# Extreme Image Coding via Multiscale Autoencoders with Generative Adversarial Optimization

Chao Huang, Haojie Liu, Tong Chen, Qiu Shen, and Zhan Ma*

Vision Lab, Nanjing University

*Abstract*—We propose a MultiScale AutoEncoder (MSAE) based extreme image coding/compression framework to offer visually pleasing reconstruction at a very low bitrate. Our method leverages the "priors" at different resolution scale to improve the compression efficiency, and also employs the generative adversarial network (GAN) with multiscale discriminators to perform the end-to-end trainable rate-distortion optimization. We compare the perceptual quality of our reconstructions with traditional compression algorithms using High-Efficiency Video Coding (HEVC) based Intra Profile and JPEG2000 on the public *Cityscapes*, *ADE20K* and *Kodak* datasets, demonstrating the significant subjective quality improvement. However, objective measurements, such as PSNR, SSIM, etc, are often deteriorated by applying the generative adversarial optimization.

## I. Introduction

Images that capture vivid scenes and events are stored and shared extensively every day. Thus image compression plays a vital role to ensure the efficient storage and sharing at the entire Internet scale. Traditional image compression methods such as JPEG, JPEG2000, HEVC based BPG, as well as recent deep neural network (DNN) based image compression methods [1]–[4] have presented significant advances in compression efficiency. Typically, these DNN-based schemes exhibit better visual quality than the traditional methodologies, at the same bit rate [5]. However, both of them fail to represent images efficiently with pleasant reconstruction quality at very low bitrates (e.g., targeting for $< 0.05$ bits per pixel (bpp)) [6].

This is mainly due to the reason that visual sensitive information (i.e., perceptual significance) can not be well preserved using conventional quality optimization criteria, such as peak signal-to-noise ratio (PSNR) and multiscale structural similarity (MS-SSIM) [8], at such extreme compression scenario. Recent explorations have shown that adversarial loss could be a tentative solution to capture global semantic information and local texture, yielding appealing reconstructions [6], [9]. Thus, Agustsson *et al.* [6] developed a GAN-based extreme image compression framework with bitrate below 0.1 bpp, resulting in the noticeable subjective quality improvement compared with the JPEG2000 [10] and BPG [11]. However, it had limitations by adopting a purely GAN-based structure. First, it was difficult to ensure the generalization of GAN to capture a variety of distributions of different datasets. In the meantime, GAN sometimes would introduce unexpected textures because of the failure of discriminator [12].

In this work, we propose a MultiScale AutoEncoder (MSAE) based extreme image compression structure where

Corresponding Author: Z. Ma.

we employ a multiscale network shown in Fig.1(a) to generate spatial scalable bitstreams. To the best of our knowledge, most learning based compression methods [1]–[4], generate a single layer bitstream at its native spatial resolution, without utilizing mutual information from other spatial scales. Different from Scalable Auto-encoder [13] that iteratively codes the pixel-level errors at the same resolution, "priors" at different spatial resolution scale that well capture the local textures, are embedded as reference to help the coarse-to-fine reconstruction and compression in our MSAE framework. Generative Adversarial loss [14] is applied in different scales for end-to-end trainable rate-distortion optimization, so as to optimize the reconstruction quality subjectively by maintaining the global semantic structure for visual significance, at a very low bit rate budget. We have our method tested on *Cityscapes*, *ADE20K* and `Kodak` datasets, yielding significant perceptual quality margins over the existing JPEG2000 and BPG.

## II. MultiScale AutoEncoder with Generative Adversarial Optimization

Fig.1(a) presents the extreme image compression framework of MSAE with generative adversarial optimization. Let $X_k$ be the original image ($k$ is the size of the input). We downscale the $X_k$ to obtain two more inputs $X_{k/s}$ and $X_{k/(s*s)}$. $s$ denotes the downscaling factor, which is set by 2 in this paper. Let $\mathbb{A}_i$ be the autoencoder network at scale $i$ ($i \in [k, k/2, k/4]$), and $\mathbb{U}$ denotes the upscaling operator. We then define the overall MSAE framework by

$$X'_{k/4} = \mathbb{A}_{k/4}(X_{k/4}), \tag{1}$$

$$X'_{k/2} = \mathbb{U}(X'_{k/4}) + \mathbb{A}_{k/2}(X_{k/2} - \mathbb{U}(X'_{k/4})), \tag{2}$$

$$X'_k = \mathbb{U}(X'_{k/2}) + \mathbb{A}_k(X_k - \mathbb{U}(X'_{k/2})). \tag{3}$$

Our proposed MSAE framework in (1), (2), and (3) has presented a coarse-to-fine reconstruction step by step. At the lowest scale $k/4$, the autoencoder $\mathbb{A}_{k/4}$ only takes $X_{k/4}$ as an input to derive the reconstructed image $X'_{k/4}$, yielding the *coarsest* representation of original $X_k$. Then $X'_{k/4}$, as the prior, is upscaled and aggregated with residuals at each scale to derive the final $X'_k$. Low resolution reconstructions are referred as "priors" to improve the overall rate-distortion performance. In addition, conditional GAN [15] is integrated into our MSAE system to do end-to-end training for visually appealing reconstruction, by enabling the multiscale discriminators for each input high-resolution images.
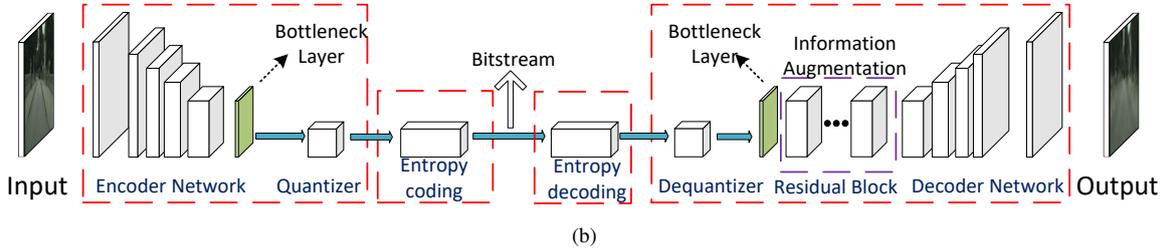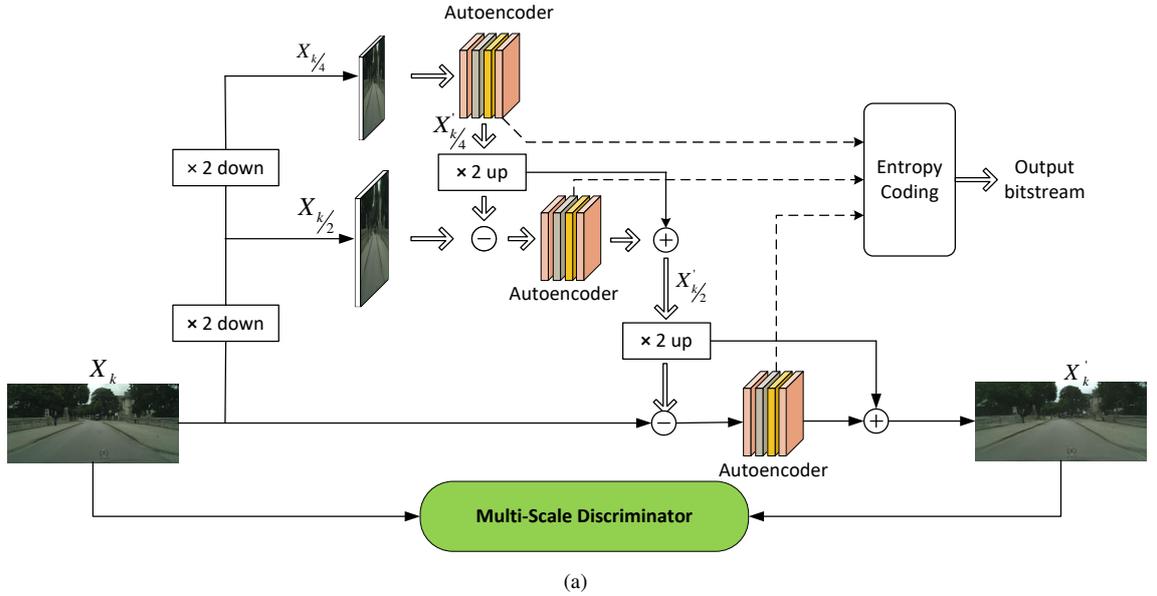
Fig. 1. Our extreme image compression framework via Multi-Scale AutoEncoder (MSAE) with GAN optimization. (a) overall structure, (b) autoencoder. The multi-scale distriminator in (a) contains three identical discriminators that are patch-based fully convolutional networks [7].The encoder network contains 1 convolutional layer with stride 1 and 4 convolutional layers with stride 2; all the residual blocks in information augmentation have the same convolutional kernel size 3 and stride 1; the decoder network is a mirror version of the encoder, which contains 4 transposed convolutional layers with stride 2 and 1 convolutional layer with stride 1. Entropy encoding and decoding denote the arithmetic encoding and decoding.

## A. AutoEncoder

The same autoencoder architecture is used in our MSAE framework at each scale. Except at the scale $k/4$ where the downscaled image $X_{k/4}$ serves as the input, residuals between upscaled priors and inputs (i.e., $X_{k/2} - \mathbb{U}(X'_{k/4})$ and $X_k - \mathbb{U}(X'_{k/2})$) at the same resolution, are fed into the autoencoder for compression. Using residuals, instead of default textures, generally boost the coding efficiency at the same bitrate budget due to better energy compaction and redundancy exploit.

Such autoencoder, shown in Fig.1(b), includes an encoder $E$ to encode the input $X$ to a set of feature maps (fMaps) $\omega$. Then the $\omega$ is passed to the quantizer $Q$ and will be quantized to a compressed representation $\hat{\omega} = Q(E(X))$. Specifically, the encoder $E$ first compresses the input with size of $W \times H \times C$ to feature maps with dimensions at $\frac{W}{16} \times \frac{H}{16} \times 480$. Usually, $W$ is for image width, $H$ is the height and $C$ is the number of color channels (e.g., $C = 3$ for RGB color space). The fMaps are then projected down to $\frac{W}{16} \times \frac{H}{16} \times C_{neck}$ at bottleneck layer prior to being quantized for $\hat{\omega}$. Note that $C_{neck}$ varies at different scale.

The decoder, denoted by the generator $G$, tries to reconstruct the image $X' = G(\hat{\omega})$ from the compressed representation $\hat{\omega}$. Within the decoder, *information augmentation module* with nine residual blocks [16] is aggregated to retrieve more information from the data to improve the reconstruction. Decoded fMaps will go through a mirror network of $E$ to obtain final reconstruction with dimensions at the same dimension, i.e., $W \times H \times C$, as the input image.

Note that the autoencoder is optimized using PSNR or MS-SSIM in default, often resulting in compression artifacts such as blocking, blurring and contouring effects at a low bitrate. To address this problem, we adopt adversarial loss [9] in training to reconstruct image $X'$ with visually pleasant quality.

## B. End-to-End Rate-Distortion Optimization

We adopt adversarial training in end-to-end optimization framework for extreme compression. This is mainly due to the reason that adversarial loss can address the blurring and contouring problems at a low bitrate level [9]. In the proposed framework, the decoder or generator $G$ is conditioned on the compressed representations and there is no necessity to add random noise for generator [15]. For discriminator $D$, we use the multiscale architecture following [14], which measures the

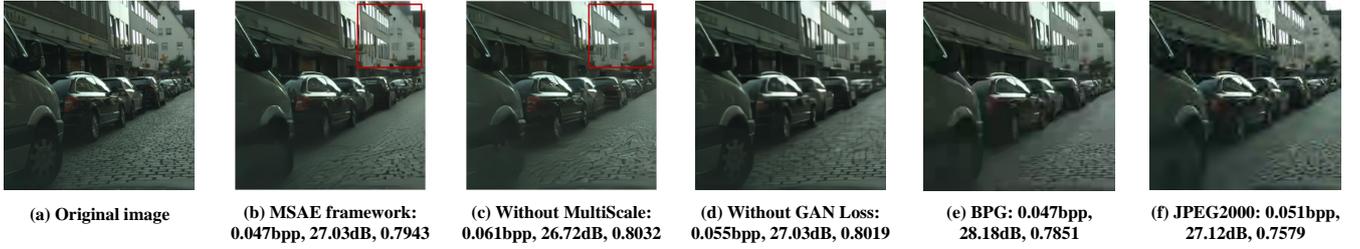| (a) Original image | (b) MSAE framework: 0.047bpp, 27.03dB, 0.7943 | (c) Without MultiScale: 0.061bpp, 26.72dB, 0.8032 | (d) Without GAN Loss: 0.055bpp, 27.03dB, 0.8019 | (e) BPG: 0.047bpp, 28.18dB, 0.7851 | (f) JPEG2000: 0.051bpp, 27.12dB, 0.7579 |

Fig. 2. Visual comparison on different architectures and loss functions, evaluated on a real-world image from Cityscapes dataset (values below each image are bitrate, PSNR and SSIM). In (c), we replace MSAE model with a single scale model. (d) MSAE with MS-SSIM loss optimization with color degradation. In (e) and (f), traditional codec frameworks make the reconstructed images have undesired blur and artifacts. Our complete model in (b) is able to produce better prediction. Compared with the result in (c), MSAE model can preserve local textures better (in red box).

divergence between real image and fake image generated by $G$ both globally and locally. Here we introduce a loss function that is closer to the perceptual similarity instead of relying on pixel-wise distortion [17], i.e.,

$$
L_f = \frac{\lambda}{W_{m,n}H_{m,n}} \sum_{x=1}^{W_{m,n}} \sum_{y=1}^{H_{m,n}} \left( \phi_{m,n}(Y_k)_{x,y} - \phi_{m,n}(Y_k')_{x,y} \right)^2,
\tag{4}
$$

with $Y_k = D(X_k)$ and $Y_k' = D(X_k')$. $\phi_{m,n}$ represents the feature map generated by the $n$-th convolution (with stride 2) of the $m$-th scale for the multiscale discriminator. $W_{m,n}$ and $H_{m,n}$ are the dimensional size of the respective feature maps. For the coefficient $\lambda$, we set it to 10.

The regular GAN [9] hypothesizes the discriminator as a classifier with the sigmoid cross entropy loss function, which may lead to gradient vanishing problem. In this paper, we use objective measures $f(y) = (y-1)^2$ and $g(y) = y^2$ developed for Least-Squares GAN [18], where $f$ and $g$ denote the scalar functions. It results in the generator loss as,

$$
L_G = \min_G f \left( D \left( G(\hat{\omega}_k) + \mathbb{U}(G(\hat{\omega}_{k/2}) + \mathbb{U}(G(\hat{\omega}_{k/4}))) \right) \right),
\tag{5}
$$

and the discriminator loss as:

$$
L_D = \min_D \left( f \left( D(X_k) \right) + g \left( D(X_k') \right) \right).
\tag{6}
$$

In order to backpropagate through the non-differentiable quantizer $Q$, we model the entropy rate following the [3] at bottleneck layer. We simply add uniform noise to ensure differentiability during training and replace it with ROUND($\cdot$) in inference. The entropy of $\hat{\omega}_i$ is evaluated using:

$$
H(\hat{\omega}) = - \sum_j \log_2(p_{\hat{\omega}_j|\psi^{(j)}}(\hat{\omega}_j \mid \psi^{(j)})),
\tag{7}
$$

where $\psi^{(j)}$ represents parameters of each univariate distribution $p_{\hat{\omega}_j}$. To balance the quality of the reconstruction and the bitrate, An entropy rate term need to be added to the training loss for optimal rate-distortion efficiency, i.e.,

$$
L_{RD} = \min_{d,H} \sum_{i \in [k, \frac{k}{s}, \frac{k}{s*s}]} \left( L_G + \alpha_i d(X_i, X_i') + L_f + \beta_i H(\hat{\omega}_i) \right).
\tag{8}
$$

As we can see, the rate-distortion trade-off is adjusted by setting the variations of $\alpha_i$ and $\beta_i$. Distortion, i.e., $d(X_i, X_i')$, is measured by the PSNR in this study, and the entropy of compressed representation, i.e., $H(\hat{\omega}_i)$, is used to approximate the encoding bitrate [3]. Such compound loss $L_{RD}$ is applied in a end-to-end trainable framework to achieve the optimal rate-distortion performance.

## III. EXPERIMENTAL STUDIES

**Datasets:** We use two public accessible datasets for training: Cityscapes [19] and ADE20K [20]. Cityscapes dataset contains 3475 images, each of them has the dimension of $2048 \times 1024 \times 3$ in RGB color space. During the training, we randomly select 2400 images for training and the rest for validation. These images are downscaled to $1024 \times 512 \times 3$ in our experiments to avoid GPU memory overflow in training. For the ADE20K dataset, we choose 4927 images. It is then segmented randomly to a training set and a validation set, with sizes from $256 \times 256 \times 3$ to $1024 \times 1024 \times 3$. For simplicity, we rescale all of them to $512 \times 512 \times 3$ for training and validation.

**Parameters:** We set $\beta_k = 100$ and $\beta_{k/4} = \beta_{k/2} = 1$. Meanwhile, $\alpha_k = 1$, $\alpha_{k/4} = \alpha_{k/2} = 100$ accordingly. The number of channels of the bottleneck layer varies between different scales. For scale $k/4$ and $k/2$, we set the $C_{neck} = 1$, while at scale $k$, $C_{neck} = 4$. This setting aims to provide sufficient "prior" information while consuming less bit overhead. Additionally, we use a learning rate of $2 \times 10^{-4}$ and the Adam optimizer for end-to-end learning.

**Performance Evaluation:** To evaluate the performance of our proposed MSAE based extreme image compression method, we compare our method with BPG and JPEG2000, as shown in Fig. 3, where both objective PSNR, SSIM and subjective snapshots of two samples are illustrated. For all the images we tested in datasets Cityscapes and ADE20K. We use basic arithmetic coding to generate actual bitstreams,

Fig. 3. Illustration of performance comparison for our proposed extreme image compression method versus BPG, JPEG2000 on ADE20K dataset with objective PSNR, SSIM and subjective snapshots.



Fig. 4. Visual comparison on Kodak image (left): example reconstructions by our MSAE framework (middle) and BPG (right). The respective bitrates, PSNR and SSIM are 0.070bpp, 19.78dB, 0.5386 vs 0.086bpp, 20.97dB, 0.6473.

and the bitrate is below 0.1 bpp. For quantitative evaluations, we compute the PSNR and SSIM between input $X$ and reconstruction $X^{'}$. But we have to mention that at such low bitrate, quantitative measurements such as PSNR or SSIM [21] become meaningless as they penalize changes in local structure rather than the preservation of the global semantics.

It is clear that our method has demonstrated noticeable perceptual quality margin over traditional JPEG2000 and BPG, even with tiny loss in objective metrics like PSNR and SSIM. Similar conclusion can be found in Fig.3. This also coincides similar observations that learning based compression can usually provide better visual quality, but worse PSNR [4]–[6].

## IV. CONCLUDING REMARKS

We have developed an extreme image compression framework via a multiscale autoencoder structure with embedded generative adversarial optimization for end-to-end training. Such multiscale authoencoder is fulfilled by downscaling the original image into various scales to capture the image statistics locally and globally. Each decoded representation at lower resolution scale is utilized as the priors for the efficient compression at higher scale. In addition to traditional pixel-wise distortion measurements (e.g., PSNR, SSIM), we have introduced the adversarial loss for pleasant image reconstruction at a very low bitrate (i.e., usually below 0.05 bpp), to preserve the image structure and global semantics. Experimental studies have demonstrated that our method has provided subjective quality improvement over existing JPEG2000 and BPG on public datasets, but objective evaluation suffers (e.g., PSNR or SSIM). This calls for the interesting studies on quality metrics to accurately capture the visual significance in semantic domain for ultra low bitrate scenario.

## REFERENCES

[1] George Toderici, Damien Vincent, Nick Johnston, Sung Jin Hwang, David Minnen, Joel Shor, and Michele Covell, "Full resolution image compression with recurrent neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 5306–5314.

[2] Oren Rippel and Lubomir Bourdev, "Real-time adaptive image compression," *arXiv preprint arXiv:1705.05823*, 2017.

[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.

[4] Haojie Liu, Tong Chen, Qiu Shen, Tao Yue, and Zhan Ma, "Deep image compression via end-to-end learning," in *2018 IEEE Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops 2018, Salt Lake City, UT, USA, June 18-22, 2018*, 2018, pp. 2575–2578.

[5] Johannes Ballé, Valero Laparra, and Eero P Simoncelli, "End-to-end optimized image compression," *arXiv preprint arXiv:1611.01704*, 2016.

[6] Eirikur Agustsson, Michael Tschannen, Fabian Mentzer, Radu Timofte, and Luc Van Gool, "Generative adversarial networks for extreme learned image compression," *arXiv preprint arXiv:1804.02958*, 2018.

[7] Jonathan Long, Evan Shelhamer, and Trevor Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3431–3440.

[8] Zhou Wang, Eero P Simoncelli, and Alan C Bovik, "Multiscale structural similarity for image quality assessment," in *The Thrity-Seventh Asilomar Conference on Signals, Systems & Computers, 2003*. Ieee, 2003, vol. 2, pp. 1398–1402.

[9] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, 2014, pp. 2672–2680.

[10] David Taubman and Michael Marcellin, *JPEG2000 image compression fundamentals, standards and practice: image compression fundamentals, standards and practice*, vol. 642, Springer Science & Business Media, 2012.

[11] Bellard, "Bpg image format," *https://bellard.org/bpg/*.

[12] Xintao Wang, Ke Yu, Chao Dong, and Chen Change Loy, "Recovering realistic texture in image super-resolution by deep spatial feature transform," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 606–615.

[13] Chuanmin Jia, Zhaoyi Liu, Yao Wang, Siwei Ma, and Wen Gao, "Layered image compression using scalable auto-encoder," in *2nd IEEE Conference on Multimedia Information Processing and Retrieval, MIPR 2019, San Jose, CA, USA, March 28-30, 2019*, 2019, pp. 431–436.

[14] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan Catanzaro, "High-resolution image synthesis and semantic manipulation with conditional gans," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8798–8807.

[15] Mehdi Mirza and Simon Osindero, "Conditional generative adversarial nets," *arXiv preprint arXiv:1411.1784*, 2014.

[16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[17] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al., "Photo-realistic single image super-resolution using a generative adversarial network," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2017, pp. 105–114.

[18] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley, "Least squares generative adversarial networks," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2794–2802.

[19] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele, "The cityscapes dataset for semantic urban scene understanding," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 3213–3223.

[20] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba, "Scene parsing through ade20k dataset," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2017, vol. 1, p. 4.

[21] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE transactions on image processing*, vol. 13, no. 4, pp. 600–612, 2004.