# Image Compression using Machine Learning: A Review

## Atul Kumar Yadav[1], Prof. Jitendra Mishra[2]

**[1]M. Tech. Scholar, Department of EC, PCST, Bhopal (India)**

**[2]Head & Professor, Department of EC, PCST, Bhopal (India)**

**Abstract-** *Apart from the existing technology on image compression represented by series of JPEG and MPEG, new technology such as neural networks, genetic algorithms, deep learning method and optimization techniques are being developed to explore the future of Image compression and coding. Successful applications of neural networks to vector quantization have now become well established, and other aspects of neural network involvement in this area are stepping up to play significant roles in assisting with those traditional technologies. This paper presents an extensive survey on the development of neural networks, deep neural network and optimization techniques for image compression.*

*Keywords:-* Image processing, Computer vision, Neural network, Convolution neural network.

## Introduction

Visual data constitutes most of the total information created and shared on the Web every day and it forms a bulk of the demand for storage and network bandwidth. It is customary to compress image data as much as possible as long as there is no perceptible loss in content. Image compression is one of the most fundamental techniques and commonly used applications in the image and video processing field. Earlier methods built a well-designed pipeline, and efforts were made to improve all modules of the pipeline by handcrafted tuning. Later, tremendous contributions were made, especially when data-driven methods revitalized the domain with their excellent modeling capacities and flexibility in incorporating newly designed modules and constraints. Despite great progress, a systematic benchmark and comprehensive analysis of end-to-end learned image compression methods are lacking. Image compression is a key technology in the development of various multimedia computer services and telecommunication applications such as teleconferencing, digital broadcast codec and video technology, etc. At present, the main core of image compression technology consists of three important processing stages: pixel transforms, quantization and entropy coding. In addition to these key techniques, new ideas are constantly appearing across different disciplines and new research fronts. As a model to simulate the learning function of human brains, neural networks have enjoyed widespread applications in telecommunication and computer science.

In recent decades, a variety of codecs have been developed to optimize the reconstruction quality with bitrate constraints. In the design of existing image compression frameworks, there are two basic principles. First, the image signal should be decorrelated, which is beneficial in improving the efficiency of entropy coding. Second, for lossy compression, the neglected information should have the least influence on the reconstruction quality, i.e., only the least important information for visual experience is discarded in the coding process.

With the emergence of deep neural networks, notably adopting the convolutional neural network (CNN) architecture, and its success in solving

complex image-based tasks, e.g. image classification, it is reasonable to expect this type of solutions to also bring benefits in terms of image coding. This may happen by using the deep neural networks as the codecs themselves or by improving standard image coding solutions, notably increasing their rate-distortion (RD) performance. The latter approach allows to benefit from compatibility with large standard image eco-systems such as the one associated to the JPEG standard. With the rapid development of deep learning, there have been many works exploring the potential of artificial neural networks to form an end-to-end optimized image compression framework. The development of these learning based methods has significant differences from traditional methods. For traditional methods, improved performance mainly comes from designing more complex tools for each component in the coding loop. Deeper analysis can be conducted on the input image, and more adaptive operations can be applied, resulting in more compact codes.

## II. Related Work

Deep learning and neural networks have achieved a great success in computer vision and image processing, especially in low-level vision problems such as image compression. In recent years deep learning has made it possible to design deep models for learning compact representations for image data**.**

[1] **J. Jiang,"Image compression with neural networks } A survey",0923-5965/99/$ - see front matter ( 1999) Elsevier Science B.V. All rights reserved,** This paper presents an extensive survey on the development of neural networks for image compression which covers three categories: direct image compression by neural networks; neural network implementation of existing techniques, and neural network based technology which provide improvement over traditional algorithms.

[2] **Mohammad Haris Baig,Vladlen Koltun,Lorenzo Torresani,"Learning to Inpaint for Image Compression",31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA,**We study the design of deep architectures for lossy image compression. We present two architectural recipes in the context of multi-stage progressive encoders and empirically demonstrate their importance on compression performance. Specifically, we show that: (a) predicting the original image data from residuals in a multi-stage progressive architecture facilitates learning and leads to improved performance at approximating the original content and (b) learning to inpaint (from neighboring image pixels) before performing compression reduces the amount of information that must be stored to achieve a high-quality approximation

[3] **Johannes Balle,Valero Laparra,Eero P. Simoncelli,"END-TO-END OPTIMIZED IMAGE COMPRESSION",Published as a conference paper at ICLR 2017,**They find that the optimized method generally exhibits better rate–distortion performance than the standard JPEG and JPEG 2000 compression methods. More importantly, we observe a dramatic improvement in visual quality for all images at all bit rates, which is supported by objective quality estimates using MS-SSIM.

[4] **Wenbin Yin,Xiaopeng Fan,Yunhui Shi,Wangmeng Zuo,"A Reference Resource Based End-to-End Image Compression Scheme",© Springer Nature Switzerland AG 2018 R. Hong et al. (Eds.): PCM 2018, LNCS 11164, pp. 534–544, 2018,** In this paper, we propose an end-to-end reference resource based image compression scheme to exploit the strong correlations with external similar images. In the proposed scheme, the side information is generated from highly correlated images in the reference resource. The features of side information can conceptually guide the compression process and assist the reconstruction process. The important map is employed to guide the allocation of local bit rate of the residual features. The proposed compression scheme is formulated as a rate distortion optimization problem in an end-to-end manner which is solved

by ADAM algorithm. Experimental results prove that the proposed compression framework greatly outperforms several image compression frameworks.

[5] **Thuong Nguyen Canh, Motong Xu, and Byeungwoo Jeon,"Rate-Distortion Optimized Quantization: A Deep Learning Approach",978-1-5386-5989-2/18/$31.00 ©2018 IEEE,**This paper is the first attempt to explore possibility of using deep learning in HEVC quantization. We set up a machine learning problem for RDOQ which predicts corresponding RDOQ quantized output upon receiving scalar quantization (SQ) result of a block as input. A residual learning framework is employed to predict the difference of SQ from RDOQ after further simplification that the residual becomes binary. By using a deep convolution neural network, the proposed deep learning based RDOQ (DL-RDOQ) is able to predict the optimal quantized levels without computing rate and distortion. Our experiments show potentially promising performance following RDOQ, especially at high bitrates.

[6] **Qing Li and Yang Chen,"Learning to Compress Using Deep AutoEncoder",2019 57th Annual Allerton Conference on Communication, Control, and Computing (Allerton) Allerton Park and Retreat Center Monticello, IL, USA, September 24-27, 2019,** A novel deep learning framework for lossy compression is proposed. The framework is based on Deep AutoEncoder (DAE) stacked of Restricted Boltzmann Machines (RBMs), which form Deep Belief Networks (DBNs). The proposed DAE compression scheme is one variant of the known fixed-distortion scheme, where the distortion is fixed and the compression rate is left to optimize. The fixed distortion is achieved by the DBN BlahutArimoto algorithm to approximate the N th-order rate distortion approximating posterior.

[7] **Yefei Wang, Dong Liu,Siwei Ma,Feng Wu,Wen Gao,,"Ensemble Learning-Based Rate-Distortion Optimization for End-to-End Image Compression",IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS FOR VIDEO TECHNOLOGY,2020IEEE Xplore,** We propose several methods to obtain multiple models. First, we adopt the boosting strategy to train multiple networks with diversity as an ensemble. Second, we train an ensemble of multiple probability distribution models to reduce the distribution gap for efficient entropy coding. Third, we present a geometric transform-based self-ensemble method. The multiple models can be regarded as the multiple coding modes, similar to those in non-deep video coding schemes. We further adopt block-level model/mode selection at the encoder side to pursue rate-distortion optimization, where we use hierarchical block partitioning to improve the adaptation ability.

[8] **M. Akin Yilmaz and A. Murat Tekalp,"End-to-End Rate-Distortion Optimization for Bi-Directional Learned Video Compression",2020 IEEE,**In this paper, we propose for the first time end-to-end optimization of a hierarchical, bi-directional motion compensated learned codec by accumulating cost function over fixed-size groups of pictures (GOP). Experimental results show that the rate-distortion performance of our proposed learned bi-directional GOP coder outperforms the state-of-the-art end-to-end optimized learned sequential compression as expected.

[9] **Paulo Eusébio,João Ascenso,Fernando Pereira,"Optimizing an Image Coding Framework with Deep Learning-based Pre- and Post-Processing",EUSIPCO 2020,**This paper targets improving image compression efficiency by designing and optimizing an image coding framework where a standard image codec, e.g. JPEG, is combined with deep neural network based pre- and post-processing. While the pre-processing CNN targets simplifying the image to make it more amenable to compression, notably involving its down-sampling, the postprocessing CNN targets enhancing the decoded image, also involving its up-sampling. To optimize the compression performance, the processing CNNs are trained involving a third CNN, so-called CNN-

FakeCodec, which targets modeling the image codec output, since the encoder-decoder pair is not differentiable, thus not allowing any training.

[10] **Yueyu Hu,Wenhan Yang,Zhan Ma,Jiaying Liu,"Learning End-to-End Lossy Image Compression: A Benchmark",IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE,9 Mar 2021,**In this paper, we first conduct a comprehensive literature survey of learned image compression methods. The literature is organized based on several aspects to jointly optimize the rate-distortion performance with a neural network, i.e., network architecture, entropy model and rate control. We describe milestones in cutting-edge learned image-compression methods, review a broad range of existing works, and provide insights into their historical development routes. With this survey, the main challenges of image compression methods are revealed, along with opportunities to address the related issues with recent advanced learning methods. This analysis provides an opportunity to take a further step towards higher-efficiency image compression.

### III. Deep Network-Based End-to-End Image Compression

In recent years, deep neural networks have achieved excellent results in image-related tasks such as image classification, super-resolution, image segmentation, object detection, and so on. Similarly, deep neural network-based image compression also gains attention. In general there are two approaches to image compression based on deep networks. First is to replace some module in the hybrid coding framework with trained network. Second is to build encoder/decoder solely upon trained network. The second approach is usually known as end-to-end (optimized) image compression. In end-to-end methods, the network has the original image as input, and outputs the reconstructed image. The RD (rate-distortion) cost is directly calculated and optimized by training network parameters. The recurrent neural network (RNN) for image compression, this method inputs an image into a shared RNN with multiple iterations, and takes the residual of each iteration as the input to the next iteration. Each iteration will output a part of the bit stream. The bit rate can be controlled by specifying the number of iterations. By accumulating the output of each iteration, the reconstructed image is obtained. This method does not jointly optimize the RD cost. Compared with RNN, convolutional neural network (CNN) has obtained more fruitful results in end-to-end image compression. Generalized divisive normalization (GDN) is used as the nonlinear activation function in the network. In order to solve the non-differentiable quantization problem, they proposed a method of adding uniform noise to approximate quantization, and directly estimated the bit rate, thus realizing the end to-end joint optimization of RD cost.

### IV. Conclusion

Data compression is a fundamental and well-studied problem in engineering, and is commonly formulated with the goal of designing codes for a given discrete data ensemble with minimal entropy. Recently, some end-to-end image compression methods have been proposed leading to a new direction of image compression. In this paper, we present a literature survey work on end-to-end reference resource based image compression scheme to exploit the strong correlations with external similar images.

**REFERENCES:**

[1] J. Jiang,"Image compression with neural networks } A survey",0923-5965/99/$ - see front matter ( 1999) Elsevier Science B.V. All rights reserved. PII: S 0 9 2 3 - 5 9 6 5 ( 9 8 ) 0 0 0 4 1 - 1.

[2] Mohammad Haris Baig,Vladlen Koltun,Lorenzo Torresani,"Learning to Inpaint for Image Compression",31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA.

[3] Johannes Balle,Valero Laparra,Eero P. Simoncelli,"END-TO-END OPTIMIZED IMAGE

COMPRESSION",Published as a conference paper at ICLR 2017.

[4] Wenbin Yin,Xiaopeng Fan,Yunhui Shi,Wangmeng Zuo,"A Reference Resource Based End-to-End Image Compression Scheme", Springer Nature Switzerland AG 2018 R. Hong et al. (Eds.): PCM 2018, LNCS 11164, pp. 534–544, 2018.

[5] Thuong Nguyen Canh, Motong Xu, and Byeungwoo Jeon,"Rate-Distortion Optimized Quantization: A Deep Learning Approach",978-1-5386-5989-2/18/\$31.00, 2018 IEEE.

[6] Qing Li and Yang Chen,"Learning to Compress Using Deep AutoEncoder",2019 57th Annual Allerton Conference on Communication, Control, and Computing (Allerton) Allerton Park and Retreat Center Monticello, IL, USA, September 24-27, 2019.

[7] Yefei Wang, Dong Liu,Siwei Ma,Feng Wu,Wen Gao,,"Ensemble Learning-Based Rate-Distortion Optimization for End-to-End Image Compression",IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS FOR VIDEO TECHNOLOGY,2020IEEE Xplore.

[8] M. Akin Yilmaz and A. Murat Tekalp,"End-to-End Rate-Distortion Optimization for Bi-Directional Learned Video Compression",2020 IEEE.

[9] Paulo Eusébio,João Ascenso,Fernando Pereira,"Optimizing an Image Coding Framework with Deep Learning-based Pre- and Post-Processing",EUSIPCO 2020,

[10] Yueyu Hu,Wenhan Yang,Zhan Ma,Jiaying Liu, "Learning End-to-End Lossy Image Compression: A Benchmark",IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE,9 Mar 2021.

**Atul Kumar Yadav** received his Bachelor`s degree in Electronics & communication, Babu Banarasi Das National Institute of Technology & Management Lucknow, (U.P.) in 2017. Currently he is pursuing Master of Technology Degree in Electronics & Comunication (Digital communication) from PCST, (RGPV), Bhopal, Madhya Pradesh India. His research area include wireless communication.

Mr. **Jitendra Mishra** he is Associate Professor and Head of the Department of Electronics and communication in PCST, Bhopal (RGPV). His received Master of Technology and Bachelor's of engineering respectively in Digital communication from BUIT, Bhopal and from RGPV, Bhopal. He has more than 12 years of teaching experience and publish 55+ papers in International journals, conferences etc. His areas of Interests are Antenna & Wave Propagation, Digital Signal Processing, Wireless Communication, Image Processing etc.