# Variable Rate Deep Image Compression with Modulated Autoencoder

Fei Yang, Luis Herranz, Joost van de Weijer, Jos A. Iglesias Guitin, Antonio M. Lpez, Mikhail G. Mozerov

*Abstract*—**Variable rate is a requirement for flexible and adaptable image and video compression. However, deep image compression methods (DIC) are optimized for a single fixed rate-distortion (R-D) tradeoff. While this can be addressed by training multiple models for different tradeoffs, the memory requirements increase proportionally to the number of models. Scaling the bottleneck representation of a shared autoencoder can provide variable rate compression with a single shared autoencoder. However, the R-D performance using this simple mechanism degrades in low bitrates, and also shrinks the effective range of bitrates. To address these limitations, we formulate the problem of variable R-D optimization for DIC, and propose modulated autoencoders (MAEs), where the representations of a shared autoencoder are adapted to the specific R-D tradeoff via a modulation network. Jointly training this modulated autoencoder and the modulation network provides an effective way to navigate the R-D operational curve. Our experiments show that the proposed method can achieve almost the same R-D performance of independent models with significantly fewer parameters.**

*Index Terms*—**Deep image compression, variable bitrate, autoencoder, modulated autoencoder**

## I. Introduction

**I**MAGE compression is a fundamental and well-studied problem in image processing and computer vision [2], [3], [4]. The goal is to design binary representations (i.e. bitstreams) with minimal entropy [5] that minimize the number of bits required to represent an image (i.e. bitrate) at a given level of fidelity (i.e. distortion) [6]. In many applications, communication networks or storage devices may impose a constraint on the maximum bitrate, which requires the image encoder to adapt to a given bitrate budget. In some applications this constraint may even change dynamically over time (e.g. video transmission). In all these cases, a bitrate control mechanism is required, and it is available in most traditional image and video compression codecs. In general, reducing the bitrate causes an increase in the distortion, i.e. there is a rate-distortion (R-D) tradeoff. This mechanism is typically based on scaling the latent representation prior to quantization to
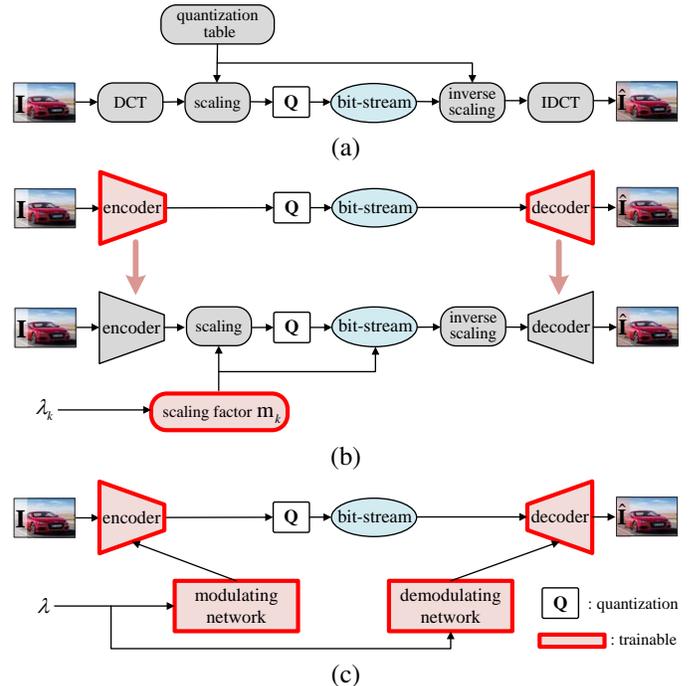
Fig. 1: Image compression and R-D control: (a) JPEG transform, (b) pre-trained autoencoder with bottleneck scaling [1], and (c) our proposed modulated autoencoder with joint training. Entropy coding/decoding is omitted for simplicity.

obtain finer or coarser quantizations, and then inverting the scaling at the decoder side (Fig. 1a).

Recent studies show that deep image compression (DIC) achieves comparable or even better results than classical image compression techniques [7], [8], [9], [10], [1], [11], [12], [13], [14], [15]. In this paradigm, the parameters of the encoder and decoder are learned from certain image data by jointly minimizing rate and distortion at a particular R-D tradeoff. However, variable bitrate requires an independent model for every R-D tradeoff. This is an obvious limitation, since it requires storing each model separately, resulting in large memory requirement.

To address this limitation, Theis et al. [1] use a single autoencoder whose bottleneck representation is scaled before quantization depending on the target bitrate (Fig. 1b). However, this approach only considers the importance of different channels from the bottleneck representation of learned autoencoders under R-D tradeoff constraint. In addition, the autoencoder is optimized for a single specific R-D tradeoff (typically high bitrate). These aspects lead to a drop in performance for low bitrates and a narrow effective range of

bitrates.

In order to tackle the limitations of multiple independent models and bottleneck scaling, we formulate the problem of variable R-D optimization for DIC, and propose the *modulated autoencoder* (MAE) framework, where the representations of a shared autoencoder at different layers are adapted to a specific R-D tradeoff via a modulating network. The modulating network is conditioned on the target R-D tradeoff, and synchronized with the actual tradeoff optimized to learn the parameters of the autoencoder and the modulating network. MAEs can achieve almost the same operational R-D points of independent models with much fewer overall parameters (i.e. just the shared autoencoder plus the small overhead of the modulating network). Multi-layer modulation does not suffer from the main limitations of bottleneck scaling, namely, drop in performance for low rates, and shrinkage of the effective range of bitrates.

## II. BACKGROUND

Almost all lossy image and video compression approaches follow the transform coding paradigm [16]. The basic structure is a transform $\mathbf{z} = f(\mathbf{x})$ that takes an input image $\mathbf{x} \in \mathbb{R}^N$ and obtains a transformed representation $\mathbf{z}$, followed by a quantizer $\mathbf{q} = Q(\mathbf{z})$ where $\mathbf{q} \in \mathbb{Z}^D$ is a discrete-valued vector. The decoder reverses the quantization (i.e. dequantizer $\hat{\mathbf{z}} = Q^{-1}(\mathbf{q})$) and the transform (i.e. inverse transform) as $\hat{\mathbf{x}} = g(\hat{\mathbf{z}})$ reconstructing the output image $\hat{\mathbf{x}} \in \mathbb{R}^N$. Before the transmission (or storage), the discrete-valued vector $\mathbf{q}$ is binarized and serialized into a *bitstream* $\mathbf{b}$. Entropy coding [17] is used to exploit the statistical redundancy in that bitstream and reduce its length.

In DIC, the handcrafted analysis and synthesis transforms are replaced by the encoder $\mathbf{z} = f(\mathbf{x}; \theta)$ and decoder $\hat{\mathbf{x}} = g(\hat{\mathbf{z}}; \phi)$ of a convolutional autoencoder, parametrized by $\theta$ and $\phi$. The fundamental difference is that the transforms are not designed but *learned* from training data. The model is typically trained by minimizing the optimization problem

$$\arg\min_{\theta,\phi} R(\mathbf{b}) + \lambda D(\mathbf{x}, \hat{\mathbf{x}}), \tag{1}$$

where $R(\mathbf{b})$ measures the rate of the bitstream $\mathbf{b}$ and $D(\hat{\mathbf{x}}, \mathbf{x})$ represents a distortion metric between $\mathbf{x}$ and $\hat{\mathbf{x}}$, and the Lagrange multiplier $\lambda$ controls the R-D tradeoff. Note that $\lambda$ is fixed in this case. The problem is solved using gradient descent and backpropagation [18].

To make the model differentiable, which is required to apply backpropagation, during training the quantizer is replaced by a differentiable proxy function [1], [7], [8]. Similarly, entropy coding is invertible, but it is necessary to compute the length of the bitstream $\mathbf{b}$. This is usually approximated by the entropy of the distribution of the quantized vector, $R(\mathbf{b}) \approx H[P_{\mathbf{q}}]$, which is a lower bound of the actual bitstream length.

In this paper, we will use scalar quantization by (element-wise) rounding to the nearest neighbor, i.e. $\mathbf{q} = \lfloor \mathbf{z} \rceil$, which will be replaced by additive uniform noise as proxy during training, i.e. $\tilde{\mathbf{z}} = \mathbf{z} + \Delta\mathbf{z}$, with $\Delta\mathbf{z} \sim \mathcal{U}\left(-\frac{1}{2}, \frac{1}{2}\right)$. There is no de-quantization in the decoder, and the reconstructed representation is simply $\hat{\mathbf{z}} = \mathbf{q}$. To estimate the entropy we

will use the entropy model described in [8] to approximate $P_{\mathbf{q}}$ by $p_{\tilde{\mathbf{z}}}(\tilde{\mathbf{z}})$. Finally, we will use mean squared error (MSE) as a distortion metric. With these particular choices, (1) becomes

$$\underset{\theta,\phi}{\arg\min}\, R(\tilde{\mathbf{z}}; \theta) + \lambda D(\mathbf{x}, \hat{\mathbf{x}}; \theta, \phi), \tag{2}$$

with

$$R(\tilde{\mathbf{z}}; \theta) = \mathbb{E}_{\mathbf{x} \sim p_{\mathbf{x}}, \Delta\mathbf{z} \sim \mathcal{U}}\left[-\log_2 p_{\tilde{\mathbf{z}}}(\tilde{\mathbf{z}})\right], \tag{3}$$

$$D(\mathbf{x}, \hat{\mathbf{x}}; \theta, \phi) = \mathbb{E}_{\mathbf{x} \sim p_{\mathbf{x}}, \Delta\mathbf{z} \sim \mathcal{U}}\left[\|\mathbf{x} - \hat{\mathbf{x}}\|^2\right]. \tag{4}$$

## III. MULTI-RATE DEEP IMAGE COMPRESSION WITH MODULATED AUTOENCODERS

### A. Problem definition

We are interested in DIC models that can operate satisfactorily on different R-D tradeoffs, and adapt to a specific one when required. Note that Eq.(2) optimizes rate and distortion for a fixed tradeoff $\lambda$. We extend that formulation to multiple R-D tradeoffs (i.e. $\lambda \in \Lambda = \{\lambda_1, \ldots, \lambda_M\}$) as the *multi-rate-distortion problem*

$$\underset{\theta,\phi}{\arg\min} \sum_{\lambda \in \Lambda} \left[R(\tilde{\mathbf{z}}; \theta, \lambda) + \lambda D(\mathbf{x}, \hat{\mathbf{x}}; \theta, \phi, \lambda)\right], \tag{5}$$

with

$$R(\tilde{\mathbf{z}}; \theta, \lambda) = \mathbb{E}_{\mathbf{x} \sim p_{\mathbf{x}}, \Delta\mathbf{z} \sim \mathcal{U}}\left[-\log_2 p_{\tilde{\mathbf{z}}}(\tilde{\mathbf{z}})\right], \tag{6}$$

$$D(\mathbf{x}, \hat{\mathbf{x}}; \theta, \phi, \lambda) = \mathbb{E}_{\mathbf{x} \sim p_{\mathbf{x}}, \Delta\mathbf{z} \sim \mathcal{U}}\left[\|\mathbf{x} - \hat{\mathbf{x}}\|^2\right], \tag{7}$$

where we are simplifying the notation by omitting features dependency on $\lambda$, i.e. $\tilde{\mathbf{z}} = \tilde{\mathbf{z}}(\lambda) = f(\mathbf{x}; \theta, \lambda)$ and $\hat{\mathbf{x}} = \hat{\mathbf{x}}(\lambda) = g(\tilde{\mathbf{z}}(\lambda); \phi, \lambda)$. This formulation can be easily extended to a continuous range of tradeoffs. Note also that these optimization problems assume that all R-D operational points are equally important. It could be possible to integrate an importance function $I(\lambda)$ to further give more importance to certain R-D operational points if required. We assume uniform importance (continuous or discrete) for simplicity.

### B. Bottleneck scaling

A possible way to make the encoder and decoder aware of $\lambda$ is simply scaling the latent representation in the bottleneck before quantization (implicitly scaling the quantization bin), and then inverting the scaling in the decoder. In this case, $\mathbf{q} = Q(\mathbf{z} \odot \mathbf{s}(\lambda))$ and $\hat{\mathbf{x}}(\lambda) = g(\hat{\mathbf{z}}(\lambda) \odot (\mathbf{1}/\mathbf{s}(\lambda)); \phi)$, where $\mathbf{s}(\lambda)$ is the specific scaling factor for the tradeoff $\lambda$. Conventional codecs use predefined tables for $\mathbf{s}(\lambda)$ (the descaling is often implicitly subsumed in the dequantization, e.g. JPEG), while other approaches [1] keep encoder and decoder fixed, optimized for a particular R-D tradeoff (Fig. 1(b)).

We observe several limitations in this approach: (1) scaling only the bottleneck features is not flexible enough to adapt to a large range of R-D tradeoffs, (2) using the inverse of the scaling factor in the decoder may also limit the flexibility of the adaptation mechanism, (3) optimizing the parameters of the autoencoder only for a single R-D tradeoff leads to suboptimal parameters for other distant tradeoffs, (4) training the autoencoder and the scaling factors separately may also be limiting. In order to overcome these limitations we propose the *modulated autoencoder (MAE)* framework.
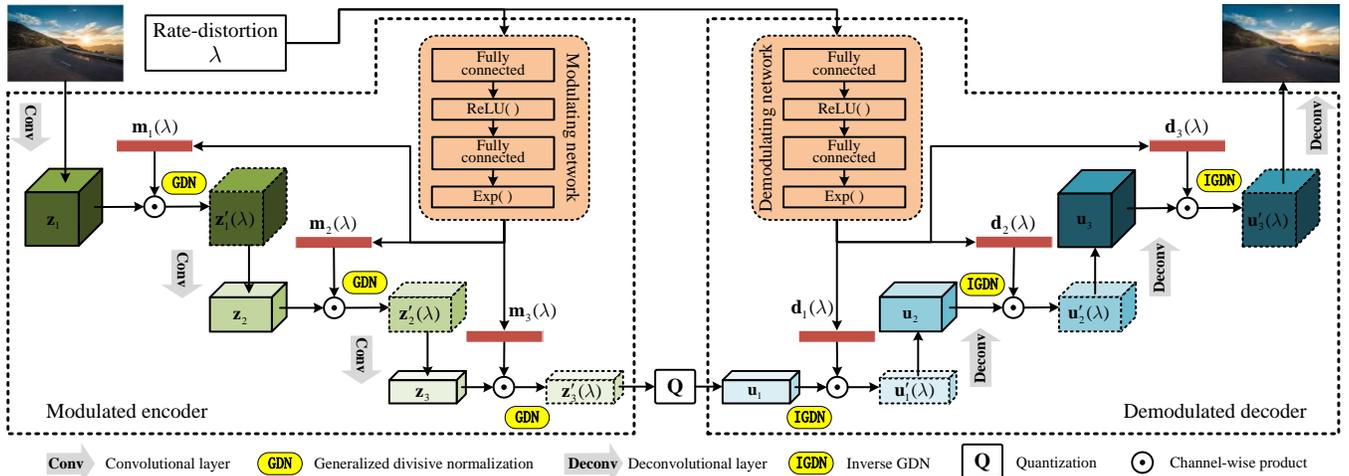
Fig. 2: Modulated autoencoder (MAE) architecture, combining modulating networks and shared autoencoder. The channel-wise product is performed before GDN in the encoder and after IGDN in the decoder.

### C. Modulated autoencoders

Variable rate is achieved in MAEs by modulating the internal representations in the encoder and the decoder (Fig. 2). Given a set of internal representations in the encoder $\mathbf{Z} = \{\mathbf{z}_1, \ldots, \mathbf{z}_K\}$ and in the decoder $\mathbf{U} = \{\mathbf{u}_1, \ldots, \mathbf{u}_L\}$, they are replaced by the corresponding modulated and demodulated versions $\mathbf{Z}' = \{\mathbf{z}_1 \odot \mathbf{m}_1(\lambda), \ldots, \mathbf{z}_K \odot \mathbf{m}_K(\lambda)\}$ and $\mathbf{U}' = \{\mathbf{u}_1 \odot \mathbf{d}_1(\lambda), \ldots, \mathbf{u}_L \odot \mathbf{d}_L(\lambda)\}$, where $\mathbf{m}(\lambda) = (\mathbf{m}_1(\lambda), \ldots, \mathbf{m}_K(\lambda))$ and $\mathbf{d}(\lambda) = (\mathbf{d}_1(\lambda), \ldots, \mathbf{d}_L(\lambda))$ are the modulating and demodulating functions.

Our MAE architecture extends the DIC architecture proposed in [8] which combines convolutional layers and GDN/IGDN layers [19]. In our experiments, we choose to modulate the outputs of the convolutional layers in the encoder and decoder, i.e. $\mathbf{Z}$ and $\mathbf{U}$, respectively.

The modulating function $\mathbf{m}(\lambda)$ for the encoder is learned by a *modulating network* as $\mathbf{m}(\lambda) = \mathbf{m}(\lambda; \vartheta)$ and the demodulating function $\mathbf{d}(\lambda)$ by the *demodulating network* as $\mathbf{d}(\lambda) = d(\lambda; \varphi)$. As a result, the encoder has learnable parameters $\{\theta, \vartheta\}$ and the decoder $\{\phi, \varphi\}$.

Finally, the optimization problem for the MAE is

$$\underset{\theta, \phi, \vartheta, \varphi}{\operatorname{argmin}} \sum_{\lambda \in \Lambda} \left[ R(\tilde{\mathbf{z}}; \theta, \vartheta, \lambda) + \lambda D(\mathbf{x}, \hat{\mathbf{x}}; \theta, \phi, \vartheta, \varphi, \lambda) \right], \quad (8)$$

which extends Eq.(5) with the modulating/demodulating networks and their corresponding parameters. All parameters are learned jointly using gradient descent and backpropagation.

This mechanism is more flexible than bottleneck scaling since it allows multi-level modulation, decouples encoder and decoder scaling and allows effective joint training of both autoencoder and modulating network, therefore optimizing jointly to all R-D tradeoffs of interest.

### D. Modulating and demodulating networks

The modulating network is a perceptron with two FC layers and ReLU [20] and exponential nonlinearities (Fig. 2). The exponential nonlinearity guarantees positive outputs which we found beneficial in training. The input is a scalar value $\lambda$ and the output is $\mathbf{m}(\lambda) = (\mathbf{m}_1(\lambda), \ldots, \mathbf{m}_K(\lambda))$. A small first

hidden layer allows learning a meaningful nonlinear function between tradeoffs and modulation vectors, which is more flexible than simple scaling factors and allows more expressive interpolation between tradeoffs. In practice, we use normalized tradeoffs as $\hat{\lambda}_k = \lambda_k / \max_{\lambda \in \Lambda}(\lambda)$. The demodulating network follows a similar architecture.

## IV. EXPERIMENTS

### A. Experimental setup

We evaluated MAE on the CLIC 2019 Professional dataset with 585 training images and 226 test images. In addition, we also test our models on Kodak dataset. Our implementation [1] is based on the autoencoder architecture of [8], which is augmented with modulation mechanisms and modulating networks (two FC layers, with 150 and $3 \times 192$ units respectively) for all the convolutional layers. We use MSE as distortion metric. The model is trained with crops of size $240 \times 240$ using Adam with a minibatch size of 8 and initial learning rates of 0.0004 and 0.002 for MAE and the entropy model, respectively. After 400k iterations, the learning rates are halved for another 150k iterations. We also tested MAEs with scale hyperpriors, as described in [21]. In our experiments, we consider seven ($\lambda \in [64, 128, 256, 512, 1024, 2048, 4096]$) and four ($\lambda \in [64, 256, 1024, 4096]$) R-D tradeoffs for the models without and with hyperprior, respectively. We consider two baselines:

**Independent models**. Each R-D operational point is obtained by training a new model with a different R-D tradeoff $\lambda$ in (2), requiring each model to be stored separately. This provides the optimal R-D performance, but also requires more memory to store all the models for different R-D tradeoffs.

**Bottleneck scaling**. The autoencoder is optimized for the highest R-D tradeoff in the range. Then it is frozen and the scaling parameters are learned for the other tradeoffs.
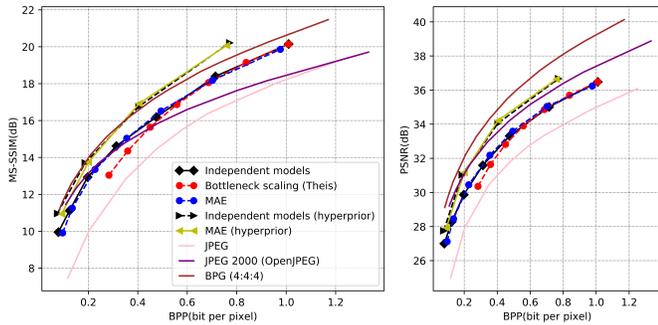
[1] https://github.com/FireFYF/modulatedautoencoder

Fig. 3: R-D curves for different methods on CLIC 2019 Professional test dataset.
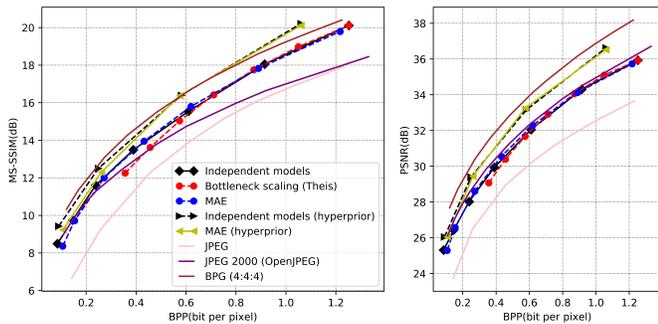


Fig. 4: R-D curves for different methods on Kodak dataset.

### B. Results

We use MS-SSIM (dB)[2] and PSNR (dB) to evaluate image distortion. Fig. 3 and Fig. 4 show the R-D operational curves for the proposed MAE and the two baselines on the CLIC 2019 Professional dataset and the Kodak dataset, respectively. In addition, it also includes JPEG [22], JPEG2000 [23] and BPG [24] coding methods to demonstrate their performance. For both of them, we can see that the best R-D performance is obtained by using independent models. Hyperprior models also have superior R-D performance. Bottleneck scaling is satisfactory for high bitrates, closer to the optimal R-D operational point of the autoencoder, but degrades for lower bitrates. Interestingly, bottleneck scaling cannot achieve as low bitrates as independent models since the autoencoder is optimized for a high bitrate. This can be observed in the R-D curve as a narrower range of bitrates. Note that our independent models results did not achieve the same performance as in [8] and [21] due to the different training datasets. The proposed MAEs can achieve an R-D performance very close to the corresponding independent models, demonstrating that multi-layer modulation with joint training is a more powerful mechanism to achieve effective variable rate compression.

The main advantage of bottleneck scaling and MAEs is that the autoencoder is shared, which results in much fewer parameters than independent models, which depend on the number of R-D tradeoffs (Table I). Both methods have a small overhead due to the modulating networks or the scaling factors (which is smaller in bottleneck scaling).

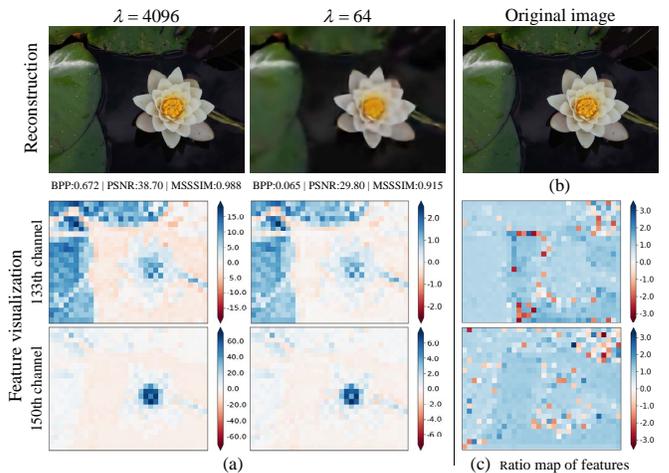[2]MS-SSIM (dB) is computed as $-10 \log_{10}(1 - \text{MS-SSIM})$ [11]



Fig. 5: Modulated feature maps: (a) reconstructed images for high ($\lambda = 4096$) and low ($\lambda = 64$) bitrates (first row), and the corresponding feature maps for two channels of the bottleneck before quantization (2nd and 3rd rows), (b) original image for comparison, and (c) element-wise ratio (logarithmic scale) between the feature maps at the two different tradeoffs.

TABLE I: Model size (millions of parameters)

| Architecture | w/o hyperprior [8] (seven R-D tradeoffs) | w/ hyperprior [21] (four R-D tradeoffs) |
|---|---|---|
| Independent models | 28.02 **M** | 40.53 **M** |
| Bottleneck scaling [1] | 4.00 **M** | - |
| Modulated AE (ours) | 4.06 **M** | 10.27 **M** |

In order to illustrate the differences between the bottleneck scaling and MAE bitrate adaptation mechanisms, we consider the image in Fig. 5b and the reconstructions for high and low bitrates in Fig. 5a. We show two of the 192 channels in the bottleneck feature before quantization (Fig. 5a), and observe that the maps for the two bitrates are similar but the range is higher for $\lambda = 4096$, so the quantization will be finer. This is also what we would expect in bottleneck scaling. However, a closer look highlights the difference between both methods. We also compute the element-wise ratio between the bottleneck features at $\lambda = 4096$ and $\lambda = 64$, and show the ratio image for the same channels of the example (Fig. 5c). We can see that the MAE learns to perform a more complex adaptation of the features beyond simple channel-wise bottleneck scaling since different areas of the ratio map show different values (the ratio map would be uniform in bottleneck scaling), which allows MAE to allocate bits more freely when optimizing for different R-D tradeoffs, especially for low bitrates.

### V. CONCLUSION

In this work, we introduce the modulated autoencoder, a novel variable rate deep image compression framework, based on multi-layer feature modulation and joint learning of autoencoder parameters. MAEs can realize variable bitrate image compression with a single model, while keeping the performance close to the upper bound of independent models that require significantly more memory. We show that MAE outperforms bottleneck scaling [1], especially for low bitrates.

## REFERENCES

[1] L. Theis, W. Shi, A. Cunningham, and F. Huszár, "Lossy image compression with compressive autoencoders," *arXiv preprint arXiv:1703.00395*, 2017.

[2] A. S. Lewis and G. Knowles, "Image compression using the 2-D wavelet transform," *IEEE Transactions on Image Processing*, vol. 1, no. 2, pp. 244–250, 1992.

[3] J. D. Villasenor, B. Belzer, and J. Liao, "Wavelet filter evaluation for image compression," *IEEE Transactions on Image Processing*, vol. 4, no. 8, pp. 1053–1060, 1995.

[4] D. Taubman and M. Marcellin, *JPEG2000 image compression fundamentals, standards and practice: image compression fundamentals, standards and practice*. Springer Science & Business Media, 2012, vol. 642.

[5] C. E. Shannon, "A mathematical theory of communication," *Bell system technical journal*, vol. 27, no. 3, pp. 379–423, 1948.

[6] A. K. Jain, *Fundamentals of digital image processing*. Englewood Cliffs, NJ: Prentice Hall, 1989.

[7] G. Toderici, S. M. O'Malley, S. J. Hwang, D. Vincent, D. Minnen, S. Baluja, M. Covell, and R. Sukthankar, "Variable rate image compression with recurrent neural networks," *arXiv preprint arXiv:1511.06085*, 2015.

[8] J. Ballé, V. Laparra, and E. P. Simoncelli, "End-to-end optimized image compression," *arXiv preprint arXiv:1611.01704*, 2016.

[9] K. Gregor, F. Besse, D. J. Rezende, I. Danihelka, and D. Wierstra, "Towards conceptual compression," in *Advances In Neural Information Processing Systems*, 2016, pp. 3549–3557.

[10] G. Toderici, D. Vincent, N. Johnston, S. Jin Hwang, D. Minnen, J. Shor, and M. Covell, "Full resolution image compression with recurrent neural networks," in *conference on Computer Vision and Pattern Recognition*, 2017, pp. 5306–5314.

[11] N. Johnston, D. Vincent, D. Minnen, M. Covell, S. Singh, T. Chinen, S. Jin Hwang, J. Shor, and G. Toderici, "Improved lossy image compression with priming and spatially adaptive bit rates for recurrent networks," in *conference on Computer Vision and Pattern Recognition*, 2018, pp. 4385–4393.

[12] D. Liu, H. Ma, Z. Xiong, and F. Wu, "Cnn-based dct-like transform for image compression," in *International Conference on Multimedia Modeling*. Springer, 2018, pp. 61–72.

[13] F. Mentzer, E. Agustsson, M. Tschannen, R. Timofte, and L. Van Gool, "Conditional probability models for deep image compression," in *conference on Computer Vision and Pattern Recognition*, 2018, pp. 4394–4402.

[14] D. Minnen, J. Ballé, and G. D. Toderici, "Joint autoregressive and hierarchical priors for learned image compression," in *Advances in Neural Information Processing Systems*, 2018, pp. 10 771–10 780.

[15] M. Li, W. Zuo, S. Gu, D. Zhao, and D. Zhang, "Learning convolutional networks for content-weighted image compression," in *conference on Computer Vision and Pattern Recognition*, 2018, pp. 3214–3223.

[16] V. K. Goyal, "Theoretical foundations of transform coding," *IEEE Signal Processing Magazine*, vol. 18, no. 5, pp. 9–21, 2001.

[17] P. A. Wintz, "Transform picture coding," *Proceedings of the IEEE*, vol. 60, no. 7, pp. 809–820, 1972.

[18] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning internal representations by error propagation," California Univ San Diego La Jolla Inst for Cognitive Science, Tech. Rep., 1985.

[19] J. Ballé, V. Laparra, and E. P. Simoncelli, "Density modeling of images using a generalized normalization transformation," *arXiv preprint arXiv:1511.06281*, 2015.

[20] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *International Conference on Machine Learning*, 2010, pp. 807–814.

[21] J. Ballé, D. Minnen, S. Singh, S. J. Hwang, and N. Johnston, "Variational image compression with a scale hyperprior," *arXiv preprint arXiv:1802.01436*, 2018.

[22] G. K. Wallace, "The jpeg still picture compression standard," *IEEE transactions on consumer electronics*, vol. 38, no. 1, pp. xviii–xxxiv, 1992.

[23] "JPEG2000 (OpenJPEG)," https://pypi.org/project/Glymur/, 2000, [Online; accessed 8-January-2020].

[24] F. Bellard, "BPG Image Format," http://bellard.org/bpg/, 2014, [Online; accessed 8-January-2020].