

---

# PHOTOSWAP: Personalized Subject Swapping in Images — Supplementary Materials

---

Anonymous Author(s)  
Affiliation  
Address  
email

## 1 A Results on Real Images

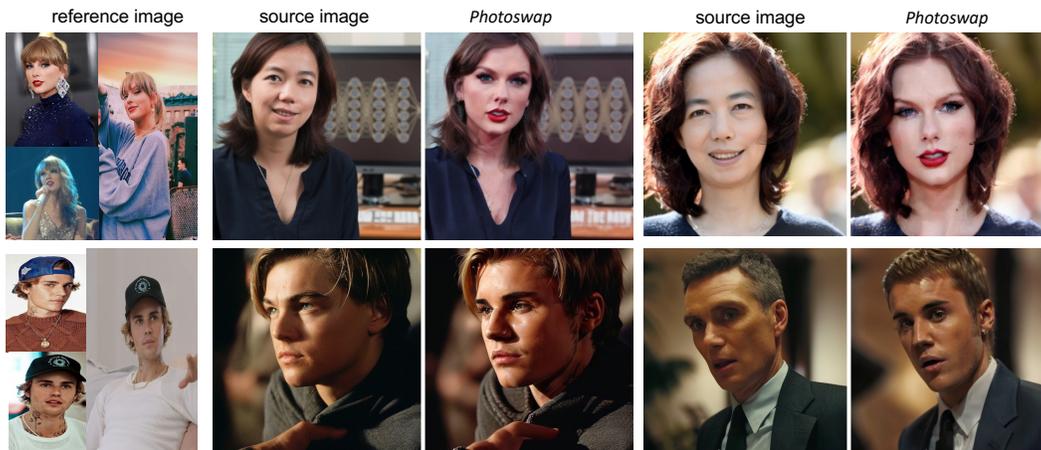


Figure 1: Face swapping on real images. Top row are real photos of Dr. Fei-Fei Li, and bottom row are real movie scenes. *Photoswap* precisely swaps human subjects while preserving the same subject pose and background context.

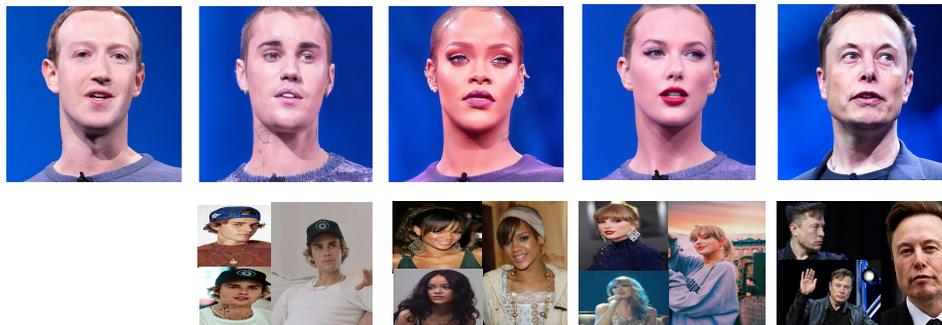


Figure 2: More results on human face swapping. Given a source human image (leftmost), *Photoswap* replaces the human identity with the reference person (bottom row) to generate the target image (top row).

2 This section showcases the practical effectiveness of our *Photoswap* method on real images in  
3 Figure 1, Figure 2, and Figure 3. These results provide a visual testament to the successful execution



Figure 3: Results of replacing Vincent Van Gogh with Elon Musk in Van Gogh’s famous self-portrait painting. *Photoswap* is able to maintain the same pose with a decent style preservation.

4 of subject swapping in real-world instances. For the implementation of subject swapping on actual  
 5 images, we require an additional process that utilizes an image inversion method, specifically the  
 6 DDIM inversion [7], to transform the image into initial noise. This inversion method relies on  
 7 a reversed sequence of sampling to achieve the desired inversion. However, there exist inherent  
 8 challenges when this inversion process is applied in text-guided synthesis within a classifier-free  
 9 guidance setting. Notably, the inversion can potentially amplify the accumulated error, which could  
 10 ultimately lead to subpar reconstruction outcomes. To fortify the robustness of the DDIM inversion  
 11 and to mitigate this issue, we further optimize the null text embedding, as detailed in Mokady *et al.*  
 12 [3]. The incorporation of this optimization technique bolsters the effectiveness and reliability of the  
 13 inversion process, consequently allowing for a more precise reconstruction.

## 14 B Results of Other Concept Learning Methods

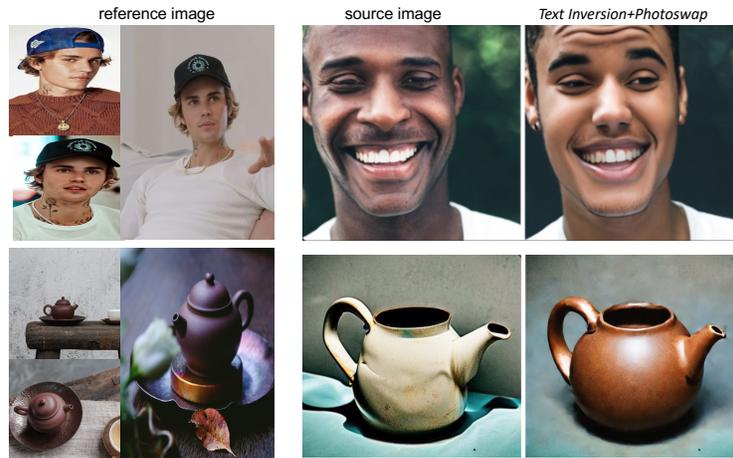


Figure 4: Results of Text Inversion [1] in *Photoswap*.

15 Our mainly use Dreambooth as the concept learning method in the experiments, primarily due to its  
 16 superior capabilities in learning subject identities [5]. However, our method is not strictly dependent  
 17 on any specific concept learning method. In fact, other concept learning methods could be effectively  
 18 employed to introduce the concept of the target subject.

19 To illustrate this, we present the results of *Photoswap* when applying Text Inversion [1]. We train the  
 20 model using 8 A100 GPUs with a batch size of 4, a learning rate of  $5e-4$ , and set the training steps to  
 21 1000. Results in Figure 4 indicate that Text Inversion also proves to be an effective concept learning  
 22 method, as it successfully captures key features of the target object. Nevertheless, we observe that  
 23 Text Inversion performance is notably underwhelming when applied to human faces. We postulate  
 24 that this is because Text Inversion focuses on learning a new embedding for the novel concept, rather  
 25 than finetuning the entire model. Consequently, the capacity to express the new concept becomes  
 26 inherently limited, resulting in its less than optimal performance in certain areas.

## 27 C Attention Swapping Step Analysis

28 In this section, we visualize the effect of the influence of swapping steps of different components.  
 29 As discussed in the main paper, self-attention output  $\phi$ , and self-attention map  $M$ , derived from the  
 30 self-attention layer, encompasses comprehensive content information from the source image, without  
 31 relying on direct computation with textual features. Previous works such as Hertz *et al.* [2] did not  
 32 explore the usage of  $\phi$  and  $M$  in the object-level image editing process.

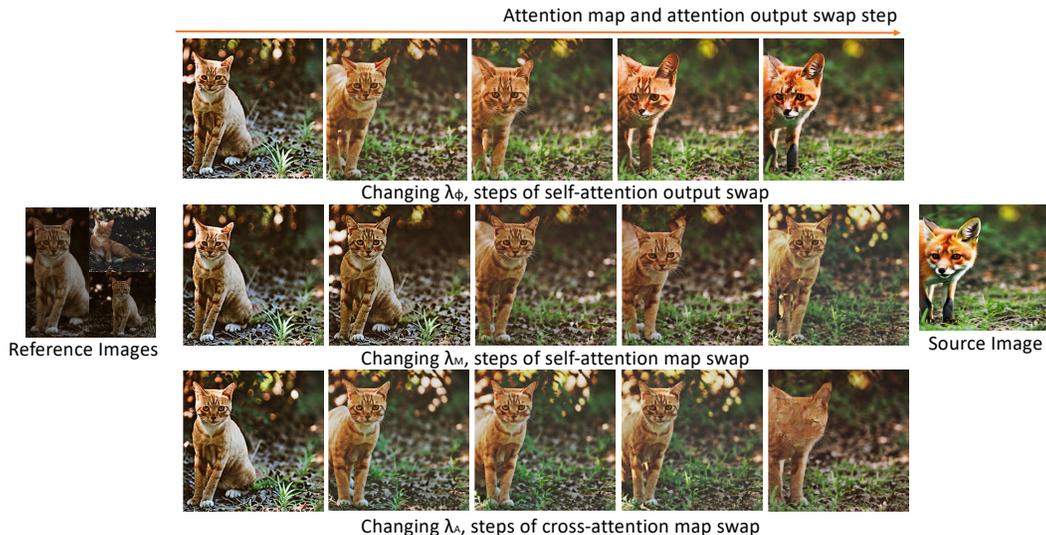


Figure 5: Results at different swapping steps. With consistent steps, swapping the self-attention output provides superior control over the layout, including the subject’s gestures and the background details. However, excessive swapping could affect the subject’s identity, as the new concept introduced through the text prompt might be overshadowed by the swapping of the attention output or attention map. This effect is more clear when swapping the self-attention output  $\lambda_\phi$ . Furthermore, we observed that replacing the attention map for an extensive number of steps can result in an image with significant noise, possibly due to a compatibility issue between the attention map and the  $v$  vector.

33 Figure 5 provides a visual representation of the effect of incrementally increasing the swapping step  
 34 for one  $\lambda$  hyperparameter while maintaining the other two at zero. Although all of them can be utilized  
 35 for subject swapping, they demonstrate varying levels of layout control. At the same swapping step,  
 36 the self-attention output  $\phi$  offers more robust layout control, facilitating better alignment of gestures  
 37 and preservation of background context. In contrast, the self-attention map  $M$  and cross-attention  
 38 map  $A$  demonstrate similar capabilities in controlling the layout.

39 However, extensive swapping can affect the subject’s identity, as the novel concept introduced via the  
 40 text prompt might be eclipsed by the swapping of the attention output or attention map. This effect  
 41 becomes particularly evident when swapping the self-attention output. This analysis further informs  
 42 the determination of the default  $\lambda_\phi$ ,  $\lambda_M$ , and  $\lambda_A$  values. While the cross-attention map  $A$  facilitates  
 43 more fine-grained generation control, given its incorporation of information from textual tokens, we  
 44 discovered that  $\phi$  offers stronger holistic generation control, bolstering the overall output’s quality  
 45 and integrity.

## 46 D Ethics Exploration

47 Like many AI technologies, text-to-image diffusion models can potentially exhibit biases reflective of  
 48 those inherent in the training data [4, 6]. Given that these models are trained on vast text and image  
 49 datasets, they might inadvertently learn and perpetuate biases, such as stereotypes and prejudices,  
 50 found within this data. For instance, should the training data contain skewed representations or  
 51 descriptions of specific demographic groups, the model may produce biased images in response to  
 52 related prompts.

53 However, *Photoswap* has been designed to mitigate bias within the generation process of a text-  
 54 to-image diffusion model. It achieves this by directly substituting the depicted subject with the

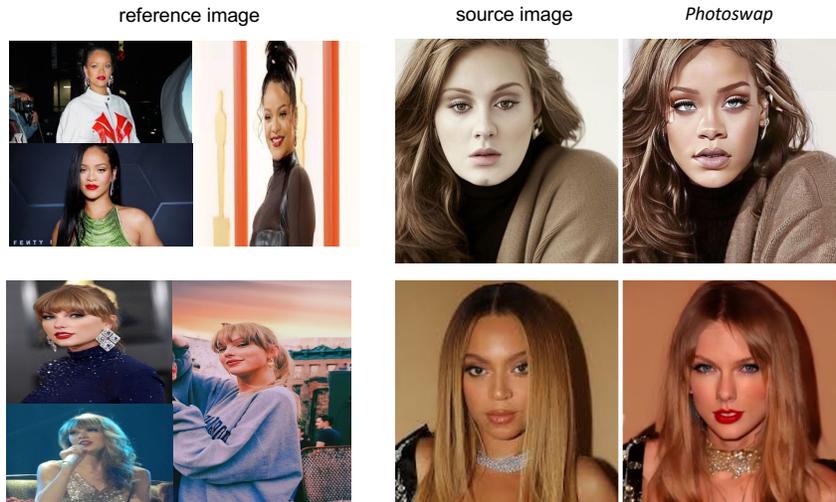


Figure 6: Human face swapping across different races. As you can see, the skin colors are also successfully transferred when swapping a white person with a black person, or vice versa.

55 intended target. In Figure 6, we present our evaluation of face swapping across various skin tones.  
 56 It is crucial to observe that when there is a significant disparity between the source and reference  
 57 images, the swapping results tend to homogenize the skin color. As a result, we advocate for the use  
 58 of *Photoswap* on subjects of similar racial backgrounds to achieve more satisfactory and authentic  
 59 outcomes. Despite these potential disparities, the model ensures the preservation of most of the target  
 60 subject’s specific facial features, reinforcing the credibility and accuracy of the final image.

## 61 E Self-Attention Map Visualization

62 In Figure 7, we show more visualization on self-attention map for real images. Here we show four  
 63 more examples of real images and synthetic images. The visualization results are consistent with  
 64 those in the main paper.

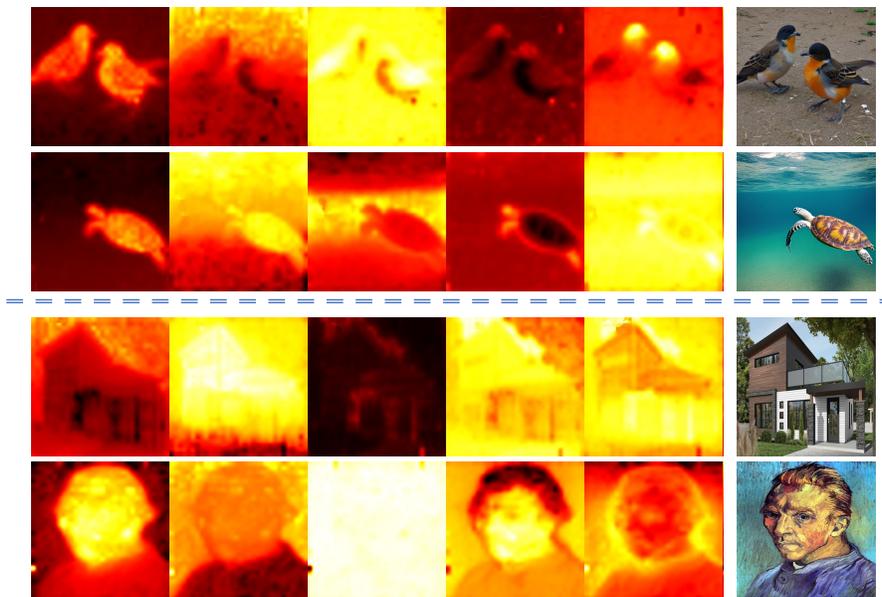


Figure 7: Self-attention visualization results. The top two rows are synthetic images and the bottom two rows are real images. There is a high correlation between self-attention maps and the images.

65 **F Failure Cases**

66 Here we highlight two common failure cases. First, the model struggles to accurately reproduce  
 67 hands. When the subject includes hands and fingers, the swapping results often fail to precisely  
 68 mirror the original hand gestures or the number of fingers. This issue could be an inherited challenge  
 69 from Stable Diffusion. Moreover, *Photoswap* can encounter difficulties when the image comprises  
 70 complex information. As illustrated in the lower row of Figure 8, *Photoswap* fails to reconstruct  
 71 the complicated formula on a whiteboard. Therefore, while *Photoswap* exhibits strong performance  
 72 across various scenarios, it’s crucial to acknowledge these limitations when considering its application  
 73 in real-world scenarios involving intricate hand gestures or complex abstract information. Future  
 74 work will focus on addressing these issues to enhance the overall performance and versatility of the  
 75 model.

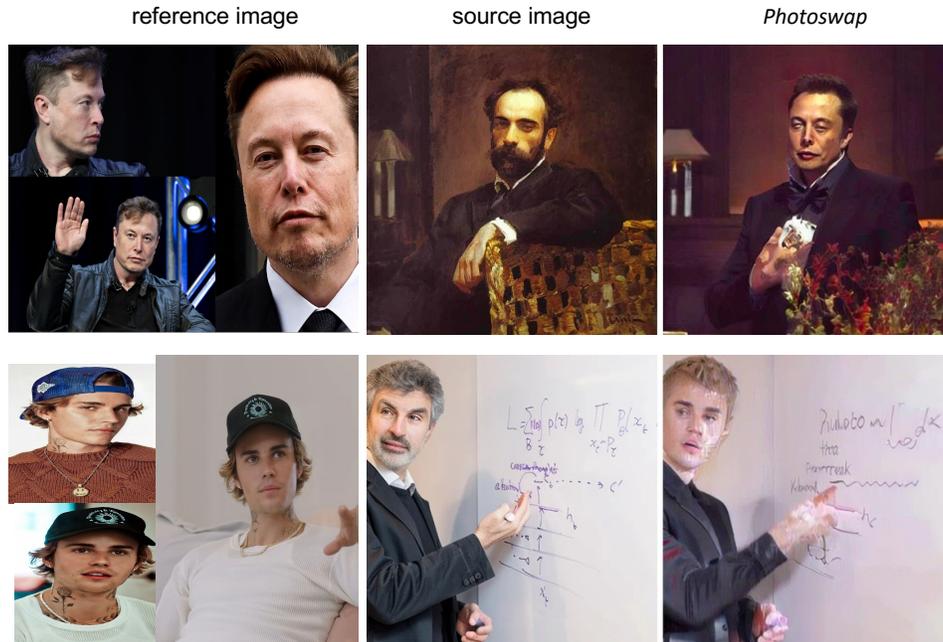


Figure 8: Failure cases.

76 **References**

77 [1] Gal, R., Alaluf, Y., Atzmon, Y., Patashnik, O., Bermano, A. H., Chechik, G., and Cohen-Or, D.  
 78 (2023). An Image is Worth One Word: Personalizing Text-to-Image Generation using Textual  
 79 Inversion. In *ICLR*.

80 [2] Hertz, A., Mokady, R., Tenenbaum, J., Aberman, K., Pritch, Y., and Cohen-Or, D. (2022).  
 81 Prompt-to-prompt image editing with cross attention control. *arXiv*.

82 [3] Mokady, R., Hertz, A., Aberman, K., Pritch, Y., and Cohen-Or, D. (2022). Null-text Inversion for  
 83 Editing Real Images using Guided Diffusion Models. In *arXiv*.

84 [4] Perera, M. V. and Patel, V. M. (2023). Analyzing bias in diffusion-based face generation models.  
 85 *arXiv preprint arXiv:2305.06402*.

86 [5] Ruiz, N., Li, Y., Jampani, V., Pritch, Y., Rubinstein, M., and Aberman, K. (2023). DreamBooth:  
 87 Fine Tuning Text-to-Image Diffusion Models for Subject-Driven Generation. In *CVPR*.

88 [6] Sasha Luccioni, A., Akiki, C., Mitchell, M., and Jernite, Y. (2023). Stable bias: Analyzing  
 89 societal representations in diffusion models. *arXiv e-prints*, pages arXiv-2303.

90 [7] Song, J., Meng, C., and Ermon, S. (2020). Denoising diffusion implicit models. In *International  
 91 Conference on Learning Representations*.