

# Towards High-resolution and Disentangled Reference-based Sketch Colorization

## Supplementary Material

### 001 1. Plugin module and injection

002 We use the plugin module to extract low-level features from  
003 the background latent and inject them into the denoising  
004 backbone through split attention injections, which means  
005 we give different KV inputs for foreground and background  
006 in self-attention and cross-attention layers when activating  
007 the plugin module. The split self-attention injection can be  
008 formulated as:

$$009 y_{sa} = \begin{cases} \text{Attention}(z_f, z, z) & \text{if } m > ts \text{ (Foreground)} \\ \text{Attention}(z_b, z_{inj}, z_{inj}) & \text{if } m \leq ts \text{ (Background)}, \end{cases} \quad (1)$$

010 where  $z, z_f, z_b, z_{inj}, m, ts$  are the forward features, fore-  
011 ground tokens of forward features, background tokens of  
012 forward features, and injection features from the plugin  
013 module, the sketch foreground mask, and the user-defined  
014 foreground threshold, respectively. Users can adjust  $ts$  to  
015 adjust the selected foreground regions or manually give a  
016 foreground mask. We perform the same split for the cross-  
017 attention. The only difference is that we jointly train back-  
018 ground LoRAs for the QKV projection layers used in cross-  
019 attention layers. And the split cross-attention is organized  
020 as:

$$021 y_{ca} = \begin{cases} \text{Softmax}\left(\frac{W^q z_f \cdot (W^k e)}{\sqrt{d}}\right)(W^v e) & \text{if } m > ts \\ \text{Softmax}\left(\frac{\hat{W}_b^q z_b \cdot (\hat{W}_b^k e_b)}{\sqrt{d}}\right)(\hat{W}_b^v e_b) & \text{if } m \leq ts. \end{cases} \quad (2)$$

022 Here,  $e, e_b$  are embeddings from reference images and  
023 reference background embeddings, respectively, and  $\hat{W}_b^*$   
024 denotes the LoRA-injected linear weight, with the LoRA  
025 rank set to 16.

026 Both split attentions are implemented through **masked**  
027 **attention**.

### 028 2. Failure case and solution

#### 029 2.1. Description

030 As discussed in the main text, the proposed Gram loss suc-  
031 cessfully disentangles spatial semantics by ensuring that the  
032 structure of the sketch-guided region (defined as the fore-  
033 ground) is derived exclusively from the sketch inputs. How-  
034 ever, we observe limitations stemming from **biased train-**  
035 **ing data**. Specifically, when a reference image lacks back-  
036 ground details, the network tends to hallucinate random  
037 content in the non-sketch regions. These failure cases are  
038 visualized in Figure 1.



Figure 1. Sketch images inherently possess rich background spatial semantics, as evidenced by our sketch-only results where no reference images are provided. By incorporating Gram loss, the network is guided to more faithfully follow the spatial semantics of the sketches. Consequently, the model tends to synthesize backgrounds when the reference image offers negligible background spatial semantics, reflecting the fact that the majority of sketch images in our training data contain background details.

### 039 2.2. Solution and discussion

040 This background issue can be mitigated by **activating the**  
041 **plugin module**, as it enables low-level feature transfer for  
042 the image. In future work, we aim to address the root cause  
043 by refining our training data and strategy for robust T2I gen-  
044 eration.

045 Crucially, this challenge is distinct from the spatial en-  
046 tanglement addressed in prior works [2, 3], which stems  
047 from foreground embedding leakage and manifests as body  
048 overflow and segmentation deterioration, visualized in the  
049 w/o Gram loss column in Figure 1. We observe that such  
050 entanglement artifacts persist—even with split-attention [2]  
051 or our proposed plugin module—at the high resolutions  
052 ( $>1024\text{px}$ ) used in this study. This persistence implies a  
053 fundamental limitation in the network’s ability to differen-  
054 tiate between foreground and background tokens.

### 055 3. More baseline for cross-content colorization comparison

056 Nano Banana is also applicable to reference-based sketch  
057 colorization. While it struggles with character sketches  
058 due to imprecise regional color transfer, it proves effec-  
059 tive for cross-content colorization. In Figure 2 and 3, we  
060 present additional qualitative results for Nano Banana and  
061 IP-Adapter. These serve as extended baselines for the com-  
062 parisons shown in Figure 9 of the main manuscript.  
063



Figure 2. Cross-content results using Nano Banana.



Figure 3. Cross-content results using ControlNet + IP-Adapter [4, 5] with style transfer weights [1].

064 **4. User Study**

065 In the user study, we present 25 sets of sketch-reference  
 066 pairs and results generated by all compared methods to  
 067 users investigated. We show the user interface in Figure 4,  
 068 and the images shown in the user study in Figure 5, Fig-  
 069 ure 6, and Figure 7.

070 **References**

- 071 [1] Haofan Wang, Qixun Wang, Xu Bai, Zekui Qin, and Anthony  
 072 Chen. Instantstyle: Free lunch towards style-preserving in  
 073 text-to-image generation. *arXiv preprint arXiv:2404.02733*,  
 074 2024. 2
- 075 [2] Dingkun Yan, Xinrui Wang, Zhuoru Li, Suguru Saito, Yusuke  
 076 Iwasawa, Yutaka Matsuo, and Jiaxian Guo. Image referenced  
 077 sketch colorization based on animation creation workflow. In  
 078 *Proceedings of the Computer Vision and Pattern Recognition*  
 079 *Conference*, pages 23391–23400, 2025. 1
- 080 [3] Dingkun Yan, Liang Yuan, Erwin Wu, Yuma Nishioka, Is-  
 081 sei Fujishiro, and Suguru Saito. Colorizediffusion: Improv-

ing reference-based sketch colorization with latent diffusion  
 model. In *Proceedings of the Winter Conference on Applica-  
 tions of Computer Vision (WACV)*, pages 5092–5102, 2025. 1

- [4] Hu Ye, Jun Zhang, Sibio Liu, Xiao Han, and Wei Yang. Ip-  
 adapter: Text compatible image prompt adapter for text-to-  
 image diffusion models. *CoRR*, abs/2308.06721, 2023. 2
- [5] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding  
 conditional control to text-to-image diffusion models. In  
*ICCV*, pages 3836–3847, 2023. 2

082  
 083  
 084  
 085  
 086  
 087  
 088  
 089  
 090  
 091



Figure 4. The user interface employed during the perceptual study. Participants were presented with input line art and a reference color sample, then asked to evaluate the generated results based on structural fidelity, style consistency, and overall visual quality.



Figure 5. Images shown during the user study.



Figure 6. Images shown during the user study.



Figure 7. Images shown during the user study.