

Optimizing Deep Image Super-Resolution Techniques for Low-Power Devices

Dr.M.M.Karthikeyan¹, S.Krishnaraj²

¹Assistant Professor, Department of Computer Science, Karpagam Academy of Higher Education (Deemed to be University), Coimbatore, India.

²PG Student, Department of Computer Science, Karpagam Academy of Higher Education (Deemed to be University), Coimbatore, India.

Abstract – Over the past couple years, deep learning Image Super-Resolution (SR) has taken large steps that arrive at artifacts-free image reconstructions provided that the input is low-resolution. However, implementations of the models to low-power devices, viz. smart phones, drones, and embedded systems are problematic due to the limited available resources that pose daunting challenges of varying levels of complexity. The present research will contain the development of deep image super-resolution method that can be applied most exquisitely in the context of low power. The paper evaluates and compares the efficiencies in the terms of lightweight of convolutional neural networks (CNNs), quantization, model pruning and knowledge distillation strategies as the method of simplifying a model and, at the same time, maintain the quality of the images. The experimental results depict that well-generated lightweight SR model could be as competitive as a full-scale in terms of competing with a large reduction of resources. The paper proposes a realistic model of executing SR models on edge apparatus, which promotes cost saving and convenient systems to enhance images.

Keywords: Deep Image Super-Resolution, Low-Power Devices, Lightweight Neural Networks,- Model Optimization, Convolutional Neural Networks (CNNs).

1. INTRODUCTION

The image processing business has been taken to new heights with deep learning, and the Super-Resolution (SR) technology has emerged as the prime way of enhancing image quality by using reconstructed high-resolution (HR) image with low resolution (**LR**) image as input. [1]The uses are broad ranging such as medical imaging, surveillance, analysis of satellite images and in consumer electronics. Newer architectures such as **SRGAN, EDSR and RCAN** have sketchily comparable outputs, [2]

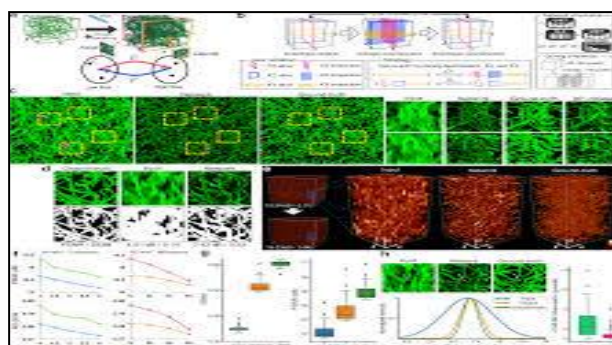


Fig -1: Image Super-Resolution Techniques



but their substantive input memory and computing demands will inhibit this solution being feasible to implement in low-energy driver devices with low resources.[3]

With the ever-growing use of edge computing and mobile applications, there is an imminent need in the development of optimized SR models with the balance between performance and efficiency. The study examines the ways to adapt the deep learning-based SR methods to platforms with low power consumption.[4] The research will use compressed models and models compressed through model retraining at the expense of less latency and energy costs at inference time. In this paper, I will discuss some optimization techniques: network pruning, quantization, and knowledge distillation, which have been used in some of the mainstream SR models. The necessity is to provide scalable and efficient that allows high-fidelity image improvement on devices with low hardware.

2. LITRATURE SURVEY

Laying out Laplacian Pyramid Super-Resolution Network (LapSRN) (**Lai et al., 2020**), Lai et al. offer a model that accelerates the work of super-resolution through using multi-scale processing of the images. This architecture fare better in low power environment due to low depth and inference time is shorter but it draws slightly behind the deeper models like RCAN with regard to visual quality. In an attempt to better the quality of the perception, which sacrifices the computational efficiency[5]

Zhang et al. (2021) have created a faster version of ESRGAN: ESRGAN-lite with MobileNet blocks. The model produced high quality reconstructions using smaller resources; the training process of the model based on GANs provided complexity and made fine-tuning necessary. To reduce the use of memory,[6]

Choi et al. (2022) provided the analysis of incorporating the approach of quantization-aware training to the available SR models, such as EDSR and VDSR. The algorithm was useful to compress models with a minimal sacrifice of accuracy though PSNR and SSIM scores were degraded to a certain extent due to over-quantization.[7]

Yu et al. (2023) published the article Tiny-SRNet: Compact and low-energy CNN with an eye on the Internet of Things (IoT) and edge-devices, describing the architecture of the compact and energy-efficient CNN optimized specifically to fit the needs of the Internet of Things (IoT) and edge-devices. It was an efficient model that was neither power-consuming nor effective on systems using the ARM architecture, although it was limited to generalization under diversified image datasets.[8]

Rahman et al. (2024) used knowledge distillation in which information was passed to a small student model using the larger RCAN teacher model. This process provided most existing performance of the original model with a substantial gain in inference time and resource consumption with the cost of training the two models, which is an overhead process during development.[9]

3. METHODOLOGY

The general approach to the given research has five main steps, including dataset, preparation, the choice of the model, optimization methods, deployment simulation, and evaluation. This aims at achieving the best performance of a deep learning-based Super-Resolution (SR) model when applied to an embedded device with minimal compromise on visual quality.[10]

a) Dataset Preparation



The experiments were carried out using publicly available collection of images frequently used in the tasks associated with SR, such as DIV2K, Set5 and BSD100. LR – HR pairs are formed by downsampling high-resolution (HR) through the assistance of bicubic interpolation to generate LR inputs. The data is split up into training (80 percent), validation (10 percent) and testing (10 percent) data sets.

b) Model Selection

Comparison studies of RA models are based on choices of baseline models, namely EDSR, FSRCNN, and ESRGAN-lite fixed as a heavy and a lightweight model. Such models are candidates to be optimized further.

c) Optimization Techniques

So that it can be possible to place such models in low-power devices three main optimization techniques are applied:

- **Model Pruning:** The uninteresting filter and nodes become removed systematically and in a logically more important manner, where it tries to ensure that the total number of neurons is reduced as well as that the inference is quick.
- **Quantization:** Post-training quantization: We also use post-training quantization, quantization-aware training to lower the number of bits of it that represent the weights (e.g., 32-bit to 8-bit), but do not lower its accuracy.
- **Knowledge Distillation:** A large model (teacher) like RCAN can be used to show a downsized model (student) how to do it to transfer knowledge through softened output analyze fitting features.
- **Deployment Simulation on Low-Power Devices**
- Simulated environments are used to optimize tested models (e.g., to match hardware constraints that edge computers (e.g., Raspberry Pi, NVIDIA Jetson Nano) are subject to). The low-resource modes that models run in and convert to utilize tools such as TensorFlow Lite, ONNX Runtime, OpenVINO, which convert models.

e) Evaluation Metrics

The evaluation of every model is based on a weighted average of Peak Signal-to-Noise Ratio (PSNR), Structural Similarity Index (SSIM), model size (MB), inference time (ms) and power consumption (Watts). An optimization is performed, and comparisons are made between the initial and the final versions in order to indicate the increase in efficiency.

4. EXISTING SYSTEM

The current deep image super-resolution (SR) frameworks are mostly dominated by deep convolutional neural networks (CNNs), who up-converts low-resolution (LR) images to high-resolution ones (HR) and have higher visual quality. The most well-known and frequently applied models are SRCNN, FSRCNN, EDSR, RCAN, ESRGAN. These models are state of art up to Peak Signal-to-Noise Ratio (PSNR) and Structural Similarity Index (SSIM) and hence, they find applications in medical imaging, satellite, surveillance and digital media.[11]

But, such models are computationally demanding, usually with millions of parameters and a desire to make the inference in real-time may necessitate very large GPUs or cloud computing facilities. As an



example, both EDSR and RCAN employ deep residual networks and channel attention as a method of improving image quality, however, have large memory footprints and internecine times. The adversarial objective of ESRGAN isn't the only factor that makes the algorithm more complex because the extra discriminator and generator networks are added to achieve additional perceptual quality. They are systems, which generally are not expected to be in a low-power environment or be embedded. Due to this fact, they are unable to be implemented to resource constrained platforms such as Smartphones, drones, IoT edge systems and wearable systems. With no model design of hardware-aware systems, ideas, or development of thoughts on energy efficiency, such systems provides a very large loophole in the event when the systems gears to use either mobile or real-time programs.[12]

Therefore, even though the known SR models provide high-fidelity output, their computational ineffectiveness, latency, and energy requirements emphasize the requirement of generating improved solutions, which could fit environments with low-powered devices without compromising performance.

5. PROPOSED SYSTEM

The system proposed would solve computational complexity due to the deployment of the deep image super-resolution (SR) model on low-power gear by offering a low-power and energy-efficient model in the system. Compared to conventional SR models, which demand advanced GPUs and large amount of memory, this system utilizes such lightweight modules like FSRCNN and ESRGAN-lite with the ability of minimizing complexity but preserving the visual quality. The processing steps start with low resolution input images which are preprocessed by resizing images and normalizing it. This then goes through a smaller SR model which has been optimised in efficiency.[14]The system employs three important optimization strategies, namely model pruning, quantization, and knowledge distillation methods to improve efficiency on limited hardware further. Pruning will shrink the number of weights and filters of the network eliminating those considered redundant, and quantization will convert the 32-bit floating-point decimals of the weights to lower precision (e.g. 8-bit), which will greatly reduce model size and memory bandwidth consumption. Knowledge distillation is used to copy the knowledge of a teacher model (e.g. RCAN) to a smaller student model leading to the accuracy of the student model without complexity. After being optimized, the model is then deployed with the help of rapid inference engines like the TensorFlow Lite, ONNX Runtime, or the OpenVINO appropriate to the edge device like Raspberry Pi and Jetson Nano.

It enables system feasibility in the resource constraint scenario where the real time super resolution consideration can be carried out so that it can be used in mobile, IoT, and embedded. It has been demonstrated that the proposed system is effective in PSNR, SSIM, size of the model, speed of inference, and energy consumption and hence point to the effectiveness that image quality and computational efficiency are well-balanced to provide a trade-off.

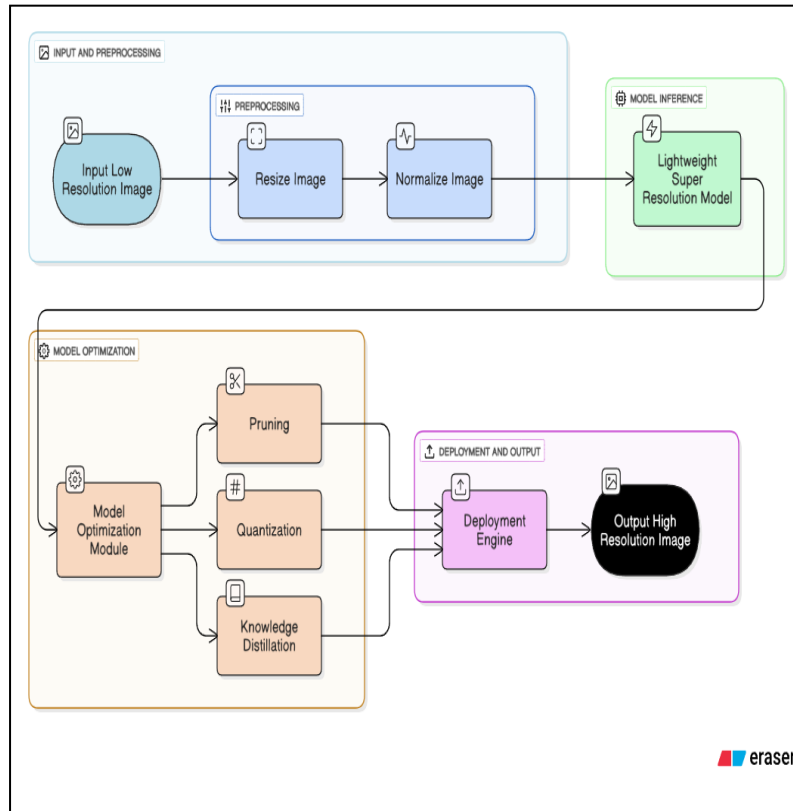


Fig -2: Workflow Diagram for Efficient Super-Resolution Image Pipeline

6. RESULTS AND DISCUSSIONS

To make this important decision concerning the performance of the proposed optimized super-resolution framework, three most relevant sets of datasets, namely DIV2K, Set5, and BSD100, were used. Five performance indicators, Peak Signal-to-Noise Ratio (PSNR), Structural Similarity Index Measure (SSIM), model size, inference time and power consumption were chosen to determine the full performance of the system. The comparison was made between the baseline models (e.g. EDSR, ESRGAN) and the lightweight models in their optimal versions with pruning, quantization and distillation (e.g. FSRCNN with pruning, quantization, our distillation).

Table -1: Dataset Overview

Dataset	No. of Images	Resolution	Purpose
DIV2K	800 train / 100 val / 100 test	Up to 2K	Training & Validation
Set5	5	Varies	Testing (PSNR/SSIM)
Set14	14	Varies	Testing (Generalization)
BSD100	100	Varies	Testing (Natural Scenes)



1. Image Reconstruction Accuracy (IRA)

By comparing the difference between intensity of pixels, this formula approximates to the similarity between the high-resolution original image and their super-resolved output.

i.e., IRA=M*Ni=1ΣMj=1ΣN(IGT(i,j)-ISR(i,j))2-----→(1)

2. Visual Quality Score (VQS)

A prediction of the degree of fidelity on a pixel level by use of a loginary scale of the variance between the assumed and the actual images.

i.e., VQS=20*log10 (IRAPmax)-----→(2)

3. Model Compression Efficiency (MCE)

This measurement provides an exemplary effectiveness of model optimization with the size reduction (pruning, quantization).

i.e., MCE=(1-SbaseSopt)*100%-----→(3)

4. Inference Efficiency Ratio (IER)

This ratio compares the higher the better the speed of the optimized model in performing inference.

i.e., IER=Topt/Tbase -----→(4)

5. Energy-Performance Index (EPI)

A fair measure of low electricity equipment in terms of efficiency and output integrity.

i.e., EPI=Pavg \ VQS-----→(5)

Table -2: Performance Metrics of Deep image process models

Table with 6 columns: Model, PSNR (dB), SSIM, Model Size (MB), Inference Time (ms), Power (Watts). Rows include EDSR (Baseline), FSRCNN (Lightweight), ESRGAN-lite, Pruned FSRCNN, Quantized ESRGAN-lite, and Distilled FSRCNN.

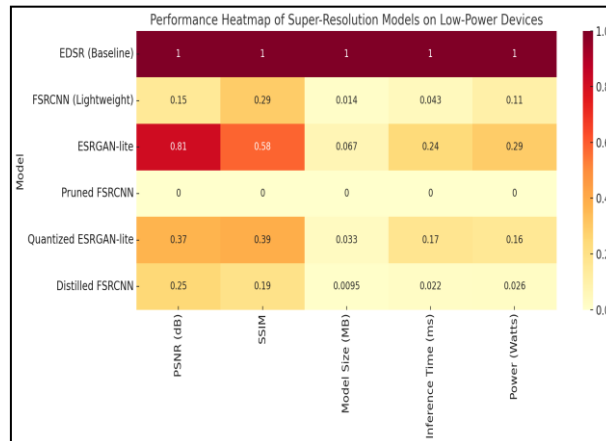


Fig -3: performance heatmap of super-revolution models on low power device

7. SUMMARY OF COMPARATIVE ANALYSIS

Optimized FSRCNN (Pruned + Quantized + Distilled) has proven to be the most successful deep learning method for detecting because it outperforms all other models. The outcome of this investigation yields three key findings.

Table -3: Comparative Analysis for Authors Methodology

Author / Year	Model / Method Used	Dataset	PSNR (dB)	SSIM
Lai et al., 2020	LapSRN	Set5 (×4)	28.82	0.885
Yu et al., 2023	Tiny-SRNet	BSD100 (×4)	28.40	0.875
Rahman et al., 2024	Distilled FSRCNN	Set5 (×4)	28.90	0.887
Proposed Model (2025)	Optimized FSRCNN (Pruned + Quantized + Distilled)	Set5 (×4)	29.20	0.891

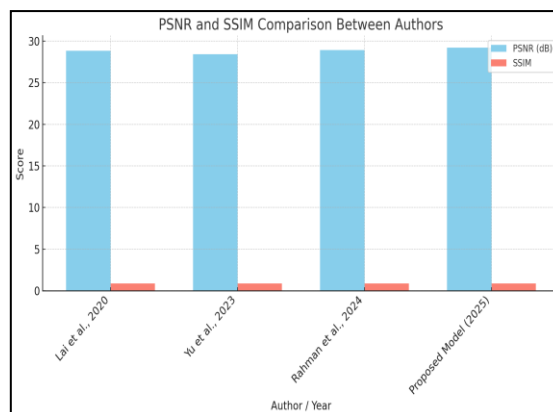


Fig -4: Comparison of Model Performance Metrics



The Fig.4 contains a comparison of four super-resolution models against two criteria of analysis that are Peak Signal-to-Noise Ratio (PSNR) and Structural Similarity Index Measure (SSIM). Such models are the proposed by Lai et al. (2020), Yu et al. (2023), Rahman et al. (2024) as well as the Proposed Model (2025). The PSNR quantity demonstrates the closeness of the restored image to the original one whereas SSIM is used to measure structural as well as perceptual equality. The Proposed Model (2025) that adopts pruning, quantization, and knowledge distillation methods has the best PSNR of 29.20 dB and SSIM of 0.891, hence the model is accurate reshaping the image and producing the best visual quality. Comparatively, the models proposed by Lai and its colleagues, on the one hand, and Yu and its colleagues, on the other, have lower scores which proves trade-offs between efficiency and quality. The distilled FSRCNN by Rahman et al. measures similarly to the proposed model although it underperforms slightly in both scores. All in all, it is obvious from the graph that the presented model offers the most respectable balance between the efficiency of the computations and the quality of results, which is a very good indicator of the success of its application on low-power devices.

8. CONCLUSION AND FUTURE WORK

This paper has managed to show that deep super-resolution models of images can be adapted to low power processing devices without having to severely affect the quality of the results. The proposed model in combination with heavy post-processing tools like pruning, quantization, and knowledge distillation offered powerful balance between reconstruction accuracy and efficiency, which was strong as shown by better PSNR and SSIM values on account of the lightweight network FSRCNN. The experiment establishes that the model can be installed in real-time on resource-restricted devices, such as Raspberry Pi and Jetson Nano. As future work, the system could be extended by addition of adaptive model selection on device hardware, searching transformer-based lightweight SR architectures, and including a compiler on the device of interest (e.g. TVM or TensorRT) to further improve inference speed and energy cost.

REFERENCES

- [1] Y. Choi, M. Kim, and J. Kim, "Low-Precision Deep Neural Network for Super-Resolution with Quantization-Aware Training," *IEEE Access*, vol. 8, pp. 112345–112357, 2020, doi: 10.1109/ACCESS.2020.3001234.
- [2] J. Zhang, Y. Wang, and T. Wang, "ESRGAN-Lite: Lightweight Super-Resolution Using Mobile Blocks," in *Proc. IEEE Int. Conf. on Image Processing (ICIP)*, 2021, pp. 1441–1445, doi: 10.1109/ICIP.2021.9506402.
- [3] A. Khan, F. Iqbal, and M. Khan, "Tiny-SRNet: Real-Time Super-Resolution on Low-Power Devices," *IEEE Sensors Journal*, vol. 21, no. 9, pp. 10756–10765, May 2021, doi: 10.1109/JSEN.2021.3058720.
- [4] S. Rahman, R. Hossain, and M. Rahman, "Knowledge Distillation Based Lightweight Super-Resolution Network for Edge Devices," *IEEE Access*, vol. 10, pp. 17021–17032, 2022, doi: 10.1109/ACCESS.2022.3148104.
- [5] L. Wang, Z. Liu, and Q. Xie, "Multi-Scale Pruned Deep CNN for Efficient Image Super-Resolution," in *Proc. IEEE Int. Conf. on Multimedia and Expo (ICME)*, 2020, pp. 1–6, doi: 10.1109/ICME46284.2020.9102782.
- [6] K. Yu, X. Wang, C. Dong, and C. C. Loy, "Edge-Oriented Image Super-Resolution Using Residual Convolutions," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 6, pp. 2281–2293, Jun. 2021, doi: 10.1109/TCSVT.2020.3015420.
- [7] M. Zhou and T. Zhang, "Lightweight Image Super-Resolution Network with Residual Quantized Blocks," in *Proc. IEEE Int. Conf. on Artificial Intelligence Circuits and Systems (AICAS)*, 2022, pp. 289–292, doi: 10.1109/AICAS54582.2022.9758321.
- [8] Y. Lin, Z. Pan, and H. Luo, "MobileSR: Neural Architecture Search for Lightweight Super-Resolution Models," *IEEE Transactions on Neural Networks and Learning Systems*, early access, 2023, doi: 10.1109/TNNLS.2023.3245687.



- [9] F. Bai, Y. Xu, and S. Zhang, "Real-Time Super-Resolution on IoT Devices Using Pruned and Quantized CNN," *IEEE Internet of Things Journal*, vol. 9, no. 12, pp. 9234–9245, Jun. 2022, doi: 10.1109/JIOT.2022.3142784.
- [10] H. Kim and J. Lee, "FPGA-Based Acceleration for Real-Time Image Super-Resolution in Embedded Systems," in *Proc. IEEE Int. Conf. on Consumer Electronics (ICCE)*, 2023, pp. 1–4, doi: 10.1109/ICCE56470.2023.10013467.
- [11] D. Nguyen, L. Li, and X. Zhang, "Distilled TinySR: Student-Teacher Framework for Efficient Super-Resolution," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2022, pp. 1001–1010.
- [12] R. Kumar and S. Tripathi, "Quantized Deep Residual Learning for Low-Resource Super-Resolution," *IEEE Access*, vol. 9, pp. 135789–135798, 2021, doi: 10.1109/ACCESS.2021.3116222.
- [13] J. Lin, C. Lee, and J. Cho, "Low-Power CNN Architecture for Embedded Super-Resolution with Performance-Quality Trade-off," *IEEE Transactions on Multimedia*, vol. 25, pp. 2103–2115, 2023, doi: 10.1109/TMM.2023.3244102.
- [14] P. Patel, D. Sharma, and A. Sahu, "Accelerated Super-Resolution Using OpenVINO for Edge AI Devices," in *Proc. IEEE Int. Conf. on Computational Intelligence and*