

1 Supcode

Algorithm: The Framework of Style Image Harmonization via Global-Local Style Mutual Guided

Input: Source image X_s , Composite image X_c , Mask image X_m

1. Style transfer: For each iteration: $L_s, L_c \leftarrow \text{Linear Projection}(X_s, X_c)$ <i>// L_s, L_c are style and content patch sequence respectively</i> $H_s \leftarrow \text{ViT Encoder}(L_s, L_m)$ <i>// H_s is output of ViT Encoder</i> For i in range(5,1,-1): VGG Encoder E_i is VGG Relu_{}_1.format(i) $sF, cF \leftarrow E_i(X_s, X_c)$ <i>// sF, cF are encoded features of style and content</i> Recons Decoder D_i is Feature_invertor_{}_1.format(i) $X_s \leftarrow D_i(WCT(sF, cF))$ <i>// X_s output of D_i after WCT operation</i> <i>// X_s and X_c are used as inputs for the next layer</i> $F \leftarrow WCT(E_i(X_s, X_c))$ <i>// F is output of WCT of last layer</i> $I \leftarrow \text{Blending Decoder}(F, H_s)$ <i>// I is output of Blending Decoder</i> <i>// The Blending Decoder structure is shown in Fig. 3</i> $X_{out} \leftarrow X_m \odot I + (1 - X_m) \odot X_s$ <i>// X_{out} is final output</i> Update $L_{transfer}$ in Eq. 7 <i>// $L_{transfer} = \lambda_2 L_s(X_{out}, X_c) + \lambda_2 L_c(X_{out}, X_s) + \lambda_{reg} L_{reg}$ set $\lambda_2 = 7, \lambda_s = 10, \lambda_{reg} = 10$</i> <i>// $L_{reg} = L_2(X_s, Tr(X_s, X_c, X_m)) + L_2(X_c, Tr(X_s, X_c, X_m))$, Tr is style transfer</i>	2. GradGAN: For each Generator iteration: $F \leftarrow \text{Encoder}(X_s, X_c)$ <i>// F is latent code after encoding</i> $I_{fake} \leftarrow \text{Decoder}(F)$ <i>// I_{fake} is output of Generator</i> Update $L_G = \lambda_{adv} L_{adv}^G(X_s, X_c) +$ $\lambda_{L2}(I_{fake}, X_c) + \lambda_{grad} L_{grad}(I_{fake}, X_m)$ <i>// calculate the L_{grad} inside the dilated mask</i> <i>// set $\lambda_{adv} = 1, \lambda_2 = 10, \lambda_{grad} = 10$</i> For each Patch Discriminator D iteration: Update $L_D = L_{adv}^D(I_{fake}, X_s)$
3. Fusion: <i>// First train 1 and 2, then perform 3</i> Given: Trained transfer module Tr , GradGAN Gr $X_{style} \leftarrow Tr(X_s, X_c, X_m), X_{grad} \leftarrow Gr(X_s, X_c, X_m)$ Update $T(X_c)$ in Eq. 5. <i>// $T(X_c) = \omega S(X_c, X_{grad}) + \phi G(X_c, X_{grad})$</i> <i>// S is style constraint, G is gradient constraint</i> <i>// set $\omega=1, \phi=1$</i>	

Fig. 1. The supcode contains three components: Style transfer, GranGAN and Fusion. First, we train global and local mutual guided style transfer, so that the source cropped area has the target style while more in line with its semantics. Second, GradGAN makes the stylized areas smoother in the gradient at the paste boundary. Finally, the gradient and style are fused using style and gradient constraint.

2 Supplementary Experiment

We show the experimental figures that cannot be shown or are difficult to be clearly shown due to the length of this paper.

Fig. 2 shows the results of the ablation experiment. We enlarge the arrangement to facilitate readers to view it more clearly. Supplementary notes: examples of results with different degrees of texture information richness are shown. The first two columns in the figure are example results with less texture, and the last two columns are rich texture example results.

Fig. 3 - Fig. 6 show the comparison results between our model and the combination models. Four image harmonization models and four style transfer models are used separately, as well as their combinations.

Some models even change the image content, such as the performance of WCT+BargainNet or WCT+Dovenet on the characters in Fig. 4. Some combinations cause serious artifacts or abrupt results, such as the fire truck in Fig. 4 and the balloon in Fig. 6. We show the superior effect on the clock in the first row of Fig. 6, while other combinations are not ideal. Although our method is not perfect in some details, like the girl's shoes in Fig. 6, in general, our method can better deal with the task of style image harmonization.

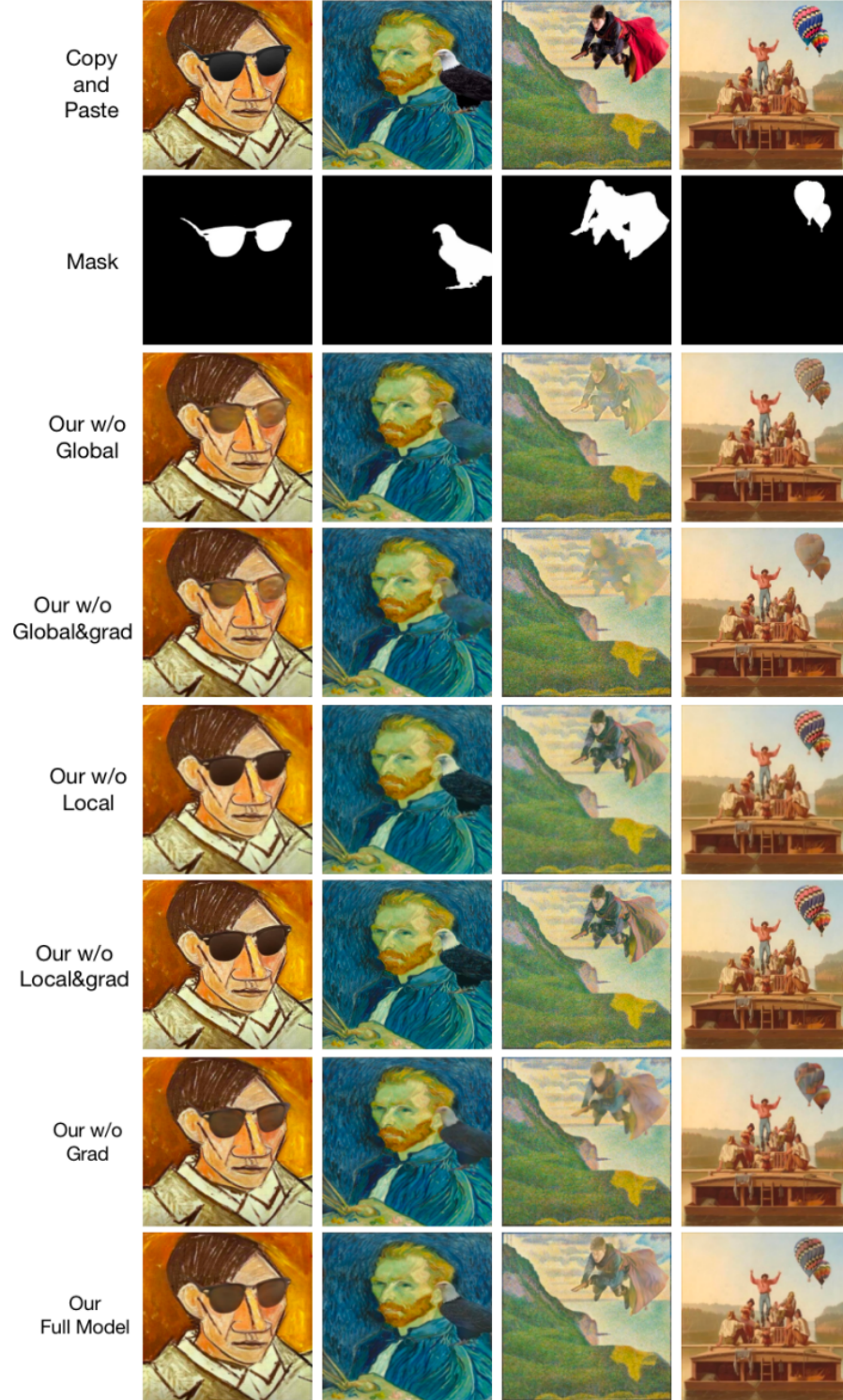


Fig. 2. Experimental study of ablation of global, local, and gradient fusion. The full model shows the best result, however other baseline models either distort the original information or are not in harmony with the surrounding, or there are problems such as sudden change in the gradient.



Fig. 3. The result of style transfer or image harmonization models.

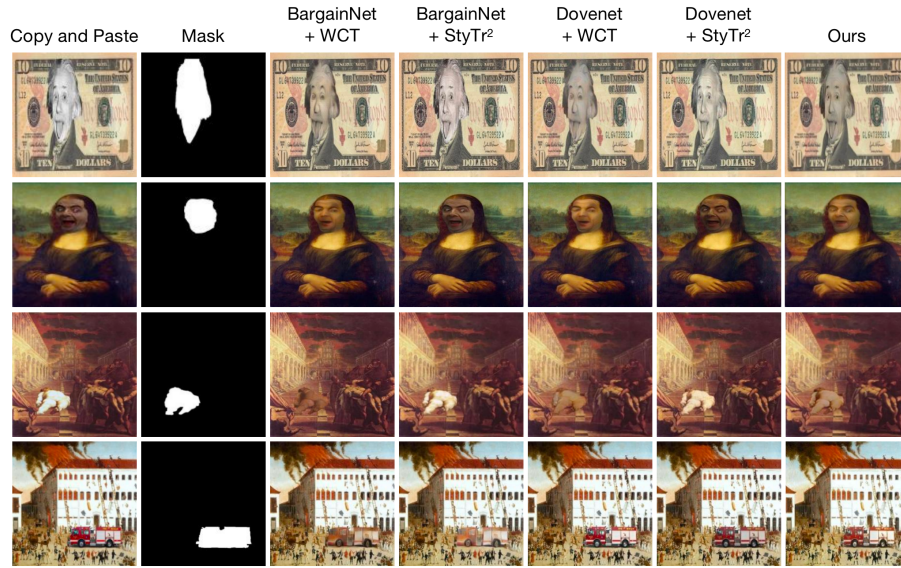


Fig. 4. Compare our model and combination models. Some models even change the image content, such as the performance of WCT+BargainNet or WCT+Dovenet on the characters. Some combinations cause serious artifacts or abrupt results, like the fire truck.

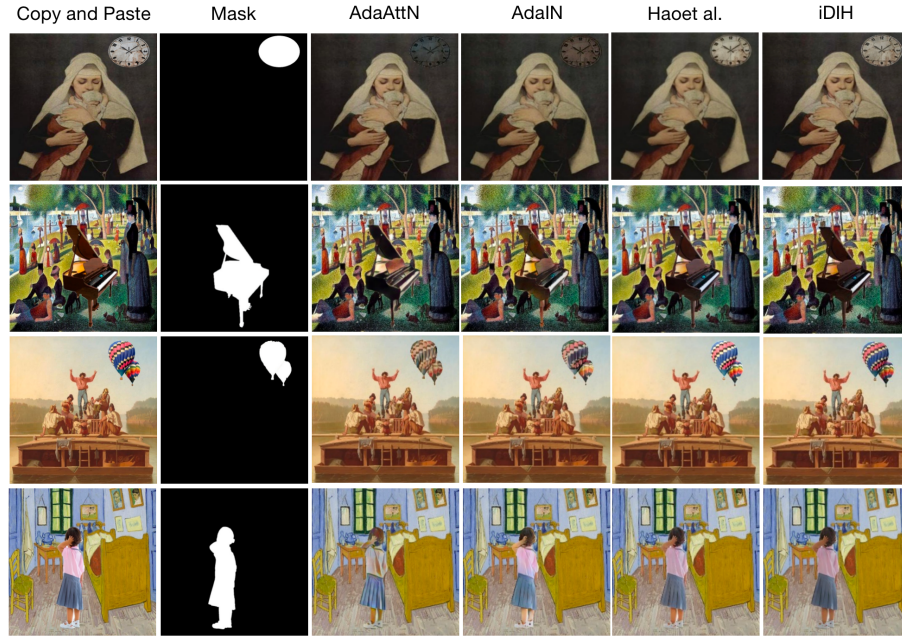


Fig. 5. The result of style transfer or image harmonization models.

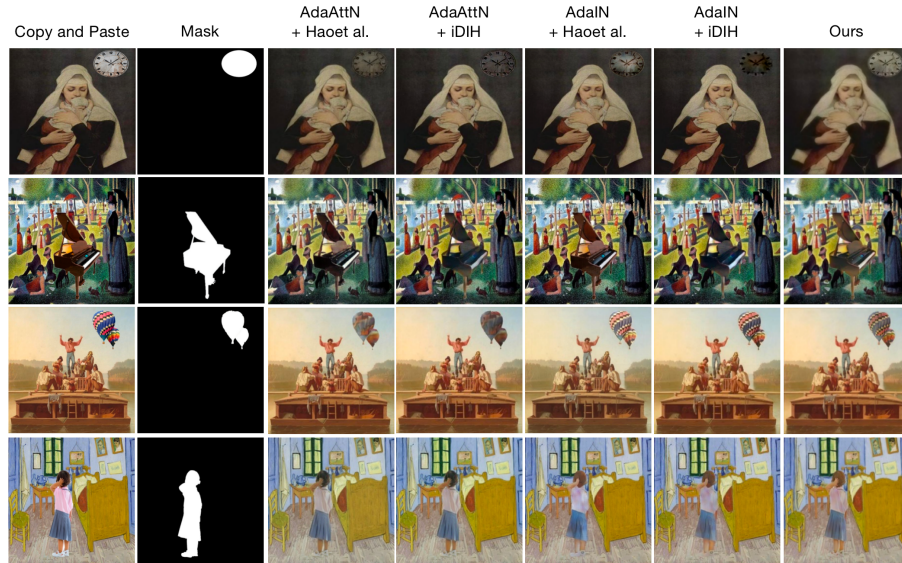


Fig. 6. Compare our model and combination models. Some combinations cause serious artifacts or abrupt results, like the balloon. We show the superior effect on the clock in the first row, while other combinations are not ideal.