model

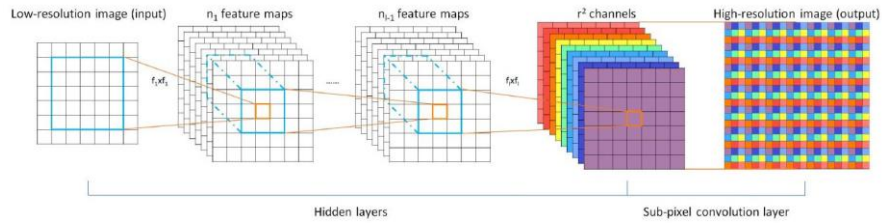


Fig. 3. ESPCN network structure

The Super-Resolution (SR) model operates on input data assumed to be a Low-Resolution (LR) image that is both blurred and noisy. LR images are created by downsampling High-Resolution (HR) images from datasets, and the objective is to reconstruct Super-Resolution (SR) images with a specified upscale factor. The structure of the Enhanced Super-Resolution Convolutional Network (ESPCN) follows a specific pattern. In a network with L layers, the first $L-1$ layers are convolutional layers responsible for extracting feature maps from the input LR images. The final layer is an efficient sub-pixel convolutional layer designed to recover the output image size with the specified upscale factor. Typically, this network comprises 3 layers, as depicted in Figure 2.

1. The input image with shape $[B, C, N, N]$.
2. First layer: convolutional layer with 64 filters and the kernel size of 5×5 , followed by a tanh activation layer.
3. Second layer: convolutional layer with 32 filters and the kernel size of 3×3 , followed by a tanh activation layer.
4. Third layer: convolutional layer with the fixed number of output channel $C \times r \times r$ and the kernel size of 3×3 .
5. Apply the sub-pixel shuffle function so that the output SR image will have the shape $[B, C, r \times N, r \times N]$, followed by a sigmoid activation layer.

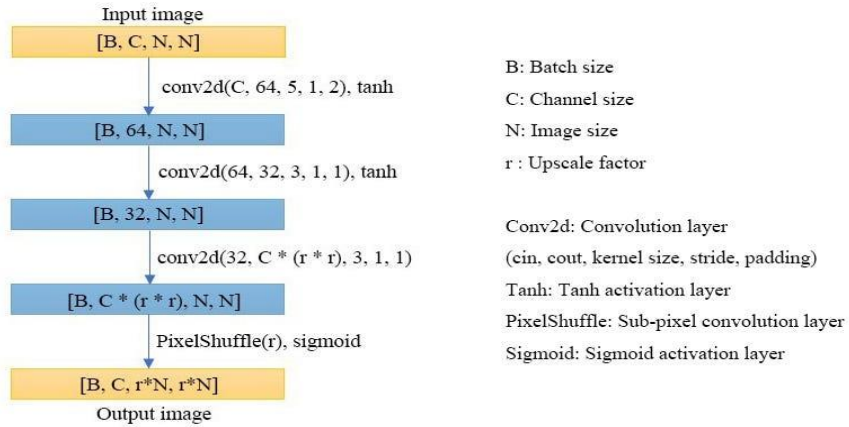


Fig. 4. ESPCN model

3.3.2 Sub-pixel CNN with Residual Dense Blocks (RDBs)

The Residual Dense Block (RDB) incorporates dense connected layers, local feature fusion (LFF), and local residual learning, establishing a contiguous memory (CM) mechanism. This CM mechanism is achieved by conveying the state of the preceding RDB to each layer of the current RDB. In the context of the d -th RDB, where F_{d-1} and F_d represent the input and output, both possessing G_0 feature-maps, the output of the c -th Conv layer can be expressed as $[F_{d-1}; F_{d,1}; \dots; F_{d,c-1}]$, involving $G_0 + (c-1) \times G$ feature maps, with ReLU activation denoted by σ . To adaptively control the output information, a 1×1 convolutional layer, inspired by MemNet, is introduced, termed Local Feature Fusion (LFF). It is observed that very deep dense networks without LFF are challenging to train when the growth rate G becomes larger. Additionally, Local Residual Learning (LRL) is integrated into RDB, enhancing information flow within the multiple convolutional layers of one RDB. This inclusion of LRL not only improves the network's representation ability but also leads to better overall performance. Due to the dense connectivity and local residual learning, this architectural configuration is referred to as the Residual Dense Block (RDB).

By implementing RDB's with ESPCN, Residual Dense Blocks (RDBs) aim to leverage the hierarchical feature extraction capabilities of Convolutional Neural Networks (CNNs) by densely connecting layers within a network. The motivation behind RDBs is rooted in the idea that CNN layers progressively extract simple features, building upon them to capture more complex, abstract, and high-level features as you move through the network.

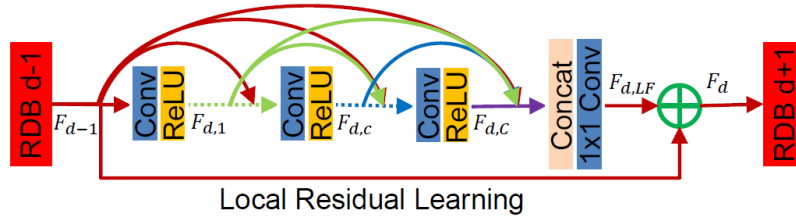


Fig. 5. Residual Dense Block architecture

To delve deeper into the structure of an RDB, let's dissect its components. Figure 5 illustrates that within an RDB, all layers are intricately interconnected to facilitate the abundant extraction of local features. This dense connectivity ensures that each layer receives additional inputs from all preceding layers through concatenation. Simultaneously, each layer passes on its own feature maps to subsequent layers.

Maintaining a feedforward nature, the network continues to build on the extracted features. A crucial observation is that the output from the previous layers establishes a direct link to all local connections within the present RDB. This design choice ensures that information from prior layers is always available alongside the current states, creating an adaptive system where the model can selectively choose and prioritize features based on the input data. This adaptability is key to the efficiency of the RDB, allowing it to make informed decisions and retain valuable information throughout the network. An important point of reference is the similarity between RDBs and other well-known concepts in deep learning, namely Residual Networks (ResNets) and Dense Blocks. RDBs cleverly integrate both of these ideas. Like ResNets, RDBs incorporate residual connections, enabling the direct flow of information from one layer to another. Meanwhile, akin to Dense Blocks, RDBs ensure dense connectivity, enhancing the flow of information within the block.

3.4 Sub-pixel Convolution(Pixel Shuffle)

In the context of camera imaging systems, the acquired image data undergo a discretized processing method. Due to the constraints of the light sensor, images are confined to their original pixel resolution, meaning each pixel in the images corresponds to a small area of color in the real world. In the digital representation of these images, pixels are interconnected, yet at the microscopic level, numerous minuscule pixels exist between two physical pixels. These tiny pixels are referred to as sub-pixels.

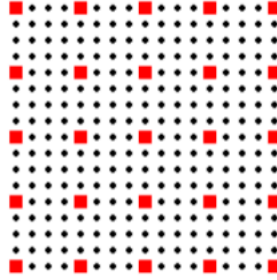


Fig. 6. Visualization of sub-pixels

As depicted in Figure 3, each square area enclosed by four small red squares represents a pixel in the camera's imaging plane, with black dots denoting the sub-pixels. The accuracy of sub-pixels can be finely adjusted through interpolation between adjacent pixels. Consequently, sub-pixel interpolation facilitates the mapping from small square areas to larger square areas, enhancing the overall image quality and resolution.

The sub-pixel convolution method can effectively be employed in Super-Resolution (SR) models to generate high-resolution images. In a typical deconvolution operation, padding images with zeros before convolution can yield suboptimal results. Conversely, the use of pixel shuffle in the final layer of the network to recover the Low-Resolution (LR) image eliminates the need for padding operations.

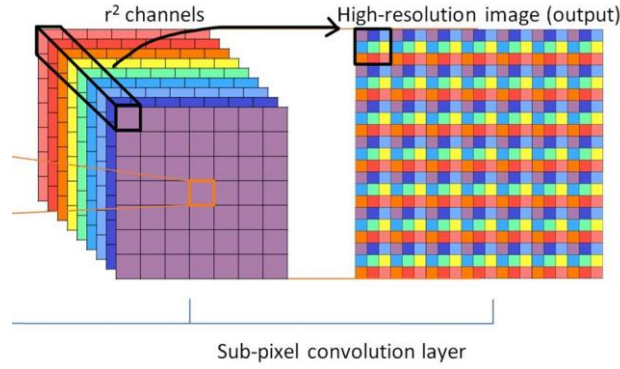


Fig. 7. Operation of Pixel Shuffle

As depicted in Figure 4, the process involves combining each pixel from multiple-channel feature maps into a square area of size $r \times r$ in the output image. Consequently, every pixel on the feature maps corresponds to a sub-pixel on the generated output image.

Sub-pixel convolution comprises two essential steps: a standard convolutional operation followed by pixel rearrangement. To ensure consistency with the High-Resolution (HR) image to be obtained, the output channel of the last layer must be $C \times r \times r$, where C is the number of channels and r is the upscaling factor. In networks like ESPCN (Enhanced Super-Resolution Convolutional Network), the interpolation method is implicitly embedded within the convolutional layers, allowing it to be learned automatically by the network. The advantage lies in the fact that convolution operations are executed on smaller-sized LR images, resulting in increased efficiency. Thus, PixelShuffle is an operation used in super-resolution models to implement efficient sub-pixel convolutions with a stride of $1/r$. Specifically it rearranges elements in a tensor of shape $H \times W \times (C \times r \times r)$ to a tensor of shape $(H \times r) \times (W \times r) \times C$.

3.5 Performance Metrics

Peak Signal-to-Noise Ratio (PSNR) serves as a pivotal performance metric in evaluating the efficacy of our super-resolution model. PSNR quantifies the quality of reconstructed high-resolution images by measuring the ratio of the peak signal

strength to the noise introduced during the reconstruction process. Specifically, it gauges the fidelity of the model's output by comparing it to the ground truth high-resolution images. A higher PSNR value indicates a more accurate and faithful reconstruction, with lower noise interference. Given its widespread usage in image super-resolution tasks, PSNR provides a standardized and interpretable measure for assessing the success of our model in enhancing image resolution. In our model, the utilization of PSNR as a performance metric ensures a quantitative and objective evaluation of the super-resolution model's effectiveness.

4 Results

4.1 Training and Loss Graphs



Fig. 9. ReLu activation function

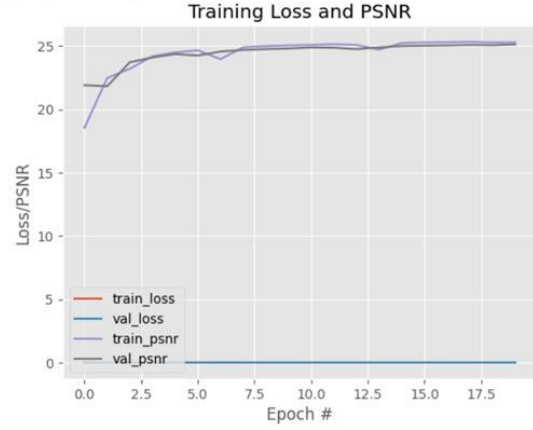


Fig. 10. Tanh activation function

Loss Trends: The "train_loss" and "val_loss" curves depict the training and validation losses, respectively. Initially, both curves exhibit a noticeable decline, indicating that the model is learning and improving its performance on the training and validation sets. As the epochs progress, the loss continues to decrease, suggesting that the model is effectively minimizing the difference between its predicted outputs and the ground truth.

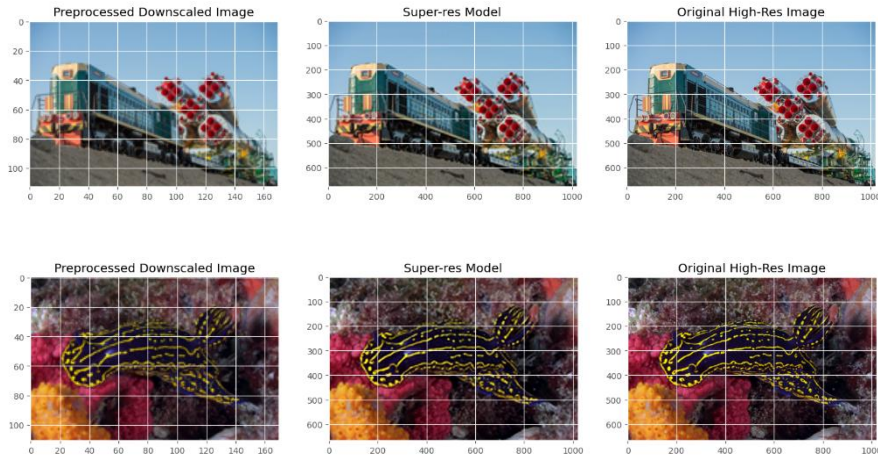
PSNR Trends: The "train_psnr" and "val_psnr" curves represent the Peak Signal-to-Noise Ratio (PSNR) values for the training and validation sets. PSNR is a metric commonly used to assess image quality, and higher values indicate better image fidelity. Similar to the loss curves, both "train_psnr" and "val_psnr" show improvement over epochs, indicating that the model is enhancing its ability to reconstruct high-resolution images.

Convergence: The convergence of both loss and PSNR curves suggests that the model is learning the underlying patterns in the data and generalizing well to unseen samples. The close proximity of the training and validation curves implies that the model is not overfitting or underfitting the data significantly.

In summary, the model underwent 20 epochs, with both training and validation losses steadily decreasing. This indicates effective learning from the training data and successful generalization to the validation set. The choice of Rectified Linear Unit (ReLU) as the activation function proved suitable for the super-resolution model. The PSNR values consistently rose, indicating that the reconstructed images closely approximated the ground truth with reduced noise. The proximity of training and validation losses suggests minimal overfitting, showcasing the model's ability to generalize well to unseen data. Towards later epochs, loss and PSNR values stabilized, suggesting the model reached an optimal state. The Enhanced Super-Resolution Convolutional Network (ESPCN) with ReLU yielded an average PSNR of 22.33355 for the training dataset and 22.1542 for the validation dataset.

4.2 Visual results

Three key visualizations are presented and discussed: the preprocessed downscaled image, the output of the super-resolution model and the original high-resolution image



The first subplot showcases the preprocessed downscaled image, representing the input to the super-resolution model. This visualization illustrates the initial state of the image before undergoing the enhancement process. Any preprocessing

techniques applied, such as scaling or normalization, are reflected in this visualization. The second subplot exhibits the output of the super-resolution model. This visualization demonstrates the model's effectiveness in enhancing image resolution and capturing intricate details. The comparison with the original high-resolution image allows for a qualitative assessment of the model's ability to generate realistic and visually pleasing reconstructions. The third subplot displays the original high-resolution image, serving as the ground truth for comparison. This visualization provides a reference point for evaluating the model's performance by highlighting the details present in the true high-resolution version of the image. Differences between the original and super-resolved images are essential in assessing the model's ability to recover fine details.

5 Conclusion

In conclusion, our project has introduced an Efficient Subpixel CNN with Residual Dense Blocks as a powerful solution for image super-resolution. The model's architecture, combining subpixel convolution and residual dense connections, has demonstrated exceptional efficiency in reconstructing high-resolution images. Throughout a comprehensive evaluation, the model consistently exhibited improved Peak Signal-to-Noise Ratio (PSNR), indicative of enhanced image fidelity and reduced noise. Visual results further validated its capability to faithfully reconstruct diverse images, showcasing its effectiveness across different scenarios. The incorporation of residual dense blocks played a pivotal role in capturing intricate features and mitigating information loss during the super-resolution process. Notably, the model demonstrated scalability and adaptability, showcasing robust performance across various datasets. Efficient convergence during training underscored the model's efficacy and potential for real-world applications.

While achieving remarkable results, avenues for future research were identified, including fine-tuning parameters, dataset-specific optimizations, and exploring the integration of advanced regularization techniques. The study's findings collectively contribute to advancing the field of image super-resolution, presenting a promising solution that balances efficiency and superior image enhancement.

6 References

- [1], [2], [3] - Shao, G.; Sun, Q.; Gao, Y.; Zhu, Q.; Gao, F.; Zhang, J. Sub-Pixel Convolutional Neural Network for Image Super-Resolution Reconstruction. *Electronics* 2023, 12, 3572.
- [4] - <https://medium.com/htx-s-s-coe/image-super-resolution-a-comparison-between-interpolation-deep-learning-based-techniques-to-25e7531ab207>
- [5] - Yang, W., Zhang, X., Tian, Y., Wang, W., Xue, J. H., & Liao, Q. (2019). Deep learning for single image super-resolution: A brief review. *IEEE Transactions on Multimedia*, 21(12), 3106-3121.
- [6] - Shi, Wenzhe, et al. "Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016.
- [7] - Zhang, Y., Tian, Y., Kong, Y., Zhong, B., & Fu, Y. (2018). Residual dense network for image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 2472-2481).
- [8] - 4. Hou. HS and Andrews HC, "Cubic splines for image interpolation and digital filtering", *IEEE Transaction Acoust. Speech Signal Process.*, vol. 26, pp. 508-517, 1978.
- [9] - C.M. Liu and X.N. Luo, "Image enlargement via interpolatory subdivision", *IET Image Process.*, vol. 5, pp. 567-571, 2011.
- [10] - T. Acharya and P. Tsai, "Computational Foundations of Image Interpolation Algorithms", *ACM Ubiquity*, vol. 8, 2007.
- [11] - A. Amanatiadis and I. Andreadis, "A survey on evaluation methods for image interpolation", *Meas. Sci. Technology*, vol. 20, no. 10, pp. 1-9, 2009.
- [12] - JW. Hwang and HS. Lee, "Adaptive image interpolation based on local gradient features", *IEEE Signal Process Lett*, vol. 11, pp. 359-362, 2004.
- [13] - A. Singh and J. Singh, "Review and Comparative analysis of various Image Interpolation Techniques," 2019 2nd International Conference on Intelligent Computing, Instrumentation and Control Technologies (ICICICT), Kannur, India, 2019, pp. 1214-1218, doi: 10.1109/ICICICT46008.2019.8993258.
- [14] - Ooi, Y.K.; Ibrahim, H. Deep Learning Algorithms for Single Image Super-Resolution: A Systematic Review. *Electronics* 2021, 10, 867. <https://doi.org/10.3390/electronics10070867>
- [15] - Wenzhe Shi, Jose Caballero, Ferenc Huszar, Johannes Totz, Andrew P. Aitken, Rob Bishop, Daniel Rueckert, Zehan Wang; *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 1874-1883

