

# Low-rate Image Compression with Super-resolution Learning

Wei Gao<sup>1,2,\*</sup>, Lvfang Tao<sup>1,2</sup>, Linjie Zhou<sup>1,2</sup>, Dinghao Yang<sup>2</sup>, Xiaoyu Zhang<sup>1,2</sup>, Zixuan Guo<sup>1,2</sup>

<sup>1</sup>School of Electronic and Computer Engineering, Shenzhen Graduate School, Peking University, Shenzhen, China

<sup>2</sup>Peng Cheng Laboratory, Shenzhen, China

{gaowei262, ltao, ljzhou, zxy2019, gzx2019}@pku.edu.cn dinghowyang@gmail.com

## Abstract

In this paper, we propose an end-to-end learned image compression framework for low-rate scenarios. Based on variational autoencoder, our method features a pair of compact-resolution and super-resolution networks, a set of hyper and main codec networks, and a conditional context model. The learning process of this framework is facilitated with integrated non-local attention modules and phase congruency priors. Multiple models are obtained from training with different hyper-parameters, and are jointly used in the image-level model selection process for rate control, which ensures that the bit-rate constraint of the CLIC challenge is satisfied. Experimental results demonstrate that the proposed method can achieve an averaged multi-scale structural similarity (MS-SSIM) score of 0.9648 with bit-rate consumption of 0.1499 bits per pixel, which outperforms the BPG image coding method significantly.

## 1. Introduction

As the popularization of image and video applications, the volume of visual data becomes increasingly huge. Therefore, lossy image compression, especially with low bit-rate, becomes a challenging task. By consuming low bit-rate, image compression algorithm should provide much smaller perceived distortions.

\*Corresponding author. This work was supported by Natural Science Foundation of China under Grant 61801303, Guangdong Basic and Applied Basic Research Foundation under Grant 2019A1515012031, Shenzhen Science and Technology Plan Basic Research Project under Grant JCYJ20190808161805519, the Open Projects Program of National Laboratory of Pattern Recognition (NLPR) under Grant 202000045, the Open Project Program of the State Key Lab of CAD&CG (Grant No. A2009), Zhejiang University, the start-up fund of Shenzhen Graduate School of Peking University under Grant 2390101081, and CCF-Tencent Open Fund under Grant IAGR20190101.

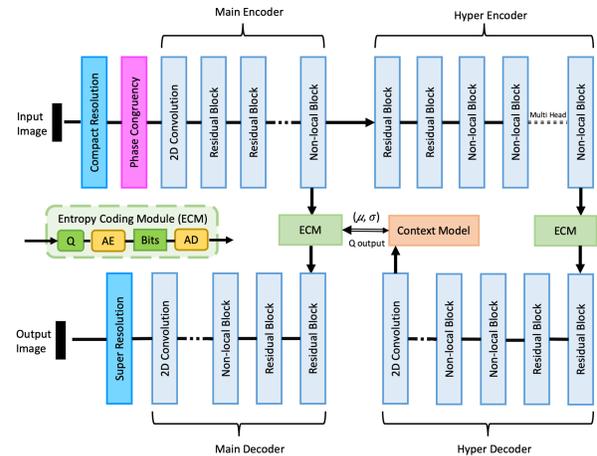


Figure 1. The overall structure of the proposed image compression method, which is mainly consist of variational autoencoder, compact resolution, super resolution, non-local and phase congruency modules. The symbols of Q, AE, AD represent quantizer, arithmetic encoder and arithmetic decoder, respectively.

Recently, with the development of neural networks, deep learning based image compression techniques have been proposed and achieve superior rate-distortion performance than traditional image codecs like JPEG [17], JPEG2000 [12], and HEVC-intra [13]. Autoencoders [1, 2, 3, 14, 9, 15] are widely used in end-to-end image compression, which include two major components, i.e., encoders and decoders. Encoders extract features from raw image to reduce the data redundancy, thereby expressing the image as a more compact feature representation. Decoders can utilize the feature expressions for image reconstruction in an inverse process. Ballé et al. [3] propose variational autoencoder (VAE) framework, where a hyper encoder is studied for better entropy modeling. Li et al. [10] add compact-resolution (CR) and super-resolution (SR) modules on traditional coding

methods as an multi-branch framework, including block-level adaptive scheme and frame-level scheme, to achieve bits saving. Jiang et al. [6] develop an end-to-end learning-based compression algorithm with compact Convolutional Neural Network (CNN) and reconstruction CNN, which shows superior results over traditional codecs. Since the perceptual quality for low-rate scenarios becomes much more important, learning-based image compression has not been fully investigated.

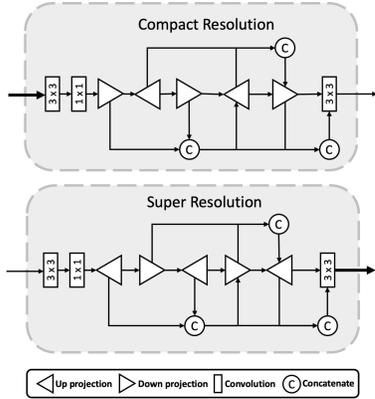


Figure 2. The structure of the compact-resolution and super-resolution modules.

In this work, we propose to improve the VAE architecture in [3] by introducing a paired CR and SR network. CR network is capable of acquiring the dense representation of images, while SR network can be jointly trained to restore the information loss during the down-scaling process in CR. Because the CR module greatly reduces the amount of pixels, the use of CR and SR is quite suitable for low-rate image compression. In addition, by using phase congruency [8] for texture evaluation, the structures of images can be better exploited for feature representation and learning in the network. The adopted non-local [18] module can capture long-range dependencies by attention mechanism to obtain global information. Finally, we also propose an iterative degradation-based rate control algorithm to balance rate and distortion trade-off on a per-image basis.

## 2. Proposed Method

### 2.1. Paired Compact-Resolution (CR) and Super-Resolution (SR) Networks

For further effective bit-rate reduction, the compact-resolution and super-resolution are employed at the beginning and end of the overall process, respectively. The original image is first down-sampled to low resolution and finally up-sampled to the original size. We would like to construct a CNN to get a better representation of an image which preserves more informative content after the down-

sampling process. The output down-sampled image is then encoded and decoded through an image compression network. Later, super-resolution network can generate the full resolution image.

The SR network used in this framework is based on deep back-projection network (DBPN) [5], which is one of the state-of-the-art single image SR methods. For light-weight model parameters and memory usage, we would like to reduce the number of network layers. More specifically, we use 3 up and down sampling units instead of 7 in the original DBPN. Additionally, the number of feature maps is also greatly cut down to simplify the network structure. For boosting the SR performance to restore the information degraded in CR, inspired by VAE [3], we intuitively construct a CR network which is symmetrical to the SR network. The structure of CR and SR network is depicted as Figure 2.

The CR network is composed of three parts, including initial feature extraction, up and down sampling projection, and final reconstruction. Initial feature extraction stage contains two convolutional layers to obtain the primary features. A convolutional layer with  $3 \times 3$  filter size firstly generates  $H_h \times W_h \times C_1$  feature map from  $H_h \times W_h \times C$  image, and then the  $1 \times 1$  filter is used to reduce the feature dimension from  $H_h \times W_h \times C_1$  to  $H_h \times W_h \times C_2$  ( $C_2 < C_1$ ). Up and down sampling projection stage includes 3 down-sampling units and 2 up-sampling units. Similar with [5], each unit is constructed as back-projection form based on residual learning. All the previous outputs are concatenated as the input to the next unit. The down-sampling unit outputs  $H_l \times W_l \times C_2$  feature from  $H_h \times W_h \times (C_2 \times n)$  feature, and the up-sampling unit outputs  $H_h \times W_h \times C_2$  feature from  $H_l \times W_l \times (C_2 \times n)$  feature, where  $n$  represents the number of concatenated features. A convolutional layer with  $3 \times 3$  filter size is used to reconstruct the final down-scaled image from  $H_l \times W_l \times C_2$  to  $H_l \times W_l \times C$  image.

The SR network is also composed of the same three parts as the CR network. The up and down sampling projection stage of SR network is symmetric with the CR network. It finally super-resolves the decoded image of the size  $H_l \times W_l \times C$  to its original size  $H_h \times W_h \times C$ .

The CR network and SR network are pre-trained before the training of the main compression network. The SR network is trained by minimizing loss function  $L_{sr}$ ,

$$L_{sr} = \|f_{sr}(g(x)) - x\|_2^2 \quad (1)$$

where  $x$  is the input image,  $f_{sr}$  represents the SR network,  $g$  is the bicubic interpolation.

Since the target image of CR is difficult to achieve, we use the SR network to assist the training of CR network. The training method in [10] is adopted, where the trained SR network is mounted after the CR network, and the weights of SR network are fixed. The loss function  $L_{cr}$

is defined as

$$L_{cr} = \|f_{sr}(f_{cr}(x)) - x\|_2^2 + \lambda \|f_{cr}(x) - g(x)\|_2^2 \quad (2)$$

where  $f_{cr}$  represents the CR network,  $\lambda$  is the parameter balancing between visual quality and the amount of contained information.

## 2.2. Non-Local Module

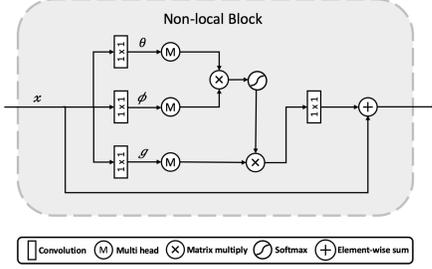


Figure 3. The structure of the adopted non-local mechanism.

Moreover, we would like to take advantage of the non-local mechanism as an enhanced attention method to better perceive image features adaptively, which computes a weighted mean of every pixel. In Figure 3, the structure of the proposed non-local module is illustrated, the input feature map  $x$  is processed into three flow  $\theta(x)$ ,  $\phi(x)$ ,  $g(x)$ , where  $\theta$ ,  $\phi$ ,  $g$  are implemented by  $1 \times 1$  convolution. Inspired by [16], we apply multi-head mechanism to learn different representation from different subspaces jointly, which are created by channel splitting. The matrix multiplication outputs of different subspaces are concatenated to aggregate information, then activated by softmax to obtain an attention mask. Furthermore, we add the global attention output with input  $x$  to get more abundant feature.

In view of the remarkable memory consumption, we only apply the non-local module in the high dimension processing, i.e., the highest layer of main encoder and decoder, all layers of hyper encoder and decoder.

## 2.3. Phase Congruency

In this paper, phase congruency (PC) [8] is employed to represent sharp transitions of image, which can evaluate the textures effectively. The PC of 2D image  $p$  with scale  $s$  and orientation  $r$  can be calculated as

$$PC = \frac{\sum_r \sum_s M_r(p) [L_{sr}(p) \Delta \Theta_{sr}(p) - N_r]}{\sum_r \sum_s L_{sr}(p) + \xi} \quad (3)$$

where  $M_r(p)$  is a metric for frequency spread, and  $L_{sr}(p)$  and  $\Delta \Theta_{sr}(p)$  are amplitude and phase deviation of  $p$ , respectively.  $N_r$  is a quantity used to reduce noise effect, while the symbol of  $\lfloor \rfloor$  means that the enclosed quantity

equals itself if the value is positive, otherwise equals zero.  $\xi$  is used for avoiding zero-division.

The PC image is down-sampled to the same size of main encoder layers via convolutions. Then multi-scale PC features are concatenated to the corresponding main encoder layers, which provides edge texture information.

## 2.4. Context and Entropy Modeling

The quantized outputs of main encoder and hyper encoder can be denoted as  $\hat{u}$  and  $\hat{v}$ , respectively. To predict the probability of  $\hat{u}$ , a context model is employed to obtain its distribution. The context model with 3D masked convolution network can predict the mean and standard deviation of  $\hat{u}$  with lower computational cost [11].

We can use two models to estimate the density of  $\hat{u}$  and  $\hat{v}$ . As [11], Gaussian distribution can be adopted to model the density of  $\hat{u}$ ,

$$p_{\hat{u}}(\hat{u}_i|\hat{v}) = \prod_i \left( \mathcal{N}(\mu_i, \sigma_i^2) * \mathcal{U}\left(-\frac{1}{2}, \frac{1}{2}\right) \right) (\hat{u}_i) \quad (4)$$

where  $\hat{u}_i$  represents the element of  $\hat{u}$ ,  $\mu_i$  and  $\sigma_i$  represent the mean and standard deviation of  $\hat{u}_i$ , respectively.

For  $\hat{v}$ , as [3], we can predict its probability without hyper-prior information by fully factorized density model as

$$p_{\hat{v}|\phi}(\hat{v}|\phi) = \prod_i \left( p_{v_i|\phi^{(i)}}(\phi^{(i)}) * \mathcal{U}\left(-\frac{1}{2}, \frac{1}{2}\right) \right) (\hat{v}_i) \quad (5)$$

where  $\hat{u}_i$  represents the element of  $\hat{u}$ , and  $\phi^{(i)}$  represents the parameters of univariate distribution  $p_{v_i|\phi^{(i)}}$ .

The bit rate of  $\hat{u}$  and  $\hat{v}$  can be calculated individually by

$$R_{\hat{u}} = - \sum_i \log_2(p_{\hat{u}}(\hat{u}_i|\hat{u}_1, \dots, \hat{u}_{i-1}, \hat{v})) \quad (6)$$

$$R_{\hat{v}} = - \sum_i \log_2(p_{\hat{v}_i|\phi^{(i)}}(\hat{v}_i|\phi^{(i)})) \quad (7)$$

where  $R_{\hat{u}}$  and  $R_{\hat{v}}$  denote the rates of  $\hat{u}$  and  $\hat{v}$ , respectively.

## 2.5. Loss Function

The loss function is designed as Eq. (8) so that the joint training of our learned model can be considered as a process of rate-distortion optimization.  $D_w$  is a weighted mixed distortion criterion, which is devised as Eq. (9) by combining mean squared error (MSE) and multi-scale structural similarity (MS-SSIM) [19] score,

$$L = D_w + R_{\hat{u}} + R_{\hat{v}} \quad (8)$$

$$D_w = \lambda_1 \times \|x - \hat{x}\|_2^2 + \lambda_2 \times (1 - \text{MS-SSIM}) \quad (9)$$

where  $\lambda_1$  and  $\lambda_2$  denote weighting coefficients,  $x$  and  $\hat{x}$  denote original and compressed images, respectively. By varying  $\lambda_1$ ,  $\lambda_2$ , models with different average bit-rates can be obtained through multiple training.

## 2.6. Image-Level Rate Control

In this work, three models will be trained to have different compression ratios, and are all employed to encode the validation image dataset. Therefore, we need to implement a rate control algorithm, which is responsible for choosing the most proper model for each individual image to be compressed. The algorithm should also ensure that the average bit-rate of all compressed images is below than, yet enough close to the target bit-rate [4, 11, 20].

Initially, we assume all the images are compressed with the highest-available quality. Then, an iterative degradation process is introduced to avoid the global bit budget being exceeded. The degradation cost defined in Eq. (10) is used to guide the degradation process,

$$Cost(i, j) = \begin{cases} \frac{Q(i, j) - Q(i, j+1)}{S(i, j) - S(i, j+1)}, & j < M \\ +\infty & j = M \end{cases} \quad (10)$$

where  $i$  and  $j$  denote the image index ( $i = 0, 1, \dots, N$ ) and its quality level ( $j = 0, 1, \dots, M$ , the lower the better).  $Q(i, j)$  denotes quality measure of the  $i$ -th image at the  $j$ -th quality level, and is calculated by MS-SSIM [19] which is reported to have better correlations with human perceptual experience.  $S(i, j)$  denotes bitstream file size of the  $i$ -th image at the  $j$ -th quality level. When a compressed image is already at the worst quality level, no further degradation could happen, hence  $Cost(i, M)$  is set to positive infinity.

For each time, one image with the lowest degradation cost among all images is selected to be downgraded to a lower quality level. The degradation process is iteratively executed until global bit constraint is satisfied.

## 3. Experimental Results

We use all images in CLIC 2020 dataset for training, including both professional and mobile sub-sets. After randomly resizing, the input image is cropped into multiple  $192 \times 192$  patches to train the networks.  $H_h = 192, W_h = 192$  is used for the input image of CR and the output image of SR, while  $H_l = 96, W_l = 96$  is used for the output image of CR and the input image of SR. We choose the scale of 2 for CR and SR, and set the parameter  $C_1 = 32, C_2 = 16$ . We set the parameter  $\lambda = 0.7$ , which shows relatively better performance in [10]. Two NVIDIA Tesla V100 GPUs are used during training and validation phase. The Adam [7] optimizer is adopted in the experiments. The initial learning rate is set to  $10^{-4}$ , and the batch size is 64.

In Figure 4, we compare rate-distortion results of our model on CLIC 2020 validation dataset with BPG (Better Portable Graphics) which is a state-of-the-art engineered image codec. Obviously, the proposed method can outperform BPG 4:2:0 within the low-rate range. For BPG codec, a QP range from 37 to 41 are used to generate bitstream and

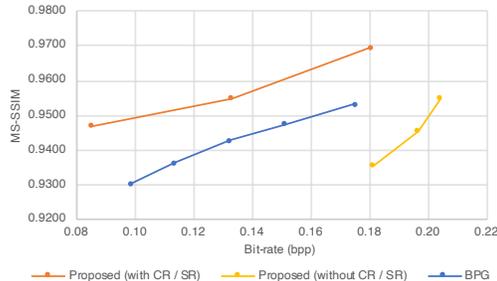


Figure 4. Rate-distortion performance of different models on CLIC 2020 validation dataset.

Table 1. Evaluation results on CLIC 2020 validation dataset.

Method	Bit-rate (bpp)	MS-SSIM			
		Mean	Max	Min	SD
Ours	0.1499	<b>0.9648</b>	<b>0.9852</b>	<b>0.9193</b>	<b>0.0142</b>
BPG	0.1498	0.9519	0.9796	0.8968	0.0164

decoded images. For fairness, we apply the same image-level rate control algorithm discussed above to these two methods, and then we can evaluate their overall performance under the same bit-rate constraint, which should be lower than 0.15 bits per pixel. The results on CLIC 2020 validation dataset are listed in Table 1, from which we can see that our proposed method outperforms BPG by higher average, maximum and minimum MS-SSIM scores. Meanwhile, our method maintains a less severe MS-SSIM fluctuation (SD: standard deviation) across 102 validation images.

To validate the effectiveness of the proposed super-resolution based method, we conduct more experiments by training and testing the proposed method with original  $192 \times 192$  image patches with removal of CR/SR models. Results of this ablation experiment are also depicted in Figure 4, from which we can find that the adoption of the paired compact-resolution and super-resolution networks shall account for the performance gain. Additionally, the performance of proposed method can be further improved if a larger image dataset could be trained and tested.

## 4. Conclusion

In this paper, we propose a novel learned image compression framework based on super-resolution learning. The use of paired compact-resolution (CR) and super-resolution (SR) networks in proposed framework shall be highlighted. Besides, efforts such as designing efficient non-local attention modules and providing phase congruency are also made to facilitate training convergence. From the ablation experiment, it can be seen that the adoption of proposed paired CR and SR networks can be beneficial for learning-based low-rate image compression tasks.

## References

- [1] Eirikur Agustsson, Fabian Mentzer, Michael Tschannen, Lukas Cavigelli, Radu Timofte, Luca Benini, and Luc V Gool. Soft-to-hard vector quantization for end-to-end learning compressible representations. In *Advances in Neural Information Processing Systems*, pages 1141–1151, 2017. [1](#)
- [2] Johannes Ballé, Valero Laparra, and Eero P Simoncelli. End-to-end optimized image compression. *arXiv preprint arXiv:1611.01704*, 2016. [1](#)
- [3] Johannes Ballé, David Minnen, Saurabh Singh, Sung Jin Hwang, and Nick Johnston. Variational image compression with a scale hyperprior. *arXiv preprint arXiv:1802.01436*, 2018. [1](#), [2](#), [3](#)
- [4] Zhengxue Cheng, Heming Sun, Masaru Takeuchi, and Jiro Katto. Deep residual learning for image compression. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2019. [4](#)
- [5] Muhammad Haris, Gregory Shakhnarovich, and Norimichi Ukita. Deep back-projection networks for super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1664–1673, 2018. [2](#)
- [6] Feng Jiang, Wen Tao, Shaohui Liu, Jie Ren, Xun Guo, and Debin Zhao. An end-to-end compression framework based on convolutional neural networks. *IEEE Transactions on Circuits and Systems for Video Technology*, 28(10):3007–3018, 2017. [2](#)
- [7] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. [4](#)
- [8] Peter Kovési. Phase congruency: A low-level image invariant. *Psychological research*, 64(2):136–148, 2000. [2](#), [3](#)
- [9] Mu Li, Wangmeng Zuo, Shuhang Gu, Debin Zhao, and David Zhang. Learning convolutional networks for content-weighted image compression. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3214–3223, 2018. [1](#)
- [10] Yue Li, Dong Liu, Houqiang Li, Li Li, Zhu Li, and Feng Wu. Learning a convolutional neural network for image compact-resolution. *IEEE Transactions on Image Processing*, 28(3):1092–1107, 2018. [1](#), [2](#), [4](#)
- [11] Haojie Liu, Tong Chen, Qiu Shen, and Zhan Ma. Practical stacked non-local attention modules for image compression. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, June 2019. [3](#), [4](#)
- [12] Majid Rabbani and Rajan Joshi. An overview of the jpeg 2000 still image compression standard. *Signal processing: Image communication*, 17(1):3–48, 2002. [1](#)
- [13] Gary J Sullivan, Jens-Rainer Ohm, Woo-Jin Han, and Thomas Wiegand. Overview of the high efficiency video coding (hevc) standard. *IEEE Transactions on circuits and systems for video technology*, 22(12):1649–1668, 2012. [1](#)
- [14] Lucas Theis, Wenzhe Shi, Andrew Cunningham, and Ferenc Huszár. Lossy image compression with compressive autoencoders. *arXiv preprint arXiv:1703.00395*, 2017. [1](#)
- [15] Robert Torfason, Fabian Mentzer, Eirikur Agustsson, Michael Tschannen, Radu Timofte, and Luc Van Gool. Towards image understanding from deep compression without decoding. *arXiv preprint arXiv:1803.06131*, 2018. [1](#)
- [16] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. [3](#)
- [17] Gregory K. Wallace. The jpeg still picture compression standard. *Communications of the Acm*, 38(1):xviii–xxxiv, 1992. [1](#)
- [18] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7794–7803, 2018. [2](#)
- [19] Zhou Wang, Eero P Simoncelli, and Alan C Bovik. Multi-scale structural similarity for image quality assessment. In *The Thirty-Seventh Asilomar Conference on Signals, Systems & Computers, 2003*, volume 2, pages 1398–1402. Ieee, 2003. [3](#), [4](#)
- [20] Lei Zhou, Zhenhong Sun, Xiangji Wu, and Junmin Wu. End-to-end optimized image compression with attention mechanism. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2019. [4](#)