

Semantic-Sparse Colorization Network for Deep Exemplar-based Colorization

Yunpeng Bai,¹ Chao Dong,² Zenghao Chai,¹ Andong Wang,¹ Zhengzhuo Xu,¹
Chun Yuan,^{1,3}

¹ Shenzhen International Graduate School, Tsinghua University

² Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences ³ Peng Cheng Laboratory
{byp20, wad20, xzz20}@mails.tsinghua.edu.cn, chao.dong@siat.ac.cn, zenghaochai@gmail.com, yuanc@sz.tsinghua.edu.cn

Abstract

Exemplar-based colorization approaches rely on reference image to provide plausible colors for target gray-scale image. The key and difficulty of exemplar-based colorization is to establish an accurate correspondence between these two images. Previous approaches have attempted to construct such a correspondence but are faced with two obstacles. First, using luminance channels for the calculation of correspondence is inaccurate. Second, the dense correspondence they built introduces wrong matching results and increases the computation burden. To address these two problems, we propose Semantic-Sparse Colorization Network (SSCN) to transfer both the global image style and detailed semantic-related colors to the gray-scale image in a coarse-to-fine manner. Our network can perfectly balance the global and local colors while alleviating the ambiguous matching problem. Experiments show that our method outperforms existing methods in both quantitative and qualitative evaluation and achieves state-of-the-art performance.

Introduction

Image colorization is a classic and appealing task that predicts the vivid colors from a gray-scale image. As there is no unique correct color for a given pixel, three classes of methods are proposed to constrain the output color space. The first one is called automatic colorization, such as (Cheng, Yang, and Sheng 2015; Zhang, Isola, and Efros 2016). These methods generally rely on the powerful convolutional networks and learn a direct mapping from a large-scale image dataset. The second class introduces additional human intervention, such as user-guided scribbles (Zhang et al. 2017; Sangkloy et al. 2017; Ci et al. 2018) and text (Manjunatha et al. 2018; Bahng et al. 2018). They require users to provide reliable color/text labels for more dedicated colorization. While the third class, denoted as exemplar-based method (Lu et al. 2020; He et al. 2018; Xiao et al. 2020; Xu et al. 2020; Lee et al. 2020; Bugeau, Ta, and Papadakis 2014; Chia et al. 2011; Gupta et al. 2012), is a trade-off between fully automatic and human intervention strategies. It adopts a reference image as guidance and generates a similar color-style image. These three kinds of methods have different applications and prior information, thus cannot be compared side-by-side. In this work, we study exemplar-based image

colorization, due to its large flexibility and excellent performance.

The difficulty of exemplar-based image colorization is to build an accurate correspondence between the gray-scale image and the color reference. Recent studies have achieved considerable progress but still have significant limitations. Some works regard colorization as a style transfer problem (Xu et al. 2020), and usually transfer the global color tones. As a result, they lack detailed color matching between semantically similar objects/parts. Other researchers (Lu et al. 2020) have attempted to construct a dense correspondence, but are faced with two obstacles. First, the correspondence is calculated between the luminance channels of the input and reference images. However, as gray-scale images do not contain enough semantic information as color images (a common knowledge in image classification (He et al. 2018)), the correspondence purely based on luminance channels is inaccurate. Second, the dense correspondence itself will also bring in unavoidable drawbacks. It not only introduces wrong matching results for semantically unrelated objects, but also increases the computation burden.

To address the above mentioned problems, we propose a new coarse-to-fine colorization framework – Semantic-Sparse Colorization Network – to transfer both the global image style and the detailed semantic-related colors to the gray-scale image. Specifically, in the coarse colorization stage, we adopt an image transfer network to obtain a preliminary colorized result. The color information of the reference image is encoded as a vector, which is then migrated to the gray-scale image by an AdaIN operation. In the fine colorization stage, we will first calculate the semantic correspondence between the coarse result and the reference image. Specially, only the semantic-significant parts and some background regions are reserved for calculation, leading to a sparse correlation matrix. Then the attention mechanism will be used to re-weight the reference image and help generate the final color result. The proposed method can perfectly balance the global and local colors while alleviating the ambiguous matching problem caused by dense correspondence. Extensive experiments have shown the priority of our network towards other state-of-the-art methods. To facilitate numerical evaluation, we also propose a unified evaluation pipeline for all exemplar-based colorization methods.

Our main contributions are summarized as follows:

- We propose to build a more accurate correspondence between a coarse-colored result and the reference image. It not only minimizes the information gap between the gray-scale input and the color reference, but also achieves better performance on details.
- We propose a sparse attention mechanism to make the model focus on the semantically significant regions in the reference image. It could produce more detailed results with a much lower computation cost.
- We collect a new test dataset from ImageNet to solve the problem of fair comparison. We also design a new quantitative evaluation metric to evaluate exemplar-based colorization methods.

Related Work

Because image colorization plays an essential role in image processing tasks such as old photo restoration and image editing, this subject has been studied for a long time (Charpiat, Hofmann, and Schölkopf 2008; Bugeau, Ta, and Papadakis 2014; Qu, Wong, and Heng 2006; Luan et al. 2007; Huang et al. 2005). Recently, many studies have used learning-based methods to solve this ill-posed problem. These approaches can be roughly grouped into three classes.

The first one is called automatic colorization, which directly maps gray-scale images to color images, such as (Cheng, Yang, and Sheng 2015) and (Zhang, Isola, and Efros 2016). They are the earliest methods to use convolutional networks to learn the mapping from a large-scale image dataset. MemoPainter (Yoo et al. 2019) uses a memory network to “memorize” rare examples, which can avoid the interference of dominant color in the dataset and make the model perform well even without sufficient data. More recently, Transformer has also been applied to address this task (Kumar, Weissenborn, and Kalchbrenner 2021). Some works take advantage of generative models to promote the diversity of image colorization results, such as (Cao et al. 2017), which adopts generative adversarial network (GAN) to generate color image and gray-scale image is taken as a condition. Variational autoencoder (VAE) architecture was used in (Deshpande et al. 2017) to learn a low dimensional embedding of color distribution, and sampled embeddings can be used to produce multiple colorizations. However, these methods sometimes will lead to unrealistic color effects, and their colorization process is uncontrollable.

The second class introduces additional human intervention, such as user-guided scribbles and text. They require users to provide reliable color/text labels for more dedicated colorization. Traditional scribble-based colorization methods (Levin, Lischinski, and Weiss 2004; Xu et al. 2009) usually propagate the local user hints to the whole image via an optimization approach, while learning-based methods (Zhang et al. 2017; Sangkloy et al. 2017; Ci et al. 2018) will combine color prior learned from large-scale image dataset with user’s intervention for colorization. These methods require a certain amount of human effort, and the quality of results depends on the user’s skills. Text-based colorization methods usually adopt image captions (Manjunatha et al. 2018) or palettes converted from the text (Bahng et al. 2018)

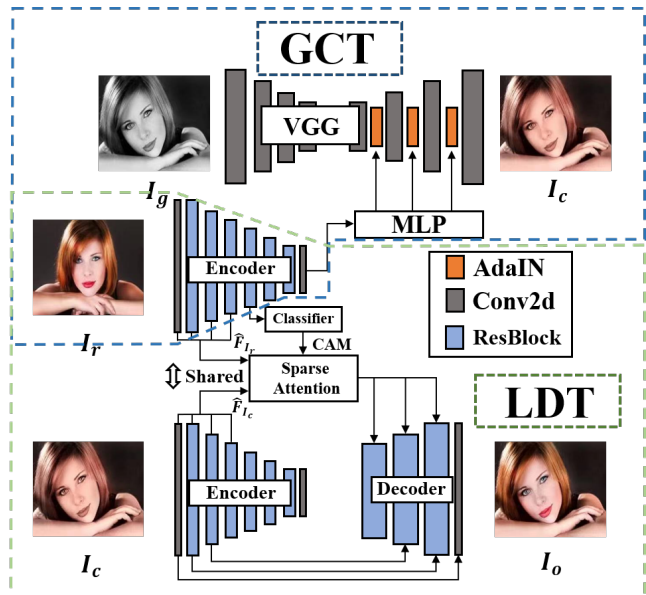


Figure 1: The illustration of the proposed two-stage image colorization framework.

as means of intervention. However, the color information represented by text is challenging to transfer to the image accurately.

The third class, denoted as exemplar-based method, is a trade-off between fully automatic and human intervention strategies. Compared to the above two classes, it adopts sample reference images to provide rich colors without requiring the user to do too much manual work. The key and difficulty of exemplar-based colorization is to establish an accurate correspondence between these two images. DEPN (Xiao et al. 2020) uses a pyramid structure to exploit multi-scale color information, but it only captures the global tones because no semantic correspondence is established. Some works (Xu et al. 2020) regard exemplar-based colorization as a style transfer problem, but cannot guarantee the correctness of semantics because they also lack a correspondence. Deep Image Analogy (Liao et al. 2017) was used in (He et al. 2018) to make the target and reference luminance channels aligned to get a coarse chrominance map for further refinement. (Lu et al. 2020) used features extracted from the luminance channel of the target and reference images to obtain dense correspondence. However, inaccuracies caused by using luminance channels to calculate correspondence and wrong matching problems introduced by dense correspondence will lead to unsatisfactory results.

Methods

Overview of the Proposed Method

The task of exemplar-based colorization can be formulated as follows. Given a gray-scale image I_g , which only contains the luminance channel l , our goal is to predict the corresponding a and b color channels in the CIE Lab color space, according to the reference color image I_r . The main chal-

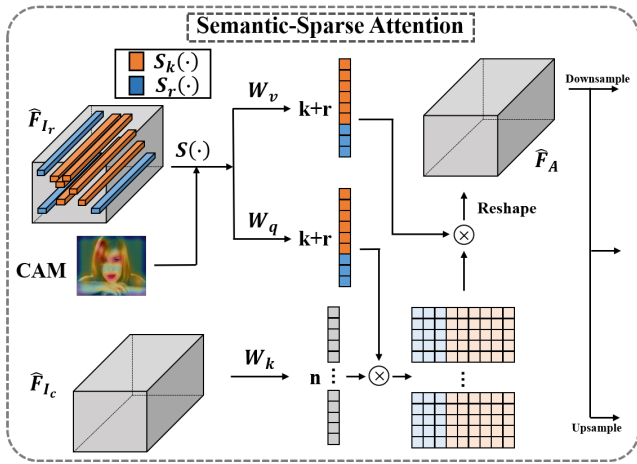


Figure 2: Structure of Semantic-Sparse attention.

lenge is to build an appropriate correspondence between the gray-scale image and the color reference. In order to make full use of the color information in the reference image, we will utilize the reference image twice in a coarse-to-fine manner during the whole colorization process. The proposed framework, namely Semantic-Sparse Colorization Network (SSCN), consists of two auxiliary modules, which transfer global and local colors in the reference image, respectively.

Specifically, taking the reference image I_r as input, our model will first encode it into features F_{I_r} . These features will be used in both global and local coloring modules. In the coarse colorization stage, the Global Color Transfer (GCT) module will use F_{I_r} to preliminarily color the gray-scale image I_g , and get a coarse-colored result I_c , which has similar global tones as I_r . Then the coarse output I_c will be further encoded into features F_{I_c} with the same encoder as F_{I_r} . In the fine colorization stage, the Local Details Transfer (LDT) module will use F_{I_r} and F_{I_c} to construct a correspondence that focuses on the semantically relevant regions of I_r . Note that these regions are sparsely selected according to their semantic levels. Based on the predicted mappings from LDT, the reference features F_{I_r} are reorganized and fused with F_{I_c} at different scales. Finally, the decoder takes the fused color features to produce the a and b channels of the input image I_g . The overall pipeline of our method SSCN is illustrated in Figure 1.

Global Color Transfer

We will first introduce the encoder of I_r , which is shared in both GCT and LDT modules. The encoder consists of six residual blocks. The last layer of F_{I_r} is passed through an MLP to form the style vector, which will be used in the GCT module for global style transfer. In GCT, the gray-scale image I_g will first be encoded into features $\{x_1, x_2, \dots, x_n\}$. Then, we perform coarse colorization in the feature space by changing feature statistics with AdaIN operation, as Formula 1:

$$AdaIN(x_i, y) = y_{s,i} \frac{x_i - \mu(x_i)}{\sigma(x_i)} + y_{b,i}, \quad (1)$$

where $\mu(x_i)$ and $\sigma(x_i)$ represent the i^{th} feature map's mean and variance, respectively. y_s and y_b are the affine parameters of the style vector, which is obtained from F_{I_r} via MLP transformation. Each feature map x_i is normalized separately and then scaled/biased using the corresponding coefficients from $y(y_s, y_b)$. After affine transformation, each feature channel will have the activation for certain color information. These features can be inverted to the Lab space by a convolutional decoder. We finally get the coarse colored result I_c of the coarse colorization stage. In our implementation, the encoder uses sub-layers of the VGG19 (Simonyan and Zisserman 2015), and the decoder is symmetric structure. AdaIN are added after CNN layers of the decoder.

Local Details Transfer

The main target of the LDT module is to build a more detailed and accurate correspondence between the coarse-colored result I_c and the reference image I_r . To begin with, we encode I_c into the corresponding features F_{I_c} , with the same encoder as F_{I_r} . To find their correspondence, we extract features from the first four layers of F_{I_r} and F_{I_c} , and resize them to the same spatial size of $1/4$ input image. Then these features are concatenated to form features \hat{F}_{I_r} and \hat{F}_{I_c} , corresponding to the latent states of coarse and reference image, respectively. Their spatial size is both $d \times H/4 \times W/4$, where d is the number of feature maps. To facilitate computation, they are further flattened in the second and third direction, and form features of size $d \times HW/16$. In this way, we segment the input image into $HW/16$ regions and represent each region with a d dimensional vector.

Based on the obtained features \hat{F}_{I_r} and \hat{F}_{I_c} , the LDT module will calculate a correlation matrix A via attention mechanism, whose element is computed by the scaled dot product (Vaswani et al. 2017) illustrated as Formula 2:

$$\alpha_{ij} = \underset{j}{softmax} \left(\frac{(W_q f_i^c) \cdot (W_k f_j^r)}{\sqrt{d}} \right). \quad (2)$$

Here, α_{ij} represents the similarity between the i -th region of \hat{F}_{I_c} and the j -th region of \hat{F}_{I_r} . \hat{F}_{I_c} is used to retrieve relevant local details from \hat{F}_{I_r} . Then, we can re-weight the features \hat{F}_{I_r} to obtain the attended feature \hat{F}_a through a weighted sum operation as Formula 3:

$$f_i^a = \sum_j \alpha_{ij} W_v f_j^r, \quad (3)$$

where W_q , W_k and W_v represent the linear transformation matrix into *query*, *key*, and *value* vectors, respectively. The attended features \hat{F}_a will be reshaped to the size of $d \times H/4 \times W/4$ and further resized into a suitable shape, fused with the features F_{I_c} at different scales and fed into the U-Net (Ronneberger, Fischer, and Brox 2015) decoder for the final detailed result of the fine colorization stage.

Semantic-Sparse Correspondence. In the above description, we use a standard attention mechanism to calculate the dense correspondence between coarse and reference images. We further propose a semantically sparse correspondence for better results with less computation cost. To be

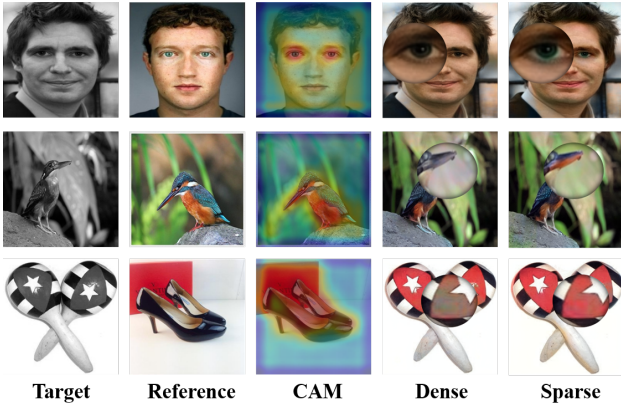


Figure 3: Comparison results of dense and sparse correspondence strategies.

specific, the reference features \hat{F}_{I_r} will go through a select operation. First, the fifth layer of F_{I_r} will be fed into a classifier and get a class activation map (CAM) (Zhou et al. 2016), which is used as the reference for selection. The CAM is flattened to $C = \{c_1, c_2, \dots, c_{HW/16}\} \in \mathbb{R}^{HW/16}$. The select operation $S(\cdot)$ contains the top- k selection $S_k(\cdot)$ and random selection $S_r(\cdot)$ implemented upon C . The $S_k(\cdot)$ selects the k largest elements of C and records their indexes \mathbf{T}_k . This encourages the attention mechanism to focus more on semantically significant areas and reduce the interference caused by insignificant parts. At the same time, the coloring of the background areas also needs reference. Thus $S_r(\cdot)$ randomly selects r more indexes \mathbf{T}_r . Finally, we obtain $S(C) = \mathbf{T}_k \cup \mathbf{T}_r$ and the semantic-sparse features $\hat{F}_{I_r}[S(C)]$. To calculate the new correspondence map, we can simply replace the features \hat{F}_{I_r} with $\hat{F}_{I_r}[S(C)]$ in Formula 2,3. The other steps remain the same as above.

Discussion

Dense Correspondence vs. Sparse Correspondence. Dense correspondence will be easily affected by irrelevant regions, especially when the reference is completely different from the gray-scale image. Even if the target region has low similarity with most reference regions, the re-weighting process will still disturb the final result. In contrast, sparse correspondence can overcome this difficulty by focusing only on semantically important regions, which can reduce the interference of other regions. Moreover, the computational complexity goes from $\mathcal{O}((HW)^2)$ to $\mathcal{O}((k+r)HW)$, while $(k+r)$ is generally 8 to 16 times smaller than HW . The comparison results of these two strategies are shown in Figure 3. It can be observed that some details are more accurately colorized after reducing the interference.

Coarse-colored vs. Gray-scale. In this work, we propose to use a coarse-colored image to build the correspondence with the reference, which is completely different from previous works. The coarse result is already consistent with the reference’s global color style, thus can produce more

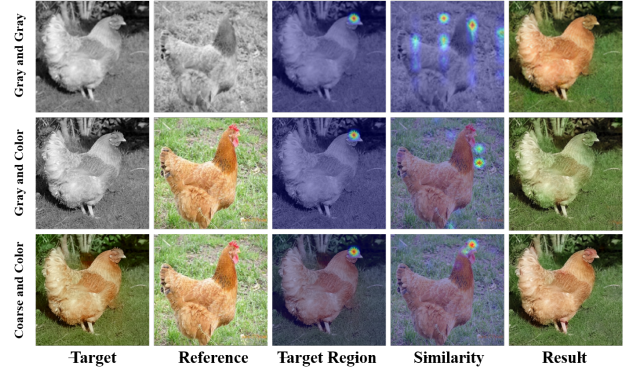


Figure 4: Comparison results of using three different data types to build the correspondence.

dedicated correspondence than directly using the gray-scale image. Moreover, the correspondence between color images is more accurate than that between gray-scale images (luminance channels). To verify this comment, we build a correlation matrix for three data types with the same operations. Figure 4 shows the comparison results of the similarity between one target region and all reference regions. It is clear that the chicken comb is correctly matched between two color images, even with different colors.

Objective Functions

Smooth-L1 Loss. To avoid simply using the average scheme for solving the ambiguity colorization problem, a widely used loss function Smooth-L1 loss is adopted in image colorization tasks. This loss is added to the results of both two stages in our architecture as L_{stage1} and L_{stage2} . The following Formula 4 can calculate the Smooth-L1 loss between T_{ab} and \hat{T}_{ab} :

$$L_{stage1,2}(T_{ab}, \hat{T}_{ab}) = \begin{cases} \frac{1}{2}(T_{ab} - \hat{T}_{ab})^2 & \text{for } |T_{ab} - \hat{T}_{ab}| \leq \delta \\ \delta |T_{ab} - \hat{T}_{ab}| - \frac{1}{2}\delta^2 & \text{otherwise.} \end{cases} \quad (4)$$

Classification Loss. There is a classification loss L_{cls} in the classifier to get a CAM as a reference for $S(\cdot)$. This loss can also improve the encoder’s ability of extracting color features. When F_{I_r} is fed into the classifier, its label vector is predicted. L_{cls} is defined as the cross-entropy between the classification vector \hat{y} and its ground truth one-hot label.

Color Histogram Loss. To transfer the color distribution of the reference image to the target image accurately, we also add a histogram loss to the final output as Formula 5. Similar to the previous work (Zhang, Isola, and Efros 2016), we treat the problem as multinomial classification. We quantify \hat{T}_{ab} output space into bins with $gridsize = 10$ and keep the in-gamut $Q = 313$. The mapping from the target image to predicted color distribution $\hat{Z} \in [0, 1]^{H \times W \times Q}$ is also learned with the decoder.



Figure 5: Qualitative comparison of colorizing results with previous methods. The target image, reference image, and each method’s colorized images are displayed from top to bottom. The proposed method outperforms other models and achieves state-of-the-art performance.

We also use the soft-encoding scheme (Zhang, Isola, and Efros 2016) to encode ground truth colors T_{ab} and get the ground truth distribution Z . Each real T_{ab} color value is represented as a convex combination of its ten nearest bin centers, weighted by a Gaussian kernel with $\sigma = 5$. We define L_{his} as cross-entropy loss for every pixel to measure the distance between predicted and ground truth distributions, and sum over all pixels.

$$L_{his}(\hat{Z}, Z) = - \sum_{h,w} \sum_q Z_{h,w,q} \log(\hat{Z}_{h,w,q}). \quad (5)$$

TV Regularization. To encourage spatial smoothness in the output result \hat{T}_{ab} , we follow previous work (Johnson, Alahi, and Fei-Fei 2016) and apply the total variation regularization $L_{TV}(\hat{T}_{ab})$ to the output of the fine colorization stage. In summary, the overall loss function for the entire network is defined as Formula 6:

$$L_{total} = \lambda_{stage1} L_{stage1} + \lambda_{stage2} L_{stage2} + \lambda_{TV} L_{TV} + \lambda_{cls} L_{cls} + \lambda_{his} L_{his}, \quad (6)$$

where λ_{stage1} , λ_{stage2} , λ_{TV} , λ_{cls} and λ_{his} are hyperparameters to constrain different loss terms.

Experiments

Implementation Details

To make the network capable of handling various colors and categories of images, we train our network on the ImageNet dataset (Deng et al. 2009), which has a training set of 1.3M images. We use ImageNet’s total training set to train the entire network with 5 epochs and set mini-batch size as 8. During training, the input image will be resized to 256×256 . We use Adam (Kingma and Ba 2015) for optimization with $\beta_1 = 0.9$, $\beta_2 = 0.999$. The learning rate of the whole network is set to 0.0001 during the entire training process. We set the coefficients for each loss function as follows: $\lambda_{stage1} = 100$, $\lambda_{stage2} = 100$, $\lambda_{cls} = 0.1$, $\lambda_{TV} = 10$, and $\lambda_{his} = 1$. For the $S(\cdot)$, both k and r are set to 256.

For the exemplar-based method, it is impossible to find enough source-reference pairs to train the network. Previous methods often utilize a self-augmented reference generated from the original image to train the network and take the original image as the ground truth.

We adopt a scheme similar to (Lee et al. 2020). The reference is generated from the original image by geometric distortion, which contains most of the contents from the original image itself, thereby providing complete information for

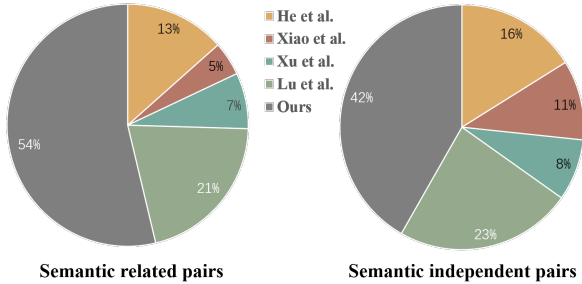


Figure 6: User preference results of five methods.

the target image. This kind of strategy encourages the network to be fully optimized and does not require manually annotated source-reference pairs.

The geometric distortion is realized by thin plate splines (TPS) transformation, a non-linear spatial transformation. The distortion is randomly applied to each image. This strategy prevents the model from lazily bringing the color in the same position from the reference while enforcing our model to build semantically accurate spatial correspondences between two images. To enable the model transfer the color effectively even without a proper reference, we deform some images violently during training, so that the augmented image has no semantic relationship with the source image. We simulate semantically unrelated image pairs for training purposes in this way.

Comparison with Previous Methods

Visual Comparison. We run all 5 models on 230 pairs of images collected from previous papers or the ImageNet validation set and show several representative results. We compare the results of our method with previous exemplar-based colorization approaches, including the color histogram-based approach (Xiao et al. 2020), stylization-based approach (Xu et al. 2020), and approaches that rely on correspondence (He et al. 2018; Lu et al. 2020). All comparison results are obtained by public available codes. We show the qualitative comparison in Figure 5.

The 4th column of Figure 5 shows the results of colorizing objects with unusual or artistic colors. Compared with methods (He et al. 2018) constrained by the perceptual loss, the proposed method can appropriately color the target image according to the user’s requirement. Since (Lu et al. 2020) tends to make the color histograms of the two images consistent, resulting in the wrong spatial distribution of colors.

In the 6th column, when there are large regions with less semantics in the image, our method can pay more attention to the semantically relevant areas, e.g., the pink area, while other methods fail to color the object or simply get a smooth result.

When the reference image is semantically unrelated to the target image (shown as 1st column in Figure 5), due to the dependence on prior color knowledge, (He et al. 2018) will ignore the colors from the reference image. Histogram-based methods (Xiao et al. 2020) can get plausible results

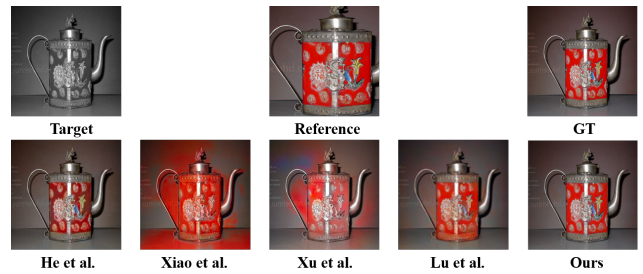


Figure 7: Comparison results of using RC reference image to color the target image.

by transferring global tones, whereas our method can yield better results. For some images with many details, (Lu et al. 2020) cannot properly color these details due to the inappropriate correspondence constructed with two gray-scale images, while the proposed method allows the target image to be colored correctly, e.g., 5th and 7th columns in Figure 5.

These experimental results show that the proposed method can construct a more accurate correspondence between target and reference and transfer color information for different image pairs effectively.

User Evaluation. We conduct user evaluation to verify the proposed method’s effectiveness subjectively. In this part, we randomly select 50 groups from the above results. Semantically dependent pairs and semantically unrelated pairs are distributed in half. Eventually, all 5×50 color images are distributed anonymously and randomly to 30 college participants.

For fairness, the images with the same reference are shown simultaneously in a random order. All participants were asked to observe the images for no more than 5 seconds and choose the image that better matches the reference. As shown in Figure 6, we show the percentage of votes for each method in the form of pie chart. It shows that images of our method are mostly preferred.

Self-Augmentation PSNR/SSIM. Unlike automatic colorization, in exemplar-based colorization setting, when given a target-reference pair, there is no ground truth that has both the target’s shape and the reference’s color. In order to make a quantitative evaluation of the colorization results, similar to the training process, we make an augmentation of the original image as the reference, so that the original image can be used as ground truth. With ground truth available for comparison, some existing evaluation metrics, such as peak signal-to-noise ratio (PSNR) and structural similarity (SSIM), can be used for evaluation.

We select 5000 images from the validation set of ImageNet to do three different data augmentation, including TPS, random rotation (RR), and random cropping (RC) as references to get different results. The quantitative comparisons of three different augmentation are reported in Table 1. Figure 7 shows an example of using a RC reference to color the target image and comparing the results with other methods. We will release this test dataset for future comparison.

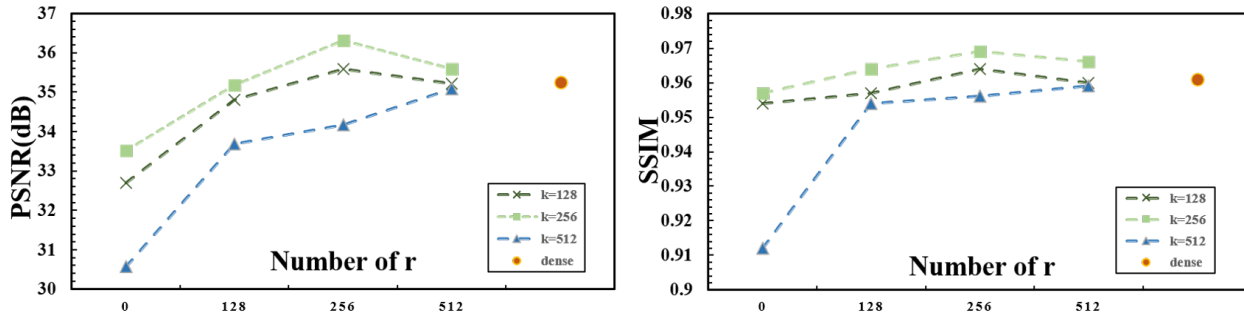


Figure 8: Ablation studies of $S_k(\cdot)$ and $S_r(\cdot)$.

Methods	TPS	RR	RC	Mean
(He et al. 2018)	28.51/0.902	28.67/0.903	27.57/0.898	28.25/0.901
(Xiao et al. 2020)	25.17/0.912	25.30/0.913	24.98/0.910	25.15/0.911
(Xu et al. 2020)	22.46/0.873	21.65/0.846	21.55/0.862	21.88/0.860
(Lu et al. 2020)	27.93/0.913	29.80/0.931	27.12/0.907	28.28/0.917
Ours	36.32/0.969	35.49/0.966	32.39/0.958	34.73/0.964

Table 1: Quantitative comparisons of self-augmentation PSNR/SSIM. A higher value indicates a better preference, while the proposed method outperforms other models.

Methods	PSNR	SSIM
dense w/o L_{cls}	33.73	0.952
dense	35.25	0.961
256-256 stage1	30.02	0.937
256-256-gray	33.02	0.955
256-256 w/o L_{his}	35.04	0.960
256-256	36.32	0.969

Table 2: Ablation studies of two-stage architecture and loss functions.

Ablation Studies

Ablation studies of $S_k(\cdot)$ and $S_r(\cdot)$. The use of sparse correspondence will lead to the question: how to select an appropriate number of regions in the process? Then we further study the effect of k and r and use TPS reference to evaluate results quantitatively as described above. When the resolution of the reference image is 256×256 , there are 4096 features available for selection. We increase k and r gradually from 128 and 0, respectively. The comparison results are shown in Figure 8. Without random selection, the value of PSNR/SSIM will be much lower because some areas of the background are incorrectly colored. Increasing r gradually can improve the results, but increasing r further will cause the result deteriorate again. In addition, it can be seen from the comparison of the three broken lines that a larger or smaller k will reduce the quality of the results.

Ablation studies of two-stage architecture. To illustrate the importance of the two-stage structure in our model, we conduct ablation studies on $k = 256, r = 256$ (256-256) version. The relevant results are shown in Table 2. First, we

evaluate the first stage results, and there is a large gap between them and the final results, thus illustrating the importance of LDT. To further validate the importance of preliminary coloring, we remove GCT from the whole architecture for comparison. Instead, we use another network with a similar structure to the encoder of I_r but with one channel input to extract the features of gray-scale image and calculate the correspondence in the same way (256-256-gray). Due to the lack of information in the gray-scale image, the PSNR/SSIM results will decrease.

Ablation studies of Loss Functions. In order to verify that the classifier cannot only provide a CAM but also help the encoder extract color features, we ablate the classification loss on the dense version. In addition, we also ablate color histogram loss of 256-256 version to analyze its effect. According to Table 2, removing either of these losses will reduce the model’s performance, especially in the L_{cls} .

Conclusions

This paper proposes a novel colorization framework named Semantic-Sparse Colorization Network (SSCN) to color the target image in a coarse-to-fine manner. Specifically, an image transfer network is adopted in the coarse colorization stage to obtain a preliminary colorized result. In the fine colorization stage, only the semantic-significant parts and some background regions of the reference image are reserved to get more accurate color details. Thus, SSCN can adequately transfer a reference image’s global color and local details onto a gray-scale image. It provides a way to obtain different levels of color information from the reference image hierarchically and accurately. Extensive experiments show that the proposed method outperforms previous state-of-the-art approaches by a large margin.

References

- Bahng, H.; Yoo, S.; Cho, W.; Park, D. K.; Wu, Z.; Ma, X.; and Choo, J. 2018. Coloring with Words: Guiding Image Colorization Through Text-Based Palette Generation. In *ECCV 2018*, 443–459. Springer.
- Bugeau, A.; Ta, V.; and Papadakis, N. 2014. Variational Exemplar-Based Image Colorization. *IEEE Trans. Image Process.*, 23(1): 298–307.

- Cao, Y.; Zhou, Z.; Zhang, W.; and Yu, Y. 2017. Unsupervised Diverse Colorization via Generative Adversarial Networks. In *ECML PKDD 2017*, 151–166. Springer.
- Charpiat, G.; Hofmann, M.; and Schölkopf, B. 2008. Automatic Image Colorization Via Multimodal Predictions. In *ECCV 2008*, 126–139. Springer.
- Cheng, Z.; Yang, Q.; and Sheng, B. 2015. Deep Colorization. In *ICCV 2015*, 415–423. IEEE Computer Society.
- Chia, A. Y. S.; Zhuo, S.; Gupta, R. K.; Tai, Y.; Cho, S.; Tan, P.; and Lin, S. 2011. Semantic colorization with internet images. *ACM Trans. Graph.*, 30(6): 156.
- Ci, Y.; Ma, X.; Wang, Z.; Li, H.; and Luo, Z. 2018. User-Guided Deep Anime Line Art Colorization with Conditional Adversarial Networks. In *MM 2018*, 1536–1544. ACM.
- Deng, J.; Dong, W.; Socher, R.; Li, L.; Li, K.; and Li, F. 2009. ImageNet: A large-scale hierarchical image database. In *CVPR 2009*, 248–255. IEEE Computer Society.
- Deshpande, A.; Lu, J.; Yeh, M.; Chong, M. J.; and Forsyth, D. A. 2017. Learning Diverse Image Colorization. In *CVPR 2017*, 2877–2885. IEEE Computer Society.
- Gupta, R. K.; Chia, A. Y. S.; Rajan, D.; Ng, E. S.; and Huang, Z. 2012. Image colorization using similar images. In *MM 2012*, 369–378. ACM.
- He, M.; Chen, D.; Liao, J.; Sander, P. V.; and Yuan, L. 2018. Deep exemplar-based colorization. *ACM Trans. Graph.*, 37(4): 47:1–47:16.
- Huang, Y.; Tung, Y.; Chen, J.; Wang, S.; and Wu, J. 2005. An adaptive edge detection based colorization algorithm and its applications. In *MM 2005*, 351–354. ACM.
- Johnson, J.; Alahi, A.; and Fei-Fei, L. 2016. Perceptual Losses for Real-Time Style Transfer and Super-Resolution. In *ECCV 2016*, 694–711. Springer.
- Kingma, D. P.; and Ba, J. 2015. Adam: A Method for Stochastic Optimization. In *ICLR 2015*.
- Kumar, M.; Weissenborn, D.; and Kalchbrenner, N. 2021. Colorization Transformer. *CoRR*, abs/2102.04432.
- Lee, J.; Kim, E.; Lee, Y.; Kim, D.; Chang, J.; and Choo, J. 2020. Reference-Based Sketch Image Colorization Using Augmented-Self Reference and Dense Semantic Correspondence. In *CVPR 2020*, 5800–5809. IEEE Computer Society.
- Levin, A.; Lischinski, D.; and Weiss, Y. 2004. Colorization using optimization. *ACM Trans. Graph.*, 23(3): 689–694.
- Liao, J.; Yao, Y.; Yuan, L.; Hua, G.; and Kang, S. B. 2017. Visual attribute transfer through deep image analogy. *ACM Trans. Graph.*, 36(4): 120:1–120:15.
- Lu, P.; Yu, J.; Peng, X.; Zhao, Z.; and Wang, X. 2020. Gray2ColorNet: Transfer More Colors from Reference Image. In *MM 2020*, 3210–3218. ACM.
- Luan, Q.; Wen, F.; Cohen-Or, D.; Liang, L.; Xu, Y.; and Shum, H. 2007. Natural Image Colorization. In *Proceedings of the Eurographics Symposium on Rendering Techniques 2007*, 309–320. Eurographics Association.
- Manjunatha, V.; Iyyer, M.; Boyd-Graber, J. L.; and Davis, L. S. 2018. Learning to Color from Language. In *NAACL-HLT 2018*, 764–769. Association for Computational Linguistics.
- Qu, Y.; Wong, T.; and Heng, P. 2006. Manga colorization. *ACM Trans. Graph.*, 25(3): 1214–1220.
- Ronneberger, O.; Fischer, P.; and Brox, T. 2015. U-Net: Convolutional Networks for Biomedical Image Segmentation. In *MICCAI 2015*, volume 9351 of *Lecture Notes in Computer Science*, 234–241. Springer.
- Sangkloy, P.; Lu, J.; Fang, C.; Yu, F.; and Hays, J. 2017. Scribbler: Controlling Deep Image Synthesis with Sketch and Color. In *CVPR 2017*, 6836–6845. IEEE Computer Society.
- Simonyan, K.; and Zisserman, A. 2015. Very Deep Convolutional Networks for Large-Scale Image Recognition. In *ICLR 2015*.
- Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, L.; and Polosukhin, I. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017*, 5998–6008.
- Xiao, C.; Han, C.; Zhang, Z.; Qin, J.; Wong, T.; Han, G.; and He, S. 2020. Example-Based Colourization Via Dense Encoding Pyramids. *Comput. Graph. Forum*, 39(1): 20–33.
- Xu, K.; Li, Y.; Ju, T.; Hu, S.; and Liu, T. 2009. Efficient affinity-based edit propagation using K-D tree. *ACM Trans. Graph.*, 28(5): 118.
- Xu, Z.; Wang, T.; Fang, F.; Sheng, Y.; and Zhang, G. 2020. Stylization-Based Architecture for Fast Deep Exemplar Colorization. In *CVPR 2020*, 9360–9369. IEEE.
- Yoo, S.; Bahng, H.; Chung, S.; Lee, J.; Chang, J.; and Choo, J. 2019. Coloring With Limited Data: Few-Shot Colorization via Memory Augmented Networks. In *CVPR 2019*, 11283–11292. Computer Vision Foundation / IEEE.
- Zhang, R.; Isola, P.; and Efros, A. A. 2016. Colorful Image Colorization. In *ECCV 2016*, 649–666. Springer.
- Zhang, R.; Zhu, J.; Isola, P.; Geng, X.; Lin, A. S.; Yu, T.; and Efros, A. A. 2017. Real-time user-guided image colorization with learned deep priors. *ACM Trans. Graph.*, 36(4): 119:1–119:11.
- Zhou, B.; Khosla, A.; Lapedriza, À.; Oliva, A.; and Torralba, A. 2016. Learning Deep Features for Discriminative Localization. In *CVPR 2016*, 2921–2929. IEEE Computer Society.