

**NOVEL DEEP LEARNING BASED  
FRAMEWORK FOR IMAGE AND VIDEO  
ENHANCEMENT**

**MANOJ SHARMA**



**DEPARTMENT OF ELECTRICAL ENGINEERING  
INDIAN INSTITUTE OF TECHNOLOGY DELHI**

**JANUARY 2020**

©Indian Institute of Technology Delhi (IITD), New Delhi, 2020

**NOVEL DEEP LEARNING BASED  
FRAMEWORK FOR IMAGE AND VIDEO  
ENHANCEMENT**

by

**MANOJ SHARMA**

DEPARTMENT OF ELECTRICAL ENGINEERING

Submitted

in fulfillment of the requirements of the degree of Doctor of Philosophy

to the



INDIAN INSTITUTE OF TECHNOLOGY DELHI

JANUARY 2020

To my family

# Certificate

This is to certify that the thesis titled **Novel Deep Learning based Framework for Image and Video Enhancement** being submitted by **Mr. Manoj Sharma** to the Department of Electrical Engineering, Indian Institute of Technology Delhi, for the award of **Doctor of Philosophy** is a record of bona-fide research work carried out by his under my guidance and supervision. In my opinion, the thesis has reached the standards fulfilling the requirements of the regulations relating to the degree. The work presented in this thesis has not been submitted elsewhere, either in part or full, for the award of any other degree or diploma.

Professor Santanu Chaudhury  
Department of Electrical Engineering  
Indian Institute of Technology Delhi  
New Delhi - 110016, India.

Professor Brejesh Lall  
Department of Electrical Engineering  
Indian Institute of Technology Delhi  
New Delhi - 110016, India.

# Acknowledgements

I would like to express my sincere gratitude to my advisors **Prof. Santanu Chaudhury** and **Prof. Brejesh Lall**. This thesis would not have been possible without their guidance, critical review and encouragement. Their intelligent ideas and comprehensive understanding helped me to establish the overall direction of the research, got me past many significant challenges and always motivated me to strive for excellence.

I thank the members of my thesis committee: Prof. Prem Kalra, Dr. Sumantra Dutta Roy and Dr. Sumeet Agarwal for their insightful suggestions which encouraged me to widen my research from various perspectives. I sincerely admire the contribution of my lab mates Ms. Swati Bhugra, Ms. Ritu Chaudhury, Ms. Anupama Ray, Mr. Raunak Gupta and Ms. Nidhi Chahal for extending their unstinting support and sympathetic attitude during the course of this study. I am grateful to my wife for supporting me during my Ph.D. My thanks to Ushaji, Manrika Mam and Amit Bhaiya for all the timely logistic support.

I acknowledge the significant contribution of my parents. Their patience and support will always be my inspiration. I would like to express my heartiest gratitude to my friends, Satyam and Sudhakar, who stood by me through all my efforts and encouraged me to pursue my Ph.D.

Manoj Sharma

# Abstract

Image enhancement is a technique to improve the interpretability or information perception in images for human understanding. The enhanced image can also be used as an input to other image processing applications. However, image enhancement is extremely challenging due to lack of sufficient prior information and is still an open problem. Enhancement encapsulates various task such as image super-resolution, denoising, contrast enhancement, noise-free super-resolution to name a few. These enhancement techniques serve as a pre-processing step in various computer vision applications such as optical character recognition (OCR), bio-medical etc. The advancement in image enhancement techniques has resulted in the requirement of efficient frameworks that can evaluate the quality of image/video data automatically. Objective quality estimation approaches can be applied to quality control systems for monitoring the image and video quality, parameter setting and optimization of benchmark algorithms for image and video processing system.

The main objective of the thesis is image and video enhancement. Among various enhancement techniques, image and video super-resolution (SR) has gained a lot of attention because of increasing availability of high-resolution (HR) displays. This is a challenging problem due to lack of prior information and image statistics. For images, we focus on spatial resolution enhancement. In addition, noise-robust resolution enhancement of noisy low-resolution (LR) images is even more challenging as both image de-noising and super-resolution tasks are ill-posed inverse problems. If an LR image contains noise or distortions, a simple resolution enhancement task gives rise to additional incremental noise or distortions. Moreover, performing de-noising followed by super-resolution does not provide HR image with fine details. Thus,

there is a need for joint optimization of image de-noising and super-resolution while preserving fine details. In case of video, resolution enhancement means enhancing resolution in both space and time. For video enhancement, our work focuses on finding a solution in order to get space-time resolution enhancement of the LR video sequences.

Image in-painting is a restoration technique that focuses on removing distorted parts or filling out missing regions to enhance images. However, in-painting of larger missing region with local details has not yet been solved. In this, we present a solution for filling larger missing regions by using super-resolution based in-painting. As previously mentioned, the solutions of the aforementioned enhancement tasks depend on image/video quality assessment. For quality assessment, past efforts have largely focused on specific kind of distortion. However, it is not practically feasible to define different type of distortions and their causes. Thus our work focuses on a technique which is independent of the type of error present and concentrates on image quality rather than on the type of specific error.

Optical character recognition (OCR) systems face severe challenges in recognizing text from document image that lack resolution. Most OCR models trained on high resolution (HR) text images fail in case of low resolution (LR), noisy or noisy-LR text images. Enhancement of noisy-LR text images, such that existing OCR engine can recognize the text from noisy-LR images is of interest to the OCR community. In the thesis, we have demonstrated an application of noise-robust resolution enhancement technique to boost existing OCR systems.

## सार

इमेज एन्हांसमेंट व्याख्या या सूचना की धारणा को बेहतर बनाने की तकनीक है मानव समझ के लिए चित्र। बढ़ी हुई छवि को अन्य के लिए एक इनपुट के रूप में भी इस्तेमाल किया जा सकता है छवि प्रसंस्करण अनुप्रयोगों। हालाँकि, छवि वृद्धि बेहद चुनौतीपूर्ण है पर्याप्त पूर्व सूचना की कमी और अभी भी एक खुली समस्या है। एन्हांसमेंट बढ़ जाता है विभिन्न कार्य जैसे कि इमेज सुपर-रिज़ॉल्यूशन, डीनोइज़िंग, कंट्रास्ट एन्हांसमेंट, नॉइज़-फ्री सुपरसेलिंग कुछ नाम है। ये एन्हांसमेंट तकनीक प्री-प्रोसेसिंग कदम के रूप में काम करती है विभिन्न कंप्यूटर दृष्टि अनुप्रयोग जैसे कि ऑप्टिकल चरित्र पहचान (ओसीआर), जैव-चिकित्सा आदि छवि वृद्धि तकनीकों में प्रगति की आवश्यकता के परिणामस्वरूप हुई है कुशल चौखटे जो स्वचालित रूप से छवि / वीडियो डेटा की गुणवत्ता का मूल्यांकन कर सकते हैं। उद्देश्य निगरानी के लिए गुणवत्ता नियंत्रण प्रणालियों के लिए गुणवत्ता अनुमान दृष्टिकोण लागू किया जा सकता है छवि और वीडियो की गुणवत्ता, पैरामीटर सेटिंग और छवि के लिए बेंचमार्क एल्गोरिदम का अनुकूलन और वीडियो प्रसंस्करण प्रणाली।

थीसिस का मुख्य उद्देश्य छवि और वीडियो वृद्धि है। विभिन्न वृद्धि के बीच तकनीक, छवि और वीडियो सुपर-रिज़ॉल्यूशन (एसआर) ने बहुत अधिक ध्यान आकर्षित किया है उच्च-रिज़ॉल्यूशन (एचआर) डिस्प्ले की बढ़ती उपलब्धता। यह एक चुनौतीपूर्ण समस्या है पूर्व सूचना और छवि आँकड़ों की कमी। छवियों के लिए, हम स्थानिक संकल्प पर ध्यान केंद्रित करते हैं वृद्धि। इसके अलावा, शोर कम संकल्प (LR) शोर-मजबूत संकल्प वृद्धि छवियां और भी चुनौतीपूर्ण हैं क्योंकि दोनों चित्र डी-नॉइज़िंग और सुपर-रिज़ॉल्यूशन कार्यों को अनप्लग कर रहे हैं उलटा समस्या। यदि LR छवि में शोर या विकृतियां होती हैं, तो एक सरल रिज़ॉल्यूशन वृद्धि कार्य अतिरिक्त वृद्धिशील शोर या विकृतियों को जन्म देता है। इसके अलावा, प्रदर्शन सुपर रिज़ॉल्यूशन के बाद डी-नॉइज़ ठीक विवरण के साथ एचआर छवि प्रदान नहीं करता है। इस प्रकार, संरक्षण करते समय छवि डी-शोर और सुपर-रिज़ॉल्यूशन के संयुक्त अनुकूलन की आवश्यकता है बारीक विवरण। वीडियो के मामले में, रिज़ॉल्यूशन एन्हांसमेंट का मतलब दोनों में रिज़ॉल्यूशन को बढ़ाना है

स्थान और समय। वीडियो एन्हांसमेंट के लिए, हमारा काम पाने के लिए एक समाधान खोजने पर ध्यान केंद्रित करता है LR वीडियो दृश्यों का स्पेस-टाइम रिज़ॉल्यूशन एन्हांसमेंट। छवि-पेंटिंग एक पुनर्स्थापना तकनीक है जो विकृत भागों को हटाने पर केंद्रित है या छवियों को बढ़ाने के लिए लापता क्षेत्रों को भरना। हालाँकि, बड़े लापता क्षेत्र की पेंटिंग स्थानीय विवरण के साथ अभी तक हल नहीं किया गया है। इसमें, हम बड़े लापता को भरने के लिए एक समाधान प्रस्तुत करते हैं इन-पेंटिंग आधारित सुपर-रिज़ॉल्यूशन का उपयोग करके क्षेत्र। जैसा कि पहले बताया गया है, के समाधान उपरोक्त संवर्द्धन कार्य छवि / वीडियो गुणवत्ता मूल्यांकन पर निर्भर करते हैं। गुणवत्ता के लिए मूल्यांकन, पिछले प्रयासों ने काफी हद तक विशिष्ट प्रकार की विकृति पर ध्यान केंद्रित किया है। हालाँकि यह है विभिन्न प्रकार की विकृतियों और उनके कारणों को परिभाषित करने के लिए व्यावहारिक रूप से संभव नहीं है। इस प्रकार हमारा काम एक ऐसी तकनीक पर ध्यान केंद्रित करता है जो वर्तमान में मौजूद त्रुटि के प्रकार से स्वतंत्र होती है और उस पर ध्यान केंद्रित करती है विशिष्ट त्रुटि के प्रकार के बजाय छवि गुणवत्ता।

ऑप्टिकल कैरेक्टर रिकग्निशन (OCR) सिस्टम टेक्स्ट को पहचानने में गंभीर चुनौतियों का सामना करते हैं दस्तावेज़ छवि से जिसमें संकल्प की कमी है। उच्च रिज़ॉल्यूशन (एचआर) पर प्रशिक्षित अधिकांश ओसीआर मॉडल पाठ चित्र कम रिज़ॉल्यूशन (LR), शोर या शोर-LR पाठ छवियों के मामले में विफल होते हैं। वृद्धि शोर-एलआर पाठ छवियों, जैसे कि मौजूदा ओसीआर इंजन शोर-एलआर से पाठ को पहचान सकता है चित्र OCR समुदाय के लिए रुचि रखते हैं। थीसिस में, हमने एक एप्लिकेशन का प्रदर्शन किया है मौजूदा ओसीआर प्रणालियों को बढ़ावा देने के लिए शोर-मजबूत संकल्प वृद्धि तकनीक।

# Contents

<b>Certificate</b>	<b>iii</b>
<b>Acknowledgements</b>	<b>v</b>
<b>Abstract</b>	<b>vii</b>
<b>List of Figures</b>	<b>xiii</b>
<b>List of Tables</b>	<b>xvii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Scope and Objectives . . . . .	2
1.2 Context . . . . .	3
1.3 Major Contributions of the Thesis . . . . .	4
<b>2 Literature Review</b>	<b>7</b>
2.1 Introduction . . . . .	7
2.1.1 Quality Assessment . . . . .	8
2.1.2 Image Super-Resolution . . . . .	8
2.1.3 Video SR . . . . .	10
2.2 Noise-free Super-resolution . . . . .	12
2.3 Super-Resolution based Inpainting . . . . .	13
2.4 Document Image Enhancement . . . . .	13

2.5	Motivation for the present work . . . . .	15
<b>3</b>	<b>Sparse Representation based Classifier to assess Video Quality</b>	<b>17</b>
3.1	Introduction . . . . .	17
3.2	Methodology . . . . .	18
3.2.1	NSS based Feature extraction . . . . .	19
3.2.2	Classification using Sparse Representation Method . . . . .	21
3.3	Experiments . . . . .	22
3.3.1	Datasets . . . . .	22
3.3.2	Results . . . . .	24
3.4	Conclusion . . . . .	25
<b>4</b>	<b>Deep-learning based Image Super-Resolution</b>	<b>27</b>
4.1	Introduction . . . . .	27
4.2	Motivation and Contribution . . . . .	28
4.3	Methodology . . . . .	28
4.4	Experiments . . . . .	33
4.4.1	Image super-resolution using CDCA . . . . .	33
4.5	Results and Analysis . . . . .	33
4.5.1	Results for image super-resolution . . . . .	33
4.5.2	Conclusion . . . . .	35
<b>5</b>	<b>Robust Super-Resolution based Enhancement</b>	<b>37</b>
5.1	Deep Learning Framework for Noise-Resilient Super-Resolution . . . . .	37
5.1.1	Introduction . . . . .	37
5.1.2	Deep-learning based Noise Resilient Image Super-Resolution . . . . .	40
5.1.3	Experiments . . . . .	43
5.1.4	Image de-noising using SSDA . . . . .	43

---

5.1.5	Image super-resolution using CDCA . . . . .	44
5.1.6	Noise-resilient image super-resolution using SSDA-CDCA . . . . .	44
5.1.7	Results and Analysis . . . . .	44
5.1.8	Results for noise resilient image super-resolution . . . . .	45
5.1.9	Conclusions . . . . .	47
5.2	An End-to-End Deep Learning Framework for Super-Resolution based Inpainting	47
5.2.1	Overview . . . . .	48
5.3	Related Work . . . . .	50
5.3.1	Methodology . . . . .	51
5.3.2	Super-Resolution using CDCA . . . . .	51
5.3.3	Experimental Results . . . . .	53
5.3.4	Experiments . . . . .	53
5.3.5	Results . . . . .	55
5.3.6	Conclusion . . . . .	60
<b>6</b>	<b>3D-Coupled Deep Convolutional Auto-Encoder for Space-Time Super Resolution of Video</b>	<b>61</b>
6.1	Introduction . . . . .	62
6.2	Space-Time Super-Resolution . . . . .	65
6.3	Results . . . . .	73
6.3.1	Experiments . . . . .	73
6.3.2	Spatial SR . . . . .	73
6.3.3	Temporal SR . . . . .	74
6.3.4	Usecase: HEVC compatible Video Super Resolution Framework . . . . .	81
6.4	Conclusions . . . . .	85
<b>7</b>	<b>A Noise-Resilient Super-Resolution framework to boost OCR performance</b>	<b>87</b>
7.1	Introduction . . . . .	88
7.2	Methodology . . . . .	90

---

7.2.1	Deep-learning based text Image Super-Resolution . . . . .	90
7.2.2	Deep-learning based Noise-Resilient Text Image Super-Resolution . . .	93
7.2.3	Deep BLSTM recognition engine . . . . .	94
7.3	Experimental Results . . . . .	94
7.3.1	Datasets . . . . .	94
7.3.2	Results . . . . .	95
7.4	Conclusions . . . . .	98
<b>8</b>	<b>Conclusion</b>	<b>99</b>
8.1	Summary . . . . .	99
8.2	Contributions . . . . .	99
8.3	Future Work . . . . .	101
	<b>Bibliography</b>	<b>103</b>
	<b>Publications</b>	<b>117</b>
	<b>Biography</b>	<b>119</b>

# List of Figures

3.1	Block Diagram of Methodology . . . . .	19
3.2	Correct Frames . . . . .	22
3.3	Distorted Frames . . . . .	22
3.4	Accuracy with Different Dictionary Size . . . . .	23
4.1	Block Diagram of CAE . . . . .	29
4.2	Block Diagram of CDCA . . . . .	30
4.3	Image SR comparison (3x) on butterfly. (a) Original. (b) Bicubic. (c) LLE. (d) CDCA. . . . .	34
5.1	Block Diagram for noise resilient SR framework . . . . .	40
5.2	SSDA-CDCA Framework . . . . .	40
5.3	Plot shows PSNR as a function of noise variance . . . . .	44
5.4	Noise Resilient Image SR (2x) comparison on lamma. (a) LR image with Gaussian Noise ( $\sigma=30$ ). (b) Conventional. (c) Proposed SSDA-CDCA. (d) Original. . . . .	45
5.5	Block Diagram of Deep Learning Framework for Super-Resolution based Inpainting . . . . .	49
5.6	Visual comparison of different blind inpainting algorithms for Image.1 (upper one), Image.2 (middle one) and Image.3 (Lower one): (a) Input, (b) BiCNN[11], (c) RED30[68], (d) Proposed, (e) Ground truth. . . . .	55
5.7	Noise Resilient Image SR (2x) comparison on lamma. (a) LR image with Gaussian Noise ( $\sigma=30$ ). (b) Conventional. (c) Proposed CAE-CDCA. (d) Original. . . . .	56

5.8	Visual comparison of different blind inpainting algorithms for Images: (a) input, (b) Conventional1, (c) state-of-the-art [93], (d) proposed, (e) ground truth. . . . .	58
5.9	Noise Resilient Image SR (3x) comparison on baboon (from left to right) (i) LR image with Gaussian Noise ( $\sigma = 30$ ) and random noise (ii) Conventional1 (iii) Proposed (iv) Original. . . . .	58
5.10	Results Comparison of Proposed Robust Enhancement with Existing frameworks on noisy LR cartoon image with Gaussian Noise ( $\sigma = 30$ ) and salt-and-pepper noise (density=0.1): (a) Bicubic Interpolation (b) Conventional1 (c) Conventional2 (d) state-of-the-art [93] (e) Proposed (f) Original. . . . .	59
5.11	Results Comparison of Proposed Robust Enhancement with Existing frameworks on noisy LR zebra image with random noise: (a) Bicubic Interpolation (b) Conventional1 (c) Conventional2 (d) State-of-the-art [93] (e) Proposed (f) Original. . . . .	59
6.1	Block Diagram of 3D-CAE . . . . .	64
6.2	Block Diagram of 3D-CDCA . . . . .	65
6.3	SR comparison (3x) on frame of calendar sequence. (a) Original. (b) Bi-cubic. (c) Bayesian [4]. (d) Enhancer [33]. (e) VSRnet [50]. (f) Proposed. . . . .	70
6.4	SR comparison (3x) on frame of city sequence. (a) Bicubic. (b) Bayesian [4]. (c) Proposed. . . . .	70
6.5	Temporal SR (3x) on Flag sequence using (a) Tri-cubic interpolation. (b) [89]. (c) Proposed. . . . .	71
6.6	Temporal SR comparison (3x) between traditional interpolation approach (left column) and our approach (right column) . . . . .	71
6.7	Temporal SR (3x) on teddy sequence using (a) Tri-cubic interpolation. (b) [89]. (c) Proposed. . . . .	72
6.8	Block Diagram of H.264 Video Encoder extension for Space-time SR . . . . .	78
6.9	H.264/AVC decoder compatible Video Space-Time SR Framework . . . . .	78

---

6.10	Block Diagram of H.264/AVC Video Decoder extension for Space-Time SR . . .	79
6.11	H.264 Video Encoder compatible Video Space-Time SR Framework . . . . .	79
6.12	Block Diagram of H.264 Video Encoder extension for Space-time SR . . . . .	80
6.13	Comparison between performance and runtime of different algorithms . . . . .	84
7.1	Block Diagram of Proposed Architecture. A: Noisy LR text image; B: De-noised Text image; C: HR text image; D: Recognized Text output . . . . .	88
7.2	Block Diagram of CDCA . . . . .	90
7.3	Block Diagram for noise resilient SR framework . . . . .	92

# List of Tables

3.1	Specification of Dataset1 . . . . .	23
3.2	Accuracy on LIVE IQA Database . . . . .	23
3.3	Accuracy on Database1 . . . . .	23
3.4	Comparison of Accuracy, Precision, Recall, F-score of the proposed framework with state-of-the-art Blind/reference-less image spatial quality evaluator (BRISQUE) [73] on LIVE IQA Dataset . . . . .	23
3.5	Comparison of Accuracy, Precision, Recall, F-score of the proposed framework with state-of-the-art Blind/reference-less image spatial quality evaluator (BRISQUE) [73] on Dataset1 . . . . .	24
4.1	Image SR Results comparison on the Set5, Set14 and BSD200 Dataset . . . . .	34
5.1	Comparison of different methods with gaussian noise for 2x noise resilient super-resolution . . . . .	46
5.2	Comparison of different methods with Gaussian noise for 3x noise resilient super-resolution . . . . .	47
5.3	Blind Image Inpainting Results comparison on different Images and Dataset . . . . .	54
5.4	Comparison of different methods with gaussian noise for 3x noise resilient super-resolution . . . . .	58
5.5	Ablation study of different losses for proposed framework on various datasets (PSNR comparison) . . . . .	58

6.1	Results comparison of 3D-CDCA without Pre-trained weights and 3D-CDCA with Pre-trained weights for SR . . . . .	69
6.2	Average PSNR and SSIM comparison of different Video SR algorithms for different sequences . . . . .	69
6.3	Average PSNR and SSIM comparison of different proposed Video SR Frameworks . . . . .	83
6.4	H.264/AVC implementation results for Raw Video . . . . .	83
7.1	Comparison of proposed noise resilient 2X super-resolution framework on OCR accuracy and PSNR . . . . .	96
7.2	Comparison of proposed noise resilient 3X super-resolution framework on OCR accuracy and PSNR . . . . .	96
7.3	Comparison of proposed noise resilient 4X super-resolution framework on OCR accuracy and PSNR . . . . .	97