PalGAN: Image Colorization with Palette Generative Adversarial Networks Supplementary Material

Yi Wang¹, Menghan Xia², Lu Qi³, Jing Shao⁴, and Yu Qiao¹

 $^1 \mathrm{Shanghai}$ AI Laboratory $\ ^2 \mathrm{Tencent}$ AI Lab $\ ^3 \mathrm{UC}$ Merced $\ ^4 \mathrm{SenseTime}$ Research

We describe the following contents in this supplementary file:

- Specific design of our given PalGAN.
- More colorization results on ImageNet, COCOStuff, and real-world legacy images.
- The discussion about limitations and failure cases.

Code and models will be released at https://github.com/shepnerd/PalGAN.

1 Network Architectures

PalGAN contains Palette Generator, Palette Assignment Generator, and Color Discriminator. We detail their topology and configuration below. For clarity, we denote Conv(k, s, c) indicates a convolutional operation whose kernel size, stride size, and output channel number of the used convolution are k, s, and c, respectively. The dilation ratio and padding size of Conv(k, s, c) are both set to 1. $FC(c_{in}, c_{out})$ denotes a fully-connected operation where c_{in} and c_{out} stands for input channel and output one, respectively. \uparrow and \otimes denote a 2× bilinear upsampling and concatenation (along with the channel dimension) operations, respectively. ResBlock(k, s, c) denotes a residual block as used in [3,5].

 $\begin{array}{l} Palette \ Generator : \mathbf{L} \rightarrow \operatorname{Conv}_1(3,2,32) \rightarrow \operatorname{LReLU} \rightarrow \operatorname{Conv}_2(3,2,64) \rightarrow \operatorname{LReLU} \\ \rightarrow \operatorname{Conv}_3(3,2,128) \rightarrow \operatorname{LReLU} \rightarrow \operatorname{Conv}_4(3,2,256) \rightarrow \operatorname{LReLU} \rightarrow \operatorname{Conv}(3,2,256) \rightarrow \\ \operatorname{LReLU} \rightarrow \operatorname{Conv}(3,2,256) \rightarrow \operatorname{LReLU} \rightarrow \operatorname{FC}(4096,256) \rightarrow \operatorname{Sigmoid} \rightarrow \mathbf{\hat{h}}. \end{array}$

 $\begin{array}{l} Palette \ Assignment \ Generator : \mathbf{L} \downarrow \otimes \mathbf{M}_4 \rightarrow \operatorname{Conv}_1(3,2,512) \rightarrow \operatorname{LReLU} \rightarrow \\ \operatorname{ResBlock}(3,1,512) \rightarrow \operatorname{ResBlock}(3,1,512) \rightarrow \uparrow [\otimes \mathbf{M}_3] \rightarrow \operatorname{ResBlock}(3,1,256) \rightarrow \\ \operatorname{ResBlock}(3,1,256) \ (\mathbf{S}) \rightarrow \uparrow [\otimes \mathbf{M}_2] \rightarrow \operatorname{ResBlock}(3,1,128) \rightarrow \\ \operatorname{ResBlock}(3,1,256) \ (\mathbf{S}) \rightarrow \uparrow [\otimes \mathbf{M}_2] \rightarrow \operatorname{ResBlock}(3,1,128) \rightarrow \\ \operatorname{ResBlock}(3,1,128) \rightarrow \\ \operatorname{ResBlock}(3,1,64) \rightarrow \\ \operatorname{CA}(\mathbf{S},\mathbf{L}) \rightarrow \uparrow \rightarrow \\ \operatorname{Conv} (3,1,2) \rightarrow \\ \operatorname{tanh} \rightarrow \mathbf{\hat{I}}, \end{array}$

where \mathbf{M}_i indicates the output from Conv_i in Palette Generator.

Color Discriminator : $\hat{\mathbf{C}} \oplus \hat{\mathbf{I}}$ or $\mathbf{C} \oplus \mathbf{I} \to \text{Conv}(3,2,32) \to \text{LReLU} \to \text{Conv}(3,2,64)$ $\to \text{LReLU} \to \text{Conv}(3,2,128) \to \text{LReLU} \to \text{Conv}(3,2,256) \to \text{LReLU} \to \text{Conv}(3,2,256)$ $\to \text{LReLU} \to \text{FC}(128,256) \to \text{proj}(\mathbf{h}) \to \text{T or F},$ where proj is the conditional projection (Eqn. 9) in the paper. 2 Y. Wang et al.

1.1 Model Capacity

The inference part of PalGAN (Palette Generator+Palette Assignment Generator) has only 22M parameters, smaller than that of UGColor [9] (34M), ColTrans [4] (74M), GPColor [8] (\geq 71M), and InstColor [6] (84M).

2 More Experimental Results and Analysis

We give more visual comparisons on ImageNet [2] (Figure 1), and COCO-Stuff [1] (Figure 2). Also, more reference-based results (Figure 3), and legacy photo colorization (Figure 4) are presented.

2.1 Limitations and Failure Cases

As discussed in the paper, PalGAN cannot handle small objects or structures (especially for high-resolution inputs) as chromatic attention is incapable of well representing these targets with small-scale feature maps, as shown in Figure 5. For crowds (top row in Figure 5) or small-scale buildings in the background (bottom row in Figure 5), all the employed methods fail to paint consistent and delicate colors according to the semantics.

2.2 Future Work

In the current design, we explicitly disentangle the color distribution from a spatial domain. Undoubtedly, predicting palettes locally will grant more controllability to users as we can regionally change the color distribution. It remains one of future works. Besides, considering the mentioned limitations on smallscale object colorization, using multi-scale attention mechanism or maintaining high-resolution representation as in [7] might address it, at the cost of large computational consumption.



Fig. 1: Visual comparison on ImageNet.

Input Automatic Colorization Reference-based Colorization Fig. 3: Our reference-based colorization results.

Fig. 4: Our legacy photo colorization results. First row: legacy images, and second row: our predictions.

(1) Input (2) Deoldify (3) UGColor (4) InstColor (5) Ours Fig. 5: Failure cases. The given methods are unable to handle small objects and structures in the give inputs, leading to obvious color bleeding and inconsistency.

6 Y. Wang et al.

References

- Caesar, H., Uijlings, J., Ferrari, V.: Coco-stuff: Thing and stuff classes in context. In: CVPR. pp. 1209–1218 (2018) 2
- Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: CVPR. pp. 248–255 (2009) 2
- Isola, P., Zhu, J.Y., Zhou, T., Efros, A.A.: Image-to-image translation with conditional adversarial networks. In: CVPR. pp. 1125–1134 (2017) 1
- Kumar, M., Weissenborn, D., Kalchbrenner, N.: Colorization transformer. arXiv preprint arXiv:2102.04432 (2021) 2
- Park, T., Liu, M.Y., Wang, T.C., Zhu, J.Y.: Semantic image synthesis with spatiallyadaptive normalization. In: CVPR. pp. 2337–2346 (2019) 1
- Su, J.W., Chu, H.K., Huang, J.B.: Instance-aware image colorization. In: CVPR. pp. 7968–7977 (2020) 2
- Wang, J., Sun, K., Cheng, T., Jiang, B., Deng, C., Zhao, Y., Liu, D., Mu, Y., Tan, M., Wang, X., et al.: Deep high-resolution representation learning for visual recognition. IEEE transactions on pattern analysis and machine intelligence (2020) 2
- Wu, Y., Wang, X., Li, Y., Zhang, H., Zhao, X., Shan, Y.: Towards vivid and diverse image colorization with generative color prior. In: ICCV. pp. 14377–14386 (2021) 2
- Zhang, R., Zhu, J.Y., Isola, P., Geng, X., Lin, A.S., Yu, T., Efros, A.A.: Realtime user-guided image colorization with learned deep priors. arXiv preprint arXiv:1705.02999 (2017) 2