

1 We sincerely thank our reviewers for the constructive feedback. We are glad to see our paper is well-written (R1, R3,
2 R4), has strength in visual quality (R1, R3, R4), and performs diverse (R1) and thorough (R3, R4) experiments.

3 **Novelty [R1, R3]** Note that we do not claim to be the first to use swapping for disentanglement. Our key contribution is
4 in the way we *constrain* this disentanglement, by separating intra-image information (co-occurrence of patches within
5 the same image) from inter-image information (visual content across images within the dataset). Typically, learning for
6 image synthesis happens by gathering information across the entire image dataset (or across an image class). However,
7 this is difficult to do correctly, as the texture statistics within each image can be quite different and incompatible. Thus
8 the recent popular “internal learning” methods, such as InGAN[58] or SinGAN[59], show superior results by learning
9 from the patches of just the input image, utilizing the commonality of texture statistics. But, of course, methods like
10 [59] can’t learn that much because their dataset is very limited (just one image!). What our proposed method is doing
11 is disentangling the texture-like information useful for internal learning from the structure-like information that can
12 profit from external learning. To separate out internal learning from external learning, we model internal information
13 (texture code) using patch co-occurrence D (which operationalizes Julesz’s 2nd order texture conjecture), which leaves
14 the structure code representing texture-invariant (compatible with texture code of any image). So, we combine the
15 best of both internal and external learning. The swapping is just a self-supervised pretext task to accomplish this
16 disentanglement. Figure 6 confirms that the quality of disentanglement is superior to generative models that don’t
17 utilize internal information (Im2StyleGAN[1], StyleGAN2[38]) or external information (STROTSS[41], WCT²[69]).
18 In the current draft, this motivational story got distributed throughout the paper, but in the next draft we will be sure to
19 summarize it concisely in the intro as well.

20 **Comparison to FUNIT [R3]** FUNIT is fundamentally different from our work; it requires labels to define the
21 disentanglement of style from structure. For example, the AnimalFaces, Birds, Flowers, and Food datasets contain
22 149, 444, 85, and 224 classes, respectively. *For this reason, FUNIT cannot be trained on any of the datasets in our*
23 *paper*. In contrast, our method is fully unsupervised, and works because co-occurrence statistics of patches of a single
24 image carry enough signal for learning a smooth embedding space for image editing. By doing so, our method can
25 be easily trained on user-collected datasets such as the Flickr Mountains and Waterfalls, which is useful for long-tail
26 image editing. In more detail, the learning objective and discriminator design are substantially different; the FUNIT
27 discriminator requires class labels, while ours performs internal learning on unlabeled images (see Novelty). While
28 substantially different to FUNIT’s discriminator in concept, we agree our discriminator shares architectural similarities
29 with FUNIT’s encoder. To clarify, the PatchEncoder of Figure 2 of SuppMat is the feature extractor of our discriminator,
30 not our encoder. We apologize for any confusion.

31 **Image Quality, Resolution, and Versatility [R1]** We show our method works well on a wider variety of datasets
32 than most GAN-based methods, on human and animal faces, outdoor and indoor scenes, user-collected datasets, and
33 paintings (Section 1.3 of Supp Mat). Moreover, our method supports HD resolution (e.g. the mountain example
34 of Figure 7, the demo video and Supp Mat). Lastly, we evaluated the quality of texture transfer against two SOTA
35 style transfer methods (non-photorealistic and photorealistic), and show that our method is competitive, if not better.
36 Therefore, we believe our method has advantage over general texture transfer methods.

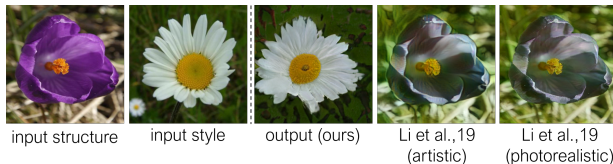
37 **Cite and compare to disentanglement / neural style transfer papers [R1]** Thank you for the suggestions. We will
38 cite the suggested works. Our end goal differs in that a real, existing image must be faithfully embedded for editing,
39 while the suggested works do not focus on accurately reconstructing input images.

40 **Comparison to StyleGAN2 [R3]** Our method aims to build an image editing model in an unsupervised setting. To
41 our knowledge, the state-of-the-art unsupervised methods for GANs-based image editing are Im2StyleGAN[1] and
42 StyleGAN2[38]. Randomly-generated StyleGAN2 images cannot be used for this application.

43 **Injection [R4]** We appreciate your feedback. We will modify “clearly injective (L119)” to “optimized toward injection”,
44 and also revise the Supp Map accordingly.

45 **About Broader Impact [R4]** We are happy to move the figure to Supp Mat based on reviewers/AC’s recommendation.

46 **Structural Limitation [R3]** We show the result of the suggested lily / daisy example. Please see the figure below.



47
48
49
50
51
52
53 Structure-style mixing results on the Oxford Flower dataset. Our
54 model can modify fine-level structure if the texture code requires it,
55 while maintaining the coarse-level structure. We additionally show
56 the result of a real-time style transfer method (Li et al., CVPR2019
57 [A6 suggested by R4], both artistic and photorealistic version) that
58 makes only small changes.

46 **Comparison to recent style transfer models [R4]**

47 Thank you for the suggestion; we will discuss our work in
48 the context of neural style transfer more explicitly. Note
49 that style transfer methods, even if real-time, are not suit-
50 able for *natural* image manipulation, because they often
51 fail to produce photorealistic images (please see left) and
52 do not allow controllable structural editing. For thorough-
53 ness, in the submission, we compared the quality of disen-
54 tlement to both SOTA non-photorealistic (STROTSS)
55 and photorealistic (WCT²) style transfer methods. While
56 WCT² does achieve better structural preservation, it does
57 so by barely changing the style (Figure 5).