

DSSLIC: Deep Semantic Segmentation-based Layered Image Compression

Mohammad Akbari, Jie Liang
 School of Engineering Science, Simon Fraser University, Canada
 akbari@sfu.ca, jiel@sfu.ca *

Jingning Han
 Google Inc.
 jingning@google.com

Abstract

Deep learning has revolutionized many computer vision fields in the last few years, including learning-based image compression. In this paper, we propose a deep semantic segmentation-based layered image compression (DSSLIC) framework in which the semantic segmentation map of the input image is obtained and encoded as the base layer of the bit-stream. A compact representation of the input image is also generated and encoded as the first enhancement layer. The segmentation map and the compact version of the image are then employed to obtain a coarse reconstruction of the image. The residual between the input and the coarse reconstruction is additionally encoded as another enhancement layer. Experimental results show that the proposed framework outperforms the H.265/HEVC-based BPG and other codecs in both PSNR and MS-SSIM metrics across a wide range of bit rates in RGB domain. Besides, since semantic segmentation map is included in the bit-stream, the proposed scheme can facilitate many other tasks such as image search and object-based adaptive image compression¹.

1. Introduction

Since 2012, deep learning has revolutionized many computer vision fields such as image classification, object detection, and face recognition. In the last couple of years, it has also made some impacts to the well-studied topic of image compression, and in some cases has achieved better performance than JPEG2000 and the H.265/HEVC-based BPG image codec [1, 2, 3, 8, 10, 11, 15, 16, 17], making it a very promising tool for the next-generation image compression.

One advantage of deep learning is that it can extract much more accurate semantic segmentation map from a

given image than traditional methods [23]. Recently, it was further shown that deep learning can even synthesize a high-quality and high-resolution image using only a semantic segmentation map as input [19], thanks to the generative adversarial networks (GAN) [6]. This suggests the possibility of developing efficient image compression using deep learning-based semantic segmentation and the associated image synthesis.

GAN architecture is composed of two networks named discriminator and generator, which are trained at the same time [6]. The generator model $G(z)$ captures the data distribution by mapping the latent z to data space, while the discriminator model $D(x) \in [0, 1]$ estimates the probability that x is a real training sample or a fake sample synthesized by G . These two models compete in a two-player minimax game in which the objective function is to find a binary classifier D that discriminates the real data from the fake (generated) ones and simultaneously encourages G to fit the true data distribution. This goal is achieved by minimizing/maximizing the binary cross entropy:

$$\mathcal{L}_{GAN} = \mathbb{E}_{x \sim p_{data}(x)}[\log D(x)] + \mathbb{E}_{z \sim p_z(z)}[\log(1 - D(G(z)))] \quad (1)$$

where G tries to minimize this objective against D that tries to maximize it.

In this paper, we employ GAN to propose a deep semantic segmentation-based layered image compression (DSSLIC) framework as shown in Figure 1. In our approach, the semantic segmentation map of the input image is extracted by a deep learning network and losslessly encoded as the base layer of the bit-stream. Next, the input image and the segmentation map are used by another deep network to obtain a low-dimensional compact representation of the input, which is encoded into the bit-stream as the first enhancement layer. After that, the compact image and the segmentation map are used to obtain a coarse reconstruction of the image. The residual between the input and the coarse reconstruction is encoded as the second enhancement layer in the bit-stream. To improve the quality, the synthesized image from the segmentation map is designed to be a residual itself, which aims to compensate the difference between the upsampled version of the compact image

*This work is supported by Google Chrome University Research program and the Natural Sciences and Engineering Research Council (NSERC) of Canada under grant RGPIN312262 and RGPAS478109.

¹The source code of the paper: <https://github.com/makbari7/DSSLIC>

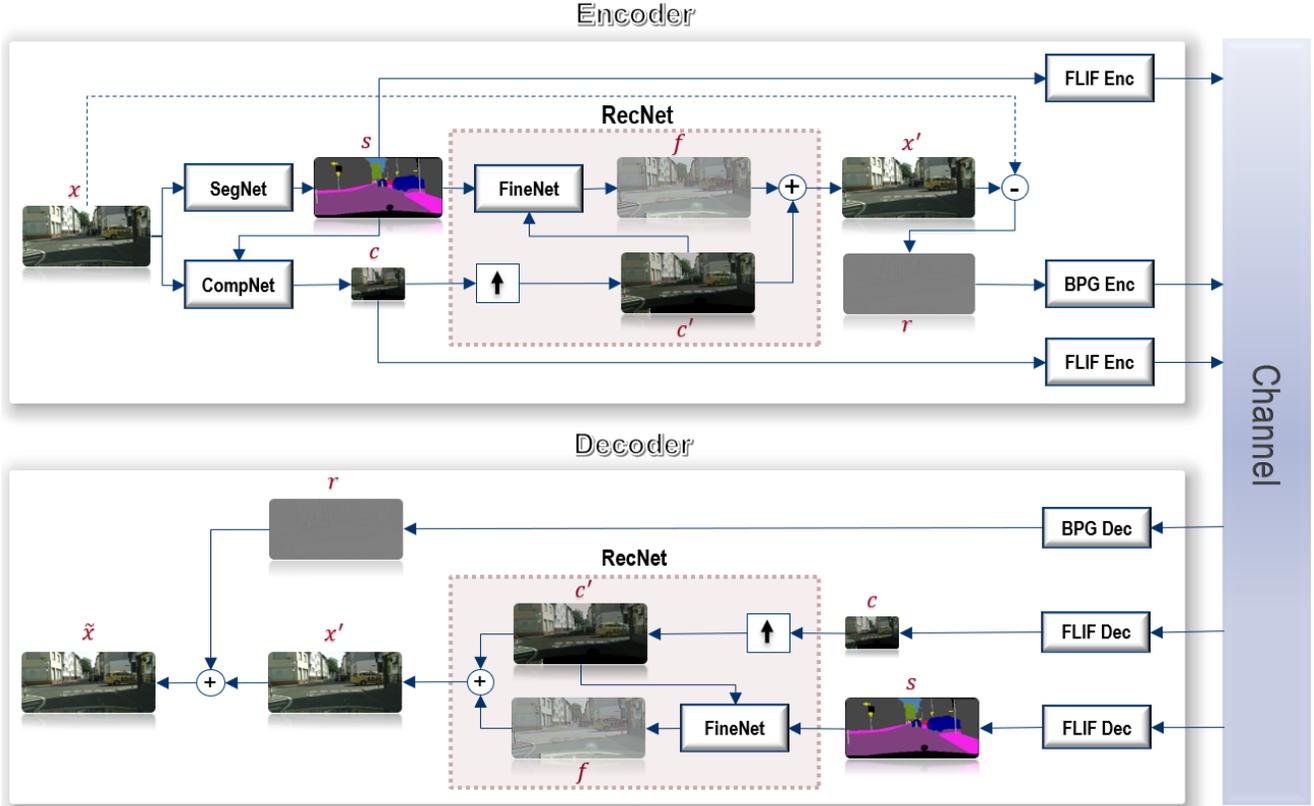


Figure 1: The overall framework of the proposed deep semantic segmentation-based layered image compression (DSSLIC) codec.

and the input image. Therefore the proposed scheme includes three layers of information.

Experimental results in the RGB (4:4:4) domain show that the proposed framework outperforms the H.265/HEVC-based BPG codec [4] in both PSNR and multi-scale structural similarity index (MS-SSIM) [21] metrics across a large range of bit rates, and is much better than JPEG, JPEG2000, and WebP [7]. For example, our method can be 4.7 dB better in PSNR than BPG for some Kodak testing images. Moreover, since semantic segmentation map is included in the bit-stream, the proposed scheme can facilitate many other tasks such as image search and object-based adaptive image compression.

The idea of semantic segmentation-based compression was already studied in MPEG-4 object-based video coding in the 1990's [14]. However, due to the lack of high-quality and fast segmentation methods, object-based image/video coding has not been widely adopted. Thanks to the rapid development of deep learning algorithms and hardware, it is now the time to revisit this approach.

This paper is organized as follows. In Section 2, the works related to learning-based image compression are briefly reviewed. The architecture of the proposed frame-

work and the corresponding formulation and objective functions are described in Section 3. In Section 4, the performance of the proposed method is evaluated and compared with the JPEG, JPEG2000, WebP, and BPG codecs.

2. Related Works

In traditional compression methods, many components such as entropy coding are hand-crafted. Deep learning-based approaches have the potential of automatically discovering and exploiting the features of the data; thereby achieving better compression performance.

In the last few years, various learning-based image compression frameworks have been proposed. In [16, 17], long short-term memory (LSTM)-based recurrent neural networks (RNNs) were used to extract binary representations, which were then compressed with entropy coding. Probability estimation in the entropy coding was also handled by LSTM convolution. Johnston et al. [8] utilized structural similarity (SSIM) loss [20] and spatially adaptive bit allocation to further improve the performance.

In [3], a scheme that involved a generalized divisive normalization (GDN)-based nonlinear analysis transform, a

uniform quantizer, and a nonlinear synthesis transform were developed. Theis et al. [15] proposed a compressive autoencoder (AE) where the quantization was replaced by a smooth approximation, and a scaling approach was used to get different rates. In [1], a soft-to-hard vector quantization approach was introduced, and a unified formulation was developed for both the compression of deep learning models and image compression.

GAN has been exploited in a number of learning-based image compression schemes. In [11], a discriminator was used to help training the decoder. A perceptual loss based on the feature map of an ImageNet-pretrained AlexNet was introduced although only low-resolution image coding results were reported in [11]. In [10], AE was embedded in the GAN framework in which the feature extraction adopted pyramid and interscale alignment. The discriminator also extracted outputs from different layers, similar to the pyramid feature generation. An adaptive training was used where the discriminator was trained and a confusion signal was propagated through the reconstructor, depending on the prediction accuracy of the discriminator.

Recently, there have also been some efforts in combining some computer vision tasks and image compression in one framework. In [9, 18], the authors tried to use the feature maps from learning-based image compression to help other tasks such as image classification and semantic segmentation although the results from other tasks were not used to help the compression part. In [2], a segmentation map-based image synthesis model was proposed, which targeted extremely low bit rates (< 0.1 bits/pixel), and used synthesized images for non-important regions.

3. Deep-semantic Segmentation-based Layered Image Compression (DSSLIC)

In this section, the proposed semantic segmentation-based layered image compression approach is described. The architecture of the codec and the corresponding deep networks used in the codec are first presented. The loss functions used for training the model are then formulated and explained.

3.1. DSSLIC Codec Framework

The overall framework of the DSSLIC codec is shown in Fig. 1. The encoder includes three deep learning networks: *SegNet*, *CompNet*, and *FineNet*. The semantic segmentation map s of the input image x is first obtained using *SegNet*. In this paper, a pre-trained PSPNet proposed in [23] is used as *SegNet*. The segmentation map is encoded to serve as side information to *CompNet* for generating a low-dimensional version c of the original image. In this paper, both s and c are losslessly encoded using the FLIF codec [13], which is a state-of-the-art lossless image codec.

Given the segmentation map s and compact image c , the *RecNet* part tries to obtain a high-quality reconstruction of the input image. Inside the *RecNet*, the compact image c is first upsampled, which, together with the segmentation map s , is fed into a *FineNet*. Note that although GAN-based synthesized images from segmentation maps are visually appealing, their details can be quite different from the original images. To minimize the distortion of the synthesized images, we modify the existing segmentation-based synthesis framework in [19] and add the upsampled version of the compact image c as an additional input. Besides, *FineNet* is trained to learn the missing fine information of the upsampled version of c with respect to the input image. This is easier to control the output of the GAN network. After adding the upsampled version of c and the *FineNet*'s output f , we get a better estimate of the input.

In our scheme, if the *SegNet* fails to assign any label to an area, the *FineNet* will ignore the semantic input and only reconstruct the image from c , which can still get good results. Therefore, our scheme is applicable to all general images. The residual r between the input and the estimate is then obtained and encoded by a lossy codec. In order to deal with negative values, the residual image r is rescaled to $[0, 255]$ with min-max normalization before encoding. The min and max values are also sent to decoder for inverse scaling. In this paper, the H.265/HEVC intra coding-based BPG codec is used [4], which is state-of-the-art in lossy coding.

As a result, in our scheme, the segmentation map s serves as the base layer, and the compact image c and the residual r are respectively the first and second enhancement layers.

At the decoder side, the segmentation map and compact representation are decoded to be used by *RecNet* to get an estimate of the input image. The output of *RecNet* is then added to the decoded residual image to get the final reconstruction of the image \tilde{x} . The pseudo code of the encoding and decoding procedures is given in Algorithm 1.

3.2. Network Architecture

The architectures of the *CompNet* (proposed in this work) and *FineNet* (modified from [19]) networks are defined as follows:

- *CompNet*: $c_{64}, d_{128}, d_{256}, d_{512}, c_3, \tanh$
- *FineNet*:
 $c_{64}, d_{128}, d_{256}, d_{512}, 9 \times r_{512}, u_{256}, u_{128}, u_{64}, c_3, \tanh$

where

- c_k : 7×7 convolution layers (with k filters and stride 1) followed by instance normalization and ReLU.
- d_k : 3×3 convolution layers (with k filters and stride 1) followed by instance normalization and ReLU.

Algorithm 1 DSSLIC Codec

```
procedure ENCODE( $x$ )
   $s \leftarrow \text{SegNet}(x)$ 
   $\triangleright$  encode  $s$  (1st enhancement layer)
   $c \leftarrow \text{CompNet}(x, s)$ 
   $\triangleright$  encode  $c$  (base layer)
   $x' \leftarrow \text{RecNet}(s, c)$ 
   $r \leftarrow x - x'$ 
   $\min, \max \leftarrow \text{Min}(r), \text{Max}(r)$ 
   $r \leftarrow \frac{r - \min}{(\max - \min)} * 255$ 
   $\triangleright$  encode  $r$  (2nd enhancement layer)

procedure DECODE( $s, c, r, \min, \max$ )
   $x' \leftarrow \text{RecNet}(s, c)$ 
   $r \leftarrow \frac{r * (\max - \min)}{255} + \min$ 
   $\tilde{x} \leftarrow x' + r$ 

function RECNET( $s, c$ )
   $c' \leftarrow \text{upsample}(c)$ 
   $f \leftarrow \text{FineNet}(s, c')$ 
   $x' \leftarrow c' + f$ 
  return  $x'$ 
```

- r_k : a residual block containing reflection padding and two 3×3 convolution layers (with k filters) followed by instance normalization.
- u_k : 3×3 fractional-strided-convolution layers (with k filters and stride $\frac{1}{2}$) followed by instance normalization and ReLU.

Inspired by [19], for the adversarial training of the proposed model, two discriminators denoted by D_1 and D_2 operating at two different image scales are used in this work. D_1 operates at the original scale and has a more global view of the image. Thus, the generator can be guided to synthesize fine details in the image. On the other hand, D_2 operates with $2 \times$ down-sampled images, leading to coarse information in the synthesized image. Both discriminators have the following architecture:

- $C_{64}, C_{128}, C_{256}, C_{512}$

where C_k denotes 4×4 convolution layers with k filters and stride 2 followed by instance normalization and LeakyReLU. In order to produce a 1-D output, a convolution layer with 1 filter is utilized after the last layer of the discriminator.

3.3. Formulation and Objective Functions

Let $x \in \mathbb{R}^{h \times w \times k}$ be the original image, the corresponding semantic segmentation map $s \in \mathbb{Z}^{h \times w}$ and the compact representation $c \in \mathbb{R}^{\frac{h}{\alpha} \times \frac{w}{\alpha} \times k}$ are generated as follows:

$$s = \text{SegNet}(x), c = \text{CompNet}(s, x), \quad (2)$$

Conditioned on s and the upsampled c , denoted by $c' \in \mathbb{R}^{h \times w \times k}$, *FineNet* (our GAN generator) reconstructs the fine information image, denoted by $f \in \mathbb{R}^{h \times w \times k}$, which is then added to c' to get the estimate of the input:

$$x' = c' + f, \text{ where } f = \text{FineNet}(s, c'). \quad (3)$$

The error between x and x' is measured using a combination of different losses including \mathcal{L}_1 , \mathcal{L}_{SSIM} , \mathcal{L}_{DIS} , \mathcal{L}_{VGG} , and GAN losses. The L1-norm loss (least absolute errors) is defined as:

$$\mathcal{L}_1 = 2\lambda \|x - x'\|_1. \quad (4)$$

It has been shown that combining pixel-wise losses such as \mathcal{L}_1 with SSIM loss can significantly improve the perceptual quality of the reconstructed images [22]. As a result, we also utilize the SSIM loss in our work, which is defined as

$$\mathcal{L}_{SSIM} = -I(x, x') \cdot C(x, x') \cdot S(x, x'), \quad (5)$$

where the three comparison functions luminance I , contrast C , and structure S are computed as:

$$I(x, x') = \frac{2\mu_x \mu_{x'} + C_1}{\mu_x^2 + \mu_{x'}^2 + C_1}, C(x, x') = \frac{2\sigma_x \sigma_{x'} + C_2}{\sigma_x^2 + \sigma_{x'}^2 + C_2}, \\ S(x, x') = \frac{\sigma_{xx'} + C_3}{\sigma_x \sigma_{x'} + C_3}, \quad (6)$$

where μ_x and $\mu_{x'}$ are the means of x and x' , σ_x and $\sigma_{x'}$ are the standard deviations, and $\sigma_{xx'}$ is the correlation coefficient. C_1 , C_2 , and C_3 are the constants used for numerical stability.

To stabilize the training of the generator and produce natural statistics, two perceptual feature-matching losses based on the discriminator and VGG networks [12] are employed. The discriminator-based loss is calculated as:

$$\mathcal{L}_{DIS} = \lambda \sum_{d=1,2} \sum_{i=1}^n \frac{1}{N_i} \|D_d^{(i)}(s, c', x) - D_d^{(i)}(s, c', x')\|_1, \quad (6)$$

where $D_d^{(i)}$ denotes the features extracted from the i -th intermediate layer of the discriminator network D_d (with n layers and N_i number of elements in each layer). Similar to [11], a pre-trained VGG network with m layers and M_j elements in each layer is used to construct the VGG perceptual loss as in below:

$$\mathcal{L}_{VGG} = \lambda \sum_{j=1}^m \frac{1}{M_j} \|V^{(j)}(x) - V^{(j)}(x')\|_1, \quad (7)$$

where V_j represents the features extracted from the j -th layer of VGG.

In order to distinguish the real training image x from the reconstructed image x' , given s and c' , the following objective function is minimized by the discriminator D_d :

$$\mathcal{L}_D = - \sum_{d=1,2} (\log D_d(s, c', x) + \log(1 - D_d(s, c', x'))), \quad (8)$$

while the generator (*FineNet* in this work) tries to fool D_d by minimizing $-\sum_{d=1,2} \log D_d(s, c', x')$. The final generator loss including all the reconstruction and perceptual losses is then defined as:

$$\mathcal{L}_G = - \sum_{d=1,2} \log D_d(s, c', x') + \mathcal{L}_1 + \mathcal{L}_{SSIM} + \mathcal{L}_{DIS} + \mathcal{L}_{VGG}. \quad (9)$$

Finally, our goal is to minimize the following hybrid loss function:

$$\mathcal{L} = \mathcal{L}_D + \mathcal{L}_G. \quad (10)$$

3.4. Training

The Cityscapes (with 30 semantic labels) [5] and ADE20K (with 150 semantic labels) [24] datasets are used for training the proposed model. For Cityscapes, all the 2974 RGB images (street scenes) in the dataset are used. All images are then rescaled to 512×1024 (i.e., $h = 512$, $w = 1024$, and $k = 3$ for RGB channels). For ADE20K, the images with at least 512 pixels in height or width are used (9272 images in total). All images are rescaled to $h = 256$ and $w = 256$ to have a fixed size for training. Note that no resizing is needed for the test images since the model can work with any size at the testing time. We set the downsampling factor $\alpha = 8$ to get the compact representation of size $64 \times 128 \times 3$ for Cityscapes and $32 \times 32 \times 3$ for ADE20K. We also consider the weight $\lambda = 10$ for \mathcal{L}_1 , \mathcal{L}_{DIS} , and \mathcal{L}_{VGG} .

All models were jointly trained for 150 epochs with mini-batch stochastic gradient descent (SGD) and a mini-batch sizes of 2 and 8 for Cityscapes and ADE20K, respectively. The Adam solver with learning rate of 0.0002 was used, which is fixed for the first 100 epochs, but gradually decreases to zero for the next 50 epochs. Perceptual feature-matching losses usually guide the generator towards more synthesized textures in the predicted images, which causes a slightly higher pixel-wise reconstruction error, especially in the last epochs. To handle this issue, we did not consider the perceptual \mathcal{L}_D and \mathcal{L}_{VGG} losses in the generator loss for the last 50 epochs. All the *SegNet*, *CompNet*, *FineNet*, and the discriminator networks proposed in this work are trained in the RGB domain.

4. Experiments

In this section, we compare the performance of the proposed DSSLIC scheme with JPEG, JPEG2000, WebP, and

the H.265/HEVC intra coding-based BPG codec [4], which is state-of-the-art in lossy image compression. Since the networks are trained for RGB images, we encode all images using RGB (4:4:4) format in different codecs for fair comparison. We use both PSNR and MS-SSIM [21] as the evaluation metric in this experiment. In this experiment, we encode the RGB components of the residual image r using lossy BPG codec with different quantization values.

The results of the ADE20K and Cityscapes test sets are given in Figures 2 and 3. The results are averaged over 50 random test images not included in the training set. As shown in the figures, our method gives better PSNR and MS-SSIM than BPG, especially when the bit rate is less than ≈ 0.9 bits/pixel/channel (bpp for short) on ADE20K and less than ≈ 0.5 bpp on Cityscapes. In particular, the average PSNR gain is more than 2dB for the ADE20k test set when the bit rate is between 0.4-0.7 bpp.

To demonstrate the generalization capability of the scheme, the ADE20K-trained model is also applied to the classical Kodak dataset (including 24 test images). The average results of the Kodak dataset are illustrated in Figure 4. For this experiment, the model trained on the ADE20K dataset is used. It is shown that when the bit rate is less than about 1.4 bpp, our scheme achieves better results than other codecs in both PSNR and MS-SSIM. For example, the average gain is about 2 dB between 0.4-0.8 bpp. This is quite promising since the proposed scheme can still be improved in many ways. This also shows that our method generalizes very well when the training and testing images are from different datasets. The average *RecNet* decoding time for Kodak images on CPU and GPU are ≈ 44 s and ≈ 0.013 s, respectively.

Some visual examples from ADE20K, Cityscapes, and Kodak test sets are given in Figures 6-12. In order to have a more clear visualization, only some cropped parts of the reconstructed images are shown in these examples. As seen in all examples, JPEG has poor performance due to the blocking artifacts. Some artifacts are also seen on JPEG2000 results. Although WebP provides higher quality results than JPEG2000, the images are blurred in some areas. The images encoded using BPG are smoother, but the fine structures are also missing in some areas.

Figure 5 and Table 1 report some ablation studies of different configurations, all are obtained without using the BPG-based residual coding, including: **upComp**: the results are obtained without considering the *FineNet* network in the pipeline, i.e., $x' = c'$ (the upsampled compact image only); **noSeg**: the segmentation maps are not considered in neither *CompNet* nor *FineNet* networks, i.e., $x' = c' + f$ where c' is the upsampled version of $c = \text{CompNet}(x)$, and $f = \text{FineNet}(c')$; **withSeg**: all the DSSLIC components shown in Figure 1 are used in this configuration (except BPG-based residual coding); **synth**;

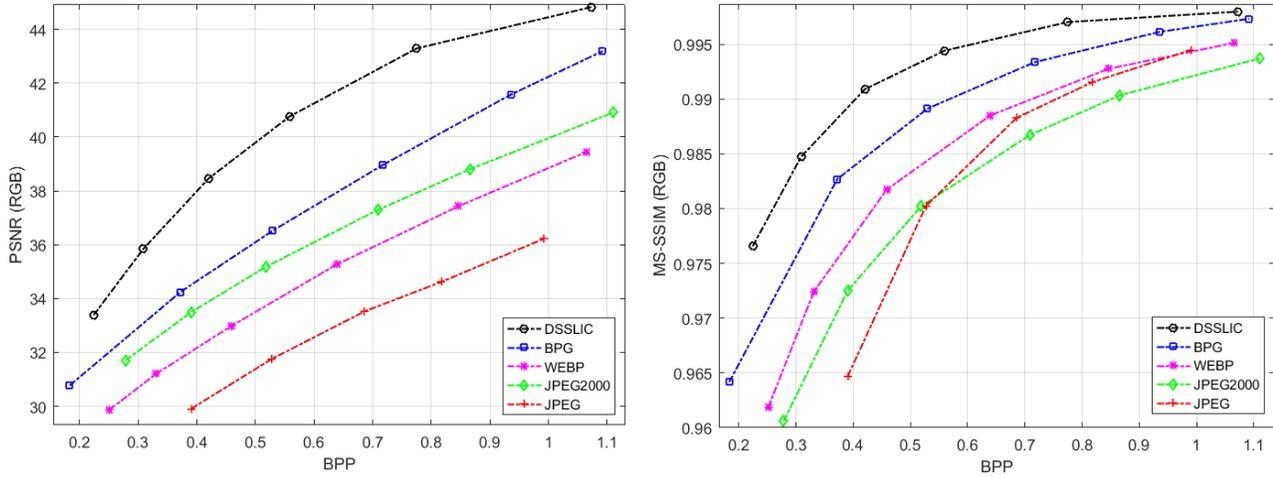


Figure 2: Comparison results on ADE20K test set in terms of PSNR (left) and MS-SSIM (right) vs. bpp (bits/pixel/channel). The results are averaged over RGB channels.

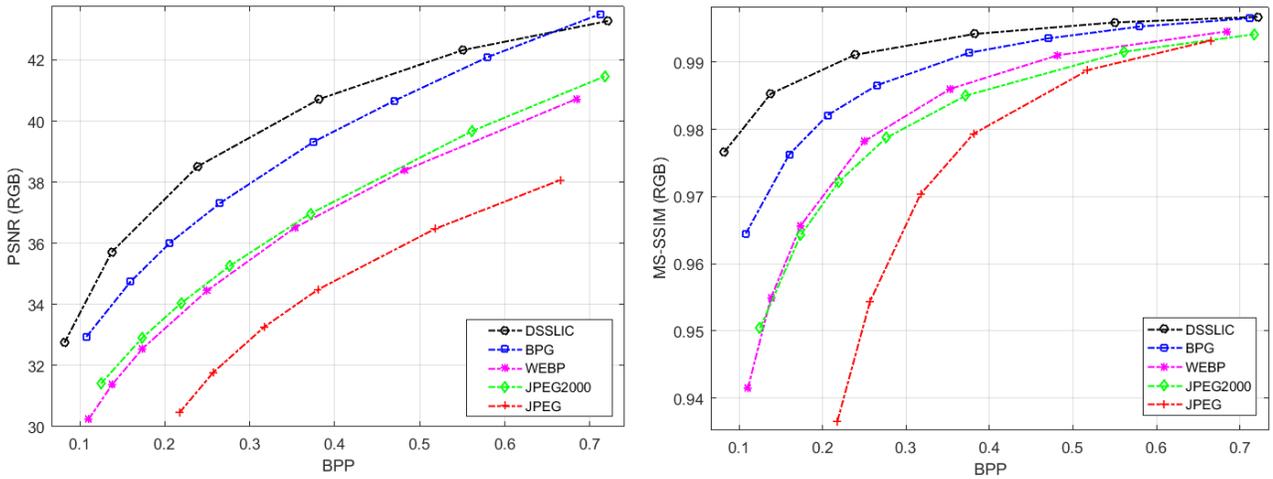


Figure 3: Comparison results on Cityscapes test set.

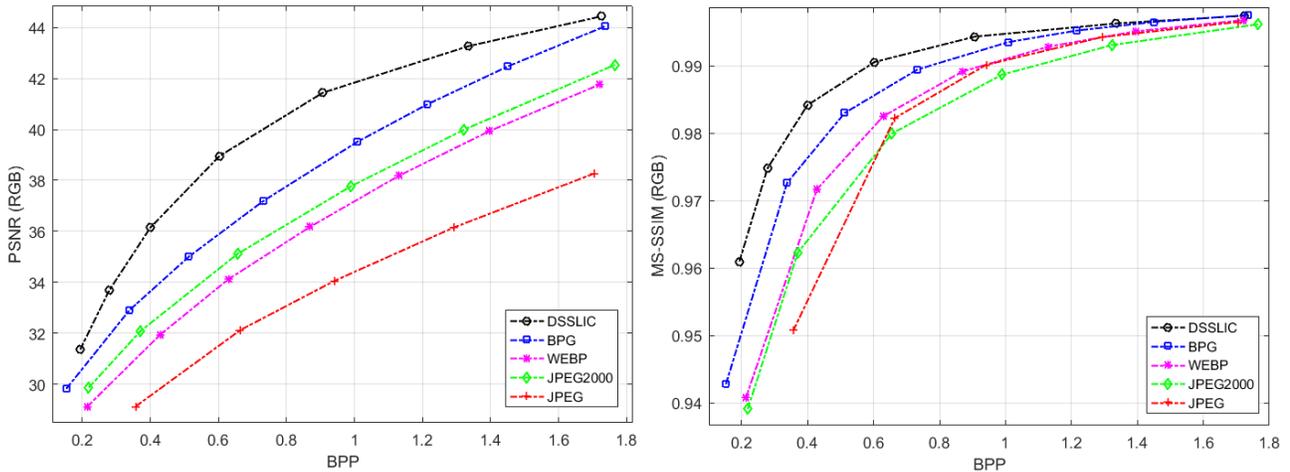


Figure 4: Comparison results on Kodak image set.



(a) Original PSNR, MS-SSIM
 (b) upComp 18.01 dB, 0.73
 (c) synth 23.68 dB, 0.84
 (d) noSeg 22.18 dB, 0.86
 (e) withSeg 25.09 dB, 0.88

Figure 5: Visual comparison of different scenarios at 0.08 BPP.

Table 1: Results of different scenarios (without BPG-based residual coding).

| | ADE20K | | | | Kodak | | | |
|----------------|--------|-------|-------|---------|--------|-------|-------|---------|
| | upComp | synth | noSeg | withSeg | upComp | synth | noSeg | withSeg |
| BPP | 0.095 | 0.092 | 0.08 | 0.095 | 0.087 | 0.088 | 0.080 | 0.087 |
| PSNR | 17.50 | 21.91 | 22.24 | 23.11 | 17.77 | 20.97 | 21.46 | 21.91 |
| MS-SSIM | 0.759 | 0.887 | 0.905 | 0.914 | 0.738 | 0.858 | 0.887 | 0.891 |

the settings in this configuration is the same as withSeg except that the perceptual losses \mathcal{L}_{VGG} and \mathcal{L}_{DIS} are considered in all training epochs. The poor performance of using only the upsampled compact images in **upComp** shows the importance of FN in predicting the missing fine information, which is also visually obvious in Figure 5. Considering perceptual losses in all training epochs (**synth**) leads to sharper and perceptually more natural images, but the PSNR is much lower. The results with segmentation maps (**withSeg**) provide slightly better PSNR than **noSeg** although the visual gain is more pronounced, e.g., the dark wall in Figure 5.

5. Conclusion

In this paper, we proposed a deep semantic segmentation-based layered image compression (DSSLIC) framework in which the semantic segmentation map of the input image was used to synthesize the image, and the residual was encoded as an enhancement layer in the bit-stream.

Experimental results showed that the proposed framework outperforms the H.265/HEVC-based BPG and the other standard codecs in both PSNR and MS-SSIM metrics in RGB (4:4:4) domain. In addition, since semantic segmentation map is included in the bit-stream, the proposed scheme can facilitate many other tasks such as image search and object-based adaptive image compression.

The proposed scheme opens up many future topics, for example, improving its high-rate performance, modifying the scheme for YUV-coded images, and applying the frame-

work for other tasks.

References

- [1] E. Agustsson, F. Mentzer, M. Tschannen, L. Cavigelli, R. Timofte, L. Benini, and L. Van Gool. Soft-to-hard vector quantization for end-to-end learning compressible representations. *arXiv preprint arXiv:1704.00648*, 2017.
- [2] E. Agustsson, M. Tschannen, F. Mentzer, R. Timofte, and L. V. Gool. Generative adversarial networks for extreme learned image compression. *arXiv preprint arXiv:1804.02958*, Apr. 2018.
- [3] J. Ballé, V. Laparra, and E. P. Simoncelli. End-to-end optimized image compression. *arXiv preprint arXiv:1611.01704*, 2016.
- [4] F. Bellard. Bpg image format (<http://bellard.org/bpg/>), 2017.
- [5] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele. The CityScapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016.
- [6] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, pages 2672–2680, 2014.
- [7] Google Inc. WebP (<https://developers.google.com/speed/webp/>), 2016.
- [8] N. Johnston, D. Vincent, D. Minnen, M. Covell, S. Singh, T. Chinen, S. J. Hwang, J. Shor, and G. Toderici. Improved lossy image compression with priming and spatially adaptive bit rates for recurrent networks. *arXiv preprint arXiv:1703.10114*, 2017.
- [9] S. Luo, Y. Yang, and M. Song. DeepSIC: Deep semantic image compression. *arXiv preprint arXiv:1801.09468*, 2018.



Figure 6: ADE20K visual example 1. (bits/pixel/channel, PSNR, MS-SSIM)



Figure 7: ADE20K visual example 2. (bits/pixel/channel, PSNR, MS-SSIM)



Figure 8: Cityscapes visual example 1. (bits/pixel/channel, PSNR, MS-SSIM)



Figure 9: Cityscapes visual example 2. (bits/pixel/channel, PSNR, MS-SSIM)



Figure 10: Kodak visual example 1. (bits/pixel/channel, PSNR, MS-SSIM)

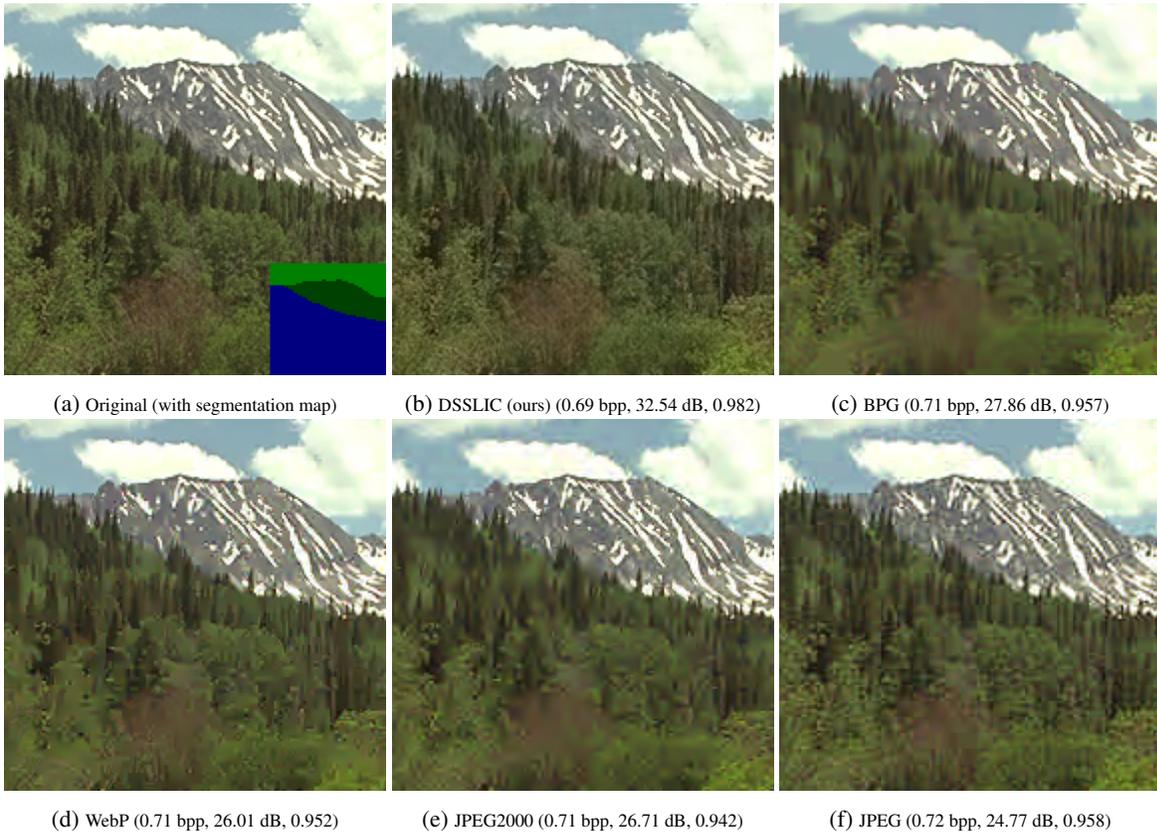


Figure 11: Kodak visual example 2. (bits/pixel/channel, PSNR, MS-SSIM)

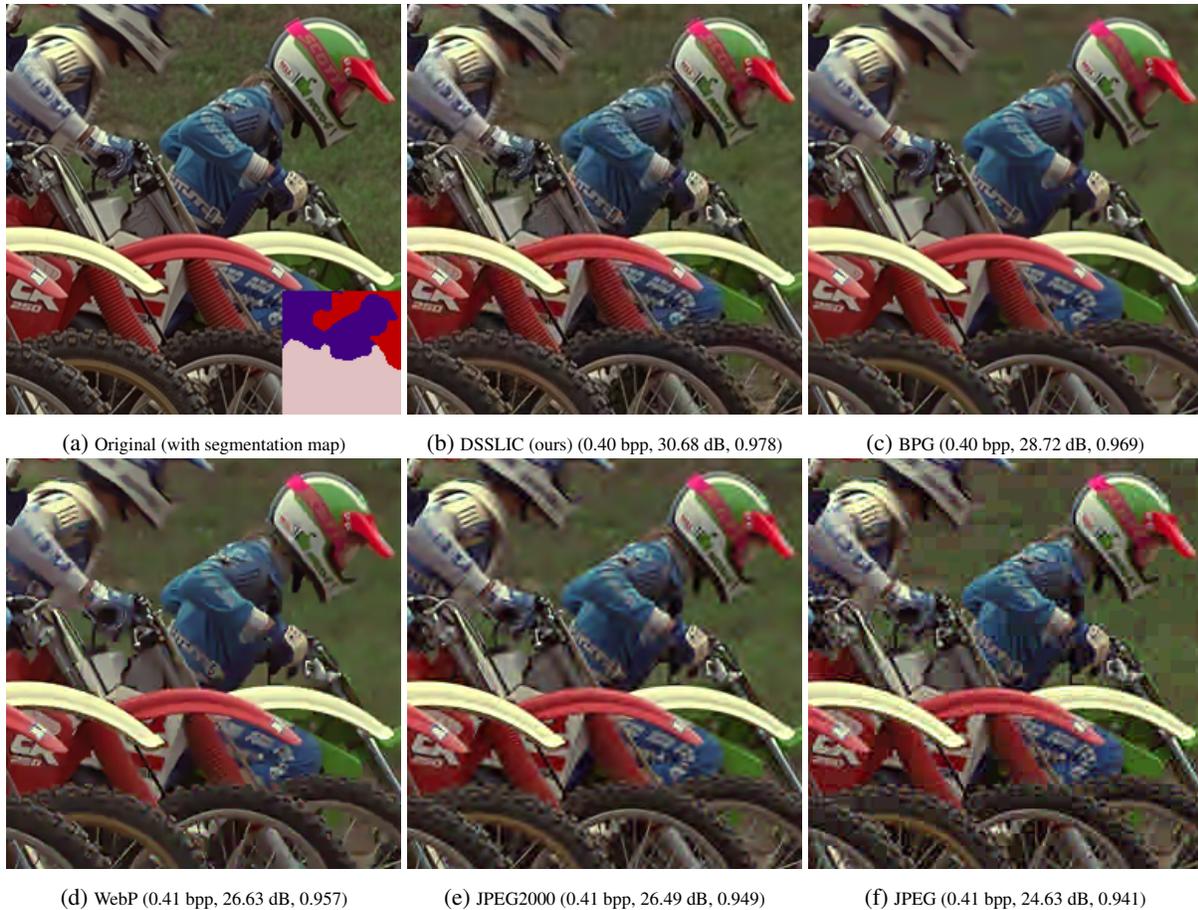


Figure 12: Kodak visual example 3. (bits/pixel/channel, PSNR, MS-SSIM)

- [10] O. Rippel and L. Bourdev. Real-time adaptive image compression. *arXiv preprint arXiv:1705.05823*, 2017.
- [11] S. Santurkar, D. Budden, and N. Shavit. Generative compression. *arXiv preprint arXiv:1703.01467*, 2017.
- [12] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [13] J. Sneyers and P. Wuille. FLIF: Free lossless image format based on maniac compression. In *Image Processing (ICIP), 2016 IEEE International Conference on*, pages 66–70. IEEE, 2016.
- [14] R. Talluri, K. Oehler, T. Bannon, J. Courtney, A. Das, and J. Liao. A robust, scalable, object-based video compression technique for very low bit-rate coding. *IEEE Trans. Circuits and Systems for Video Tech.*, 7(1):221–233, 1997.
- [15] L. Theis, W. Shi, A. Cunningham, and F. Huszár. Lossy image compression with compressive autoencoders. *arXiv preprint arXiv:1703.00395*, 2017.
- [16] 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.
- [17] G. Toderici, D. Vincent, N. Johnston, S. J. Hwang, D. Minnen, J. Shor, and M. Covell. Full resolution image compression with recurrent neural networks. In *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*, pages 5435–5443. IEEE, 2017.
- [18] R. Torfason, F. Mentzer, E. Agustsson, M. Tschannen, R. Timofte, and L. Van Gool. Towards image understanding from deep compression without decoding. *arXiv preprint arXiv:1803.06131*, 2018.
- [19] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, A. Tao, J. Kautz, and B. Catanzaro. High-resolution image synthesis and semantic manipulation with conditional GANs. *arXiv preprint arXiv:1711.11585*, 2017.
- [20] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
- [21] Z. Wang, E. P. Simoncelli, and A. C. Bovik. Multiscale structural similarity for image quality assessment. In *Signals, Systems and Computers, 2004. Conference Record of the Thirty-Seventh Asilomar Conference on*, volume 2, pages 1398–1402. Ieee, 2003.

- [22] H. Zhao, O. Gallo, I. Frosio, and J. Kautz. Loss functions for image restoration with neural networks. *IEEE Transactions on Computational Imaging*, 3(1):47–57, 2017.
- [23] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia. Pyramid scene parsing network. In *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 2881–2890, 2017.
- [24] B. Zhou, H. Zhao, X. Puig, S. Fidler, A. Barriuso, and A. Torralba. Scene parsing through ade20k dataset. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, volume 1, page 4. IEEE, 2017.