STYLEGUIDE: PREVENT CONTENT LEAKAGE USING NEGATIVE QUERY GUIDANCE

Anonymous authors

Paper under double-blind review

ABSTRACT

In the domain of text-to-image generation, diffusion models have emerged as powerful tools. Recently, studies on visual prompting, where images are used as prompts, have enabled more precise control over style and content. However, existing methods often suffer from content leakage, where undesired elements from the visual style prompt are transferred along with the intended style (content leakage). To address this issue, we 1) extend classifier-free guidance (CFG) to utilize swapping self-attention and propose 2) negative visual query guidance (NVQG) to reduce the transfer of unwanted contents. NVQG employs negative score by intentionally simulating content leakage scenarios which swaps queries instead of key and values of self-attention layers from visual style prompts. This simple yet effective method significantly reduces content leakage. Furthermore, we provide careful solutions for using a real image as a visual style prompts and for image-to-image (I2I) tasks. Through extensive evaluation across various styles and text prompts, our method demonstrates superiority over existing approaches, reflecting the style of the references and ensuring that resulting images match the text prompts.

029

024

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

Text-to-image diffusion models (T2I DMs) excel at synthesizing images that correspond to given
 text prompts (Rombach et al., 2022; Ramesh et al., 2021). However, relying solely on text prompts
 may not allow for precise control over the desired output. Even with highly detailed text prompts,
 controlling the exact style of the resulting images remains challenging (Figure 1 (a) and (b)). Text
 prompts fail to specify precise style elements such as color, shading, line details, surface texture, or
 polygon density.

To address this issue, there has been significant research into using reference images as visual style prompts. These approaches include fine-tuning the diffusion model with a set of images containing the same theme (Ruiz et al., 2023; Kumari et al., 2023), learning new text embeddings (Gal et al., 2022; Han et al., 2023a; Avrahami et al., 2023), and adapting cross-attention modules to incorporate image features (Ye et al., 2023; Wang et al., 2023). However, these methods require costly training and often let the content from the visual style prompts leak to the result, i.e., content leakage (Sohn et al., 2023).

043 In contrast, training-free methods (Hertz et al., 2023; Chung et al., 2024b; Alaluf et al., 2024) swap 044 features in the self-attention layer: the key and value from the visual style prompt replace the ones in the original process. Their motivation that the query carries the content, and the key-value carries the style allows promising performance for style reflection. However, this decomposition is not always 046 satisfactory, leading to trade-off relationship between style and content (e.g., content leakage or 047 poor style reflection). Moreover, when they tackle image-to-image (I2I, i.e., style transfer) where the 048 content and style are specified by visual content/style prompts, the style from the content prompt leaks to the result along with the content from the style prompt leaking to the result as shown in Figure 12. Hence, we need more than naive sampling with query and key-value swapping for I2I. 051 Appendix A further discusses related work. 052

In this study, we propose a method to more effectively extract the desired elements, whether style or content, from visual prompts and a text prompt. Our approach builds on classifier-free guidance



Figure 1: **Ambiguity of text prompts vs. visual style prompting.** (a) Ambiguity of text leads to different results within the same style description. (b) Even a detailed style description does not guarantee the generation of the same style images since it has many variants that can hardly be constrained using only text prompts. (c) Reference images can specify detailed visual elements.

(CFG) (Ho & Salimans, 2022) combined with swapping self-attention, enabling precise style transfer
 while maintaining the content specified by the text prompt. To address content leakage, we introduce
 negative visual query guidance, ensuring a clear separation between content and style. We also
 incorporate stochastic encoding for better style alignment and color calibration to match the final
 output to the reference image's color statistics.

We analyze where to apply swapping self-attention, identifying the optimal layers for balancing style
transfer and content fidelity. Additionally, our method can effectively remove content that is difficult
to eliminate with key and value swapping alone, working successfully even in cases with significant
structural gaps between the style image and content text, as shown in Figure 3 (e.g., a complex scene
of "a woman walking two dogs" and a single object "cat"). We extend the method to ControlNet for
I2I style transfer, further enhancing its flexibility. Qualitative and quantitative evaluations show our
method outperforms state-of-the-art approaches, providing precise control over content and style
without content leakage. Our approach is both robust and efficient, ideal for visual style prompting
tasks.¹

2 VISUAL STYLE PROMPTING

We propose StyleGuide which receives a text prompt and a visual style prompt to generate new images. The results contain the content and style specified by the text prompt and the visual style prompt, respectively, with variations due to initial noises. The overview of our method is illustrated in Figure 2. First, we explain the swapping self-attention in the aspect of style transfer literature. StyleGuide consists of classifier-free guidance (CFG) with swapping self-attention, negative visual query guidance (NVQG), optimal layer choice, stochastic encoding of real visual style prompts, and generalization to ControlNet for real content images. We explain the first three components in text-to-image (T2I) scenario with generated visual style prompt. Then we proceed to T2I with real visual style prompt and image-to-image (I2I) scenario.

2.1 SWAPPING SELF-ATTENTION IN STYLE TRANSFER LITERATURE

Modern diffusion models consist of a number of self-attention and cross-attention blocks (Vaswani et al., 2017). Both of them employ the attention mechanism, which first obtains an attention map using similarity between query features Q and key features K, then aggregates value features V using the attention map as weights: Attention $(Q, K, V) = \text{Softmax}(\frac{QK^T}{\sqrt{d}})V$. Opposed to the cross-attention

¹Our code will be released for reproducibility.



Figure 2: **Overveiw of StyleGuide**. Our proposed method includes 4 key components, highlighted in red boxes. First, stochastic encoding (Section 3.1) converts reference images into suitable latents for the visual style prompting task. Second, swapping self-attention (Section 2.2, 2.4) ensures the reference image's style is accurately reflected. Third, negative visual query guidance (Section 2.3) reduces content leakage from the reference image, allowing the desired text content (e.g., "Moose") to be better represented. Lastly, color calibration (Section 3.1) minimizes errors during the denoising process, helping to produce a cleaner final image.

129 130 131

132

133

134

135

136

layer, self-attention layer receives key and values coming from the main denoising process which has spatial dimensions with more freedom to represent spatially varying visual elements. As our goal is to reflect style elements from a reference image that are not easily represented by textual description, we opt to borrow key and values of self-attention layers in the reference process to the original process, namely swapping self-attention z (Figure 2). In addition, swapping self-attention has a strong connection with style transfer literature (Sheng et al., 2018; Park & Lee, 2019; Yao et al., 2019; Liu et al., 2021; Deng et al., 2022). where the attention mechanism reassembles visual features of a style image (key, value) on a content image (query).

137 138 139

140

2.2 CFG SAMPLING WITH SWAPPING SELF-ATTENTION FOR T2I

141 We propose CFG with swapping self-attention to reflect a visual style prompt in the T2I results. CFG 142 (Ho & Salimans, 2022) is essential to guide the generated images toward given text prompts. For 143 given a score $\epsilon_{\theta}(x_t, c)$ conditioned on *c* and unconditional score $\epsilon_{\theta}(x_t, \emptyset)$, the CFG score is defined 144 by $\tilde{\epsilon}_{\theta} = (1 + w)\epsilon_{\theta}(x_t, c) - w\epsilon_{\theta}(x_t, \emptyset)$ where *w* controls the guidance strength. ² CFG with w > 1145 improves image quality and text alignment but excessive *w* induces mode collapse (Chung et al., 146 2024a). Notably, CFG has not been explored in context of reflecting denoising process with modified 147 features.

Assuming there exists the desired but hidden content h^{content} and style h^{style} of a given condition c, we formulate our task to model $p_{\theta}(x_0|h_{\text{text}}^{\text{content}}, h_{\text{visual}}^{\text{style}})$ using $\epsilon_{\theta}(x_t, c_{\text{text}})$ and $\epsilon_{\theta}(x_t, c_{\text{visual}})$ which are the original score leading to the original T2I-generated image $x_0^{\text{ori}} \sim p_{\theta}(x_0|c_{\text{text}}) =$ $p_{\theta}(x_0|h_{\text{text}}^{\text{content}}, h_{\text{text}}^{\text{style}})$ and the reference score leading to the visual style prompt $x_0^{\text{visual}} \sim$ $p_{\theta}(x_0|c_{\text{visual}}) = p_{\theta}(x_0|h_{\text{visual}}^{\text{content}}, h_{\text{visual}}^{\text{style}})$, respectively. We design the CFG score toward the result with the desired but hidden content h^{content} and style h^{style} :

$$\tilde{\epsilon}_{\theta}(x_t, h_{\text{text}}^{\text{content}}, h_{\text{visual}}^{\text{style}}) = (1+w)\ddot{\epsilon}_{\theta}(x_t, Q_{\text{text}}, KV_{\text{visual}}) - w\epsilon_{\theta}(x_t, \emptyset), \tag{1}$$

where $\ddot{\epsilon}_{\theta}(x_t, Q_{\text{text}}, KV_{\text{visual}})$ denotes a KV-injected denoising score as below. We use $\ddot{\epsilon}_{\theta}$ to indicate the score is not from a single condition but is estimated by feature manipulation.

For given two denoising processes, one as original and another as a reference, borrowing the key-value in self-attention from the reference to the original process, i.e., key-value (KV) injection, tends to

161

²We omit the diffusion timestep t in the arguments and abuse x_t instead of z_t from latent diffusion model.



Figure 3: The effect of CFG and the proposed 172 negative visual query guidance on image genera-173 tion. The reference images provide the style for each 174 generated output. Without NVQG, content leakage 175 occurs, and the generated images fail to fully capture 176 the intended content. In contrast, using NVQG en-177 sures better alignment with both the reference style 178 and the 'Cat' prompt, reducing content distortion and 179 improving quality. 180



Figure 4: The effect of swapping selfattention across different blocks. Swapping self-attention on the bottleneck and downblocks causes content leakage, resulting in cat-like images despite a dog prompt, while swapping on downblocks disrupts resulting images. We only apply swapping selfattention in the upblocks to reflect style elements effectively.

produces results with the content from the original process and the style from the reference process
with limited control (Alaluf et al., 2024; Chung et al., 2024b; Xu et al., 2023).

We define the KV-injected score by $\ddot{\epsilon}_{\theta}(x_t, Q_{\text{text}}, KV_{\text{visual}})$ where Q_{text} and KV_{visual} denote the query from the original score³ $\epsilon_{\theta}(x_t, c_{\text{text}})$ and the key, value from the reference score $\epsilon_{\theta}(x_t, c_{\text{visual}})$. Then KV-injected-Attention becomes:

Attention
$$(Q_{\text{text}}, K_{\text{visual}}, V_{\text{visual}}) = \text{Softmax}(\frac{Q_{\text{text}}K_{\text{visual}}^{\mathsf{T}}}{\sqrt{d}})V_{\text{visual}}.$$
 (2)

190 We omit the layer index for brevity.

Naive denoising process with $\ddot{\epsilon}_{\theta}(x_t, Q_{\text{text}}, KV_{\text{visual}})$ provides limited control in generating images with content and style specified by a text c_{text} and a visual style prompt, respectively. Moving forward, our CFG with swapping self-attention in Eq. (1) enjoys higher image quality and more accurate text alignment than the naive denoising process as in the original CFG for T2I generation. The results are deferred to Section 4.2.

196 197

207

212 213

214 215

181

187 188

189

2.3 NEGATIVE VISUAL QUERY GUIDANCE

We propose negative visual query guidance (NVQG) to further reduce the content $h_{\text{visual}}^{\text{content}}$ from the visual style prompt appearing in the results. Briefly, NVQG negates the CFG of a score $\tilde{\epsilon}(x_t, Q_{\text{visual}}, KV_{\text{text}})$.

In Liu et al. (2022), a complex text condition c is factorized as a set of conditions $\{c_0, c_1, ...\}$ and Bayes' rule induces $p_{\theta}(x_t | c_0, c_1, ...) \propto \prod \frac{p_{\theta}(x_t | c_i)}{p_{\theta}(x_t)}$. Then, the score of the complex text condition cbecomes $\epsilon_{\theta}(x_t, c) = \epsilon_{\theta}(x_t, \emptyset) + \prod(\epsilon_{\theta}(x_t, c_i) - \epsilon_{\theta}(x_t, \emptyset))$. It allows reducing a specific concept \tilde{c} with composition by negating its guidance with scale w_{neg} :

$$\epsilon_{\theta}(x_t, c, \operatorname{not} \tilde{c}) = \epsilon_{\theta}(x_t, c) - w_{\operatorname{neg}}(\epsilon_{\theta}(x_t, \tilde{c}) - \epsilon_{\theta}(x_t, \emptyset))$$
(3)

208 209 Although we design $\ddot{\epsilon}_{\theta}(x_t, Q_{\text{text}}, KV_{\text{visual}})$ to predict the score toward $p_{\theta}(x_0 | h_{\text{text}}^{\text{content}}, h_{\text{visual}}^{\text{style}})$, 210 $\ddot{\epsilon}(x_t, Q_{\text{text}}, KV_{\text{visual}})$ still contain $h_{\text{visual}}^{\text{content}}$. Assuming a hidden factorization $KV_{\text{visual}} = \{KV_{\text{visual}}^{\text{content}}, KV_{\text{visual}}^{\text{style}}\}$, Bayes' rule induces

$$p_{\theta}(x_t | Q_{\text{text}}, KV_{\text{visual}}^{\text{style}}, KV_{\text{visual}}^{\text{content}}) \propto p_{\theta}(x_t) \frac{p_{\theta}(x_t | Q_{\text{text}}, KV_{\text{visual}}^{\text{style}})}{p_{\theta}(x_t)} \frac{p_{\theta}(x_t | Q_{\emptyset}, KV_{\text{visual}}^{\text{content}})}{p_{\theta}(x_t)}.$$
 (4)

³The original score and its query within are recursively altered by KV injection along the denoising process.

216 217 218 Note that, $\epsilon_{\theta}(x_t) = \epsilon_{\theta}(x_t, Q_{\emptyset}, KV_{\emptyset})$ where the source of the Q,K,V is \emptyset . Then, we get the desired 219 conditional score of $\hat{c} = \{Q_{\text{text}}, KV_{\text{visual}}^{\text{style}}\}$:

219
$$\epsilon_{\theta}(x_t, \hat{c}) \leftarrow w_{\text{visual}}(\ddot{\epsilon}_{\theta}(x_t, Q_{\text{text}}, KV_{\text{visual}}) - \epsilon_{\theta}(x_t)) - w_{\text{content}}(\ddot{\epsilon}_{\theta}(x_t, Q_{\emptyset}, KV_{\text{visual}}^{\text{content}}) - \epsilon_{\theta}(x_t)) + \epsilon_{\theta}(x_t)$$
220 (5)

where w_{visual} and w_{content} sets the strength of each classifier. By borrowing the ability of query injection which successfully conveys content Tumanyan et al. (2023); Alaluf et al. (2024); Chung et al. (2024b); Xu et al. (2023), we approximate $\ddot{\epsilon}_{\theta}(x_t, Q_{\emptyset}, KV_{\text{visual}}^{\text{content}}) \approx \ddot{\epsilon}_{\theta}(x_t, Q_{\text{visual}}, KV_{\emptyset})$ Lastly, we insert Eq. (5) into Eq. (1). Empirically, we find that replacing $\epsilon_{\theta}(x_t)$ to $\ddot{\epsilon}_{\theta}(x_t, Q_{\text{visual}})$ brings similar results. Finally, we can simply reformulate diffusion sampling with $w' = w_{\text{visual}}(w + 1)$:

$$\epsilon_{\theta}(x_t, \hat{c}) \leftarrow (w'+1)(\ddot{\epsilon}_{\theta}(x_t, Q_{\text{text}}, KV_{\text{visual}}) - w'\ddot{\epsilon}_{\theta}(x_t, Q_{\text{visual}})$$
(6)

Since we can highly relate the $w' \tilde{\epsilon}_{\theta}(x_t, Q_{\text{visual}})$ to the concept negation in equation 3 which guides the negation of concept with a scale w_{neg} , we named the query term as negative visual query guidance.

2.4 CHOOSING BLOCKS FOR SWAPPING SELF-ATTENTION

226

230

231

249

250 251

252

256

257 258

232 Here we explore the depth of the self-attention blocks to be swapped in the sense of granularity 233 of visual elements. Modern architecture of diffusion models roughly consists of three sections in a sequence: downblocks, bottleneck blocks, and upblocks. Given that the bottleneck of diffusion models 234 contains content elements of the image (Kwon et al., 2023; Jeong et al., 2024; Park et al., 2023), we 235 opt not to apply swapping self-attention to bottleneck blocks to prevent transferring contents in a 236 reference image. Figure 4 shows that not swapping the bottleneck blocks prevents content leaking 237 from the reference image. However, the synthesized images show disrupted results with seriously 238 scattered objects. Furthermore, while swapping self-attention implements the reassembling operation, 239 simply applying to all self-attention layers exposes a content leakage problem, where the content of 240 the reference image influences the resulting image, as shown in the first row of Figure 4. I.e., the 241 results contain cats even though the prompts specify "a dog". We conjecture that this phenomenon 242 happens because feature maps of downblocks have unclear content layout (Cao et al., 2023; Meng 243 et al., 2024), so substituting features based on this inaccurate layout causes the disrupted results. To 244 avoid injecting unnecessary features, we choose to swap the key and value of self-attention only in 245 upblocks.

We note that Hertz et al. (2023) applies self-attention operation at the all blocks and suffers content leakage. The last row of Figure 4 shows the success of our strategy.

3 REAL IMAGES AS REFERENCES

3.1 REAL IMAGES AS VISUAL STYLE PROMPTS

So far, we have assumed a *generated* visual style prompt $x_0^{\text{visual}} \sim p_\theta(x_0|c_{\text{visual}})$. Here, we allow *real* visual style prompts by 1) stochastic encoding and 2) color calibration.

We propose stochastic encoding to obtain $x_t^{\text{visual}} \sim q(x_t | x_0^{\text{visual}})$ by adding a random noise on x_0^{visual} following the forward process of DMs (Ho et al., 2020):

$$\epsilon_t \sim \mathcal{N}(0, I), \ x_t^{\text{visual}} = \sqrt{\alpha_t} \cdot x_0^{\text{visual}} + \sqrt{1 - \alpha_t} \cdot \epsilon_t$$
(7)

At each timestep, we samples ϵ_t to encode x_t^{visual} . It ensures that x_t^{visual} lies on the learned trajectory and does not suffer from accumulative numerical error due to iterative process of DDIM inversion used by previous methods (Hertz et al., 2023; Chung et al., 2024b). Furthermore, it does not need to store the intermediate latents as opposed to DDPM inversion used by Alaluf et al. (2024).

Although stochastic encoding performs better than DDIM inversion, a subtle color discrepancy occurs between the resulting images and the visual style prompts. We introduce color calibration at x_t^{ori} in the original process to match the statistics of predicted x_0^{ori} to predicted x_0^{visual} . In Gatys (2015), distance between channel-wise statistics is employed as a style loss for style transfer. In Song et al. (2020), predicted $x_0 (= \frac{x_t^{\text{visual}} - \sqrt{1 - \alpha_t} \cdot \epsilon_{\theta}(\hat{x}_t)}{\sqrt{\alpha_t}})$ allows to estimate x_0 with high probability at intermediate timesteps using deterministic sampling. Inheriting the advantage, we execute adain operation to match mean&std of predicted x_0^{ori} with those of x_0^{visual} . It allows precise color calibration



Figure 5: Analysis on the optimal range of upblocks for swapping self-attention. We find the optimal range of upblocks for a balanced trade-off between different aspects. The images on the right illustrate the visual results for different upblock layer indices, with the red cross indicating poor diversity and misalignment to the text prompt, the red triangle indicating a lack of style similarity, and the yellow star indicating the optimal results. Please refer to Section 2.4 for details.

rather than directly matching x_t^{visual} in Alaluf et al. (2024); Chung et al. (2024b). Furthermore, 292 ours differentiates from Chung et al. (2024b) in that using predicted x_0 at intermediate timestep 293 $t \in (0,T)$ other than x_T by reducing cumulative sampling error after the operation. Furthermore, Chung et al. (2024b) executes AdaIN at timestep T, inducing lengthy cumulative error.

We provide supportive experiments that show the effectiveness of the proposed method in Section 4.4 and a detailed Algorithm in Appendix B.2.

3.2 **REAL IMAGE AS A CONTENT FOR STYLE TRANSFER**

300 Our method can be extended not only to T2I tasks but also to I2I tasks, where users want to control 301 the content using an image. In this I2I scenario, we adopt an approach where structural information 302 from the content image is injected using ControlNet (Zhang et al., 2023). 303

Compared to our work, most existing self-attention variants (Chung et al., 2024b; Alaluf et al., 304 2024) for I2I style transfer employs query injection in self-attention to specify contents. However, 305 the query from the content image contains not only the content elements (e.g., structure, layout, 306 components) but also high nuance style elements (e.g., texture, pattern, and mesh) of the given image. 307 As a result, style leakage can occur, transferring unwanted style elements from the content images. In 308 the subsection 4.5, we demonstrate that our approach is more robust to style leakage issues compared 309 to existing self-attention methods when dealing with real content images.

310 311

284

285

286

287

288

289 290

291

295

296

297 298

299

- 4 EXPERIMENTS
- 312 313

314 In this section, we describe the effects of our proposed methods: CFG with swapping self-attention, 315 Negative visual query guidance (NVQG), stochastic encoding, and color calibration. For swapping self-attention, we provide a detailed analysis through experiments to determine the optimal layers for 316 swapping. The impact of NVQGis demonstrated through qualitative results. Additionally, we show 317 why stochastic encoding outperforms DDIM inversion when inverting real images, and we highlight 318 the benefits of color calibration through experimental results. 319

320 We also conducted both quantitative and qualitative comparisons of our method against various 321 competitors, including StyleAligned (Hertz et al., 2023), IP-Adapter (Ye et al., 2023), Dreambooth-LoRA Ruiz et al. (2023); Ryu (2023), StyleDrop (Sohn et al., 2023), DEADiff (Qi et al., 2024), CSGO 322 (Xing et al., 2024) and InstantStyle (-plus) (Wang et al., 2024a;b). The details of these comparisons, 323 along with the experimental setup and metrics, are described in the Appendix B.1.



Figure 6: Attention map visualization over late and early upblock layers. The late upblock better focuses on the style-corresponding region than the early upblock, leading to more freedom to reassemble small parts. The early upblock attends larger region leading to content leakage.

4.1 ANALYSIS FOR SWAPPING SELF-ATTENTION

Optimal layers in upblocks Since recent large T2I DMs consist of many layers, we further analyze
 the behavior by changing the start of the swapping while the end of the swapping is fixed at the end.
 We use four key metrics as shown in Figure 5, there is a layer where all four metrics abruptly change
 (red line). Notably, this point remains consistent regardless of the reference image. We choose this
 layer as the optimal start of the swapping for a balanced trade-off among all aspects. We provide
 qualitative results with detailed split of layers in Figure A2 and A3.

347

324

325

326

327

328

330

331332333334

335

336

337 338 339

340

348 **Visualizing Attention maps** Figure 6 compares average attention maps from the late upblock 349 and the early upblock applying swapping self-attention. Using late upblock has more freedom to reassemble the reference style elements leading to more doggy results than early upblock which 350 produces some cat-like attributes. The right two columns visualize the attention weight of query points 351 marked as red stars and yellow dots. Swapping self-attention on late upblock reassembles features 352 from a style correspondence, e.g., texture and color. On the other hand, swapping self-attention on 353 early upblock reassembles features from a wider area with different styles. This comparison clarifies 354 the reasons for using only late upblock. Please refer to Figure A5 for a detailed analysis. 355

356 357

4.2 EFFECTIVENESS OF CFG AND NVQG

358 This section analyzes the effects of Classifier-Free Guidance (CFG) and Negative Visual Query 359 Guidance (NVQG) on image generation, with a focus on text alignment and content leakage. Figure 3 360 shows the results of the three configurations. In the 1st row, without CFG and NVQG, the generated 361 images suffer from severe artifacts. The absence of CFG causes poor image quality resulting in 362 significant misalignment with the prompt. In the 2nd row, CFG with swapping self-attention improves 363 the text misalignment by boosting image quality. Here, the "cat" in target text prompt becomes clearer 364 in the generated images. However, content leakage from the reference image still remains where unwanted elements (layouts, structure, and composition) from the reference affecting the results. In 366 the 3rd row, NVQG releases the content leakage and produces the best results closely matching the 367 text prompt while reflecting style reference.

Overall, Figure 3 demonstrates the critical role of NVQG in reducing content leakage, enjoying the
 quality boosting of CFG. Together, they ensure that the generated images align to both the target text
 prompt and the visual style prompts, resulting in high-quality, coherent outputs.

We qualitatively showcase diversity of results within a text prompt in Figure A7 and text alignment with complex text prompts in Figure A22.

374 375

376

4.3 COMPARISON AGAINST COMPETITORS

We compare ours with StyleAligned Hertz et al. (2023), IP-Adapter Ye et al. (2023), Dreambooth-LoRA Ruiz et al. (2023); Ryu (2023), and StyleDrop Sohn et al. (2023).



Figure 7: Qualitative comparison across various styles and text prompts. StyleGuide faithfully reflects style elements in reference images without causing content leakage from the reference images.





Figure 8: Quantitative comparison. We compare the results for text alignment (CLIP score) and tween other methods (blue points) and our method i.e., content leakage from the reference. (orange point).

Figure 9: Qualitative comparison with the same style. Competitors face challenges in generating style similarity (DINO embedding similarity) be- images with diverse layouts and compositions,

Style & content control We provide a qualitative comparison in Figure 7, focusing on controlling style and content. Our method faithfully synthesizes content from the text prompt with the style of the reference image. In contrast, other methods add elements like color or texture not in the reference (e.g., feathers, bricks, iron, skin) and often suffer from content leakage (e.g., layout, screaming person, castle), which compromises text prompt faithfulness. Quantitative results in Figure 8 support these findings: IP-Adapter shows higher style similarity but neglects text prompts significantly. We provide additional comparisons with DEADiff (Qi et al., 2024), CSGO (Xing et al., 2024) and InstantStyle (-plus) (Wang et al., 2024a;b) in Figure A11.

Diversity within a text specification Starting from different initial noises, the diffusion models trained on a large dataset produce diverse results within the specification of a text prompt. Figure 9 shows that our results have various poses and viewpoints while others barely change, i.e., other methods limit the diversity of the pre-trained model. Figure A7 provides more examples.



478 479

480

4.4 COMPARISON OF DDIM INVERSION WITH OUR STRATEGY

Figure 11 shows that StyleGuide can take real images as style reference with our strategies. As shown
in Figure 11, stochastic encoding outperforms DDIM inversion and color calibration improves color
consistency with the reference image. We provide more results in Figure A15. Moreover, we show that
our strategy can boost the performance of the other self-attention variants Hertz et al. (2023); Cao et al.
(2023) in Figure A16, A17 and color calibration can be used for generation settings in Figure A18. We
provide an ablation study in Figure A19 for each configuration (swapping self-attention, NVQG, and



Figure 12: **Qualitative comparison in I2I style transfer task.** We compare our method for I2I style transfer task where the content image is given to control the content directly. Compared to the previous methods, our method transfer the reference style more accurately without style(e.g. color) leakage from the content image.

color calibration) in both real reference and generated reference settings. In both settings, swapping self-attention and NVQG improve style similarity and text alignment, while color calibration helps improve style similarity. Additionally, stochastic encoding demonstrates better performance than DDIM inversion.

4.5 COMPARISON IN STYLE TRANSFER

In Figure 12, we present a qualitative comparison between our method using ControlNet and existing state-of-the-art methods, CrossAttn (Alaluf et al., 2024) and StyleID (Chung et al., 2024b), for the I2I style transfer task where a content image is provided. Both CrossAttn and StyleID inject the query obtained by inversion from the content images. As discussed in subsection 3.2, the obtained query includes style elements from the content images, which results in the reference style not being properly reflected in the output. CrossAttn often fails to transfer the detailed representation of the reference style, leading to a rough, blocky appearance. For instance, when comparing the center examples in Figure 12, the style details are absent or inaccurately represented compared to our result. Similarly, StyleID struggles with style leakage, particularly with color information. In the center example, there is clear color leakage from the entire content image, and the same issue is visible in the rightmost example. In contrast, our method effectively reflects the details of the reference style image, with no noticeable transfer of the color values from the content images.

5 CONCLUSION AND LIMITATION

In this paper, we introduce StyleGuide using swapping self-attention, which effectively applies the style of reference images without content leakage in a training-free manner. CFG with swapping self-attention captures the reference image's style accurately and allows for direct content generation from the text, making it superior to other approaches. This integration with CFG enhances performance by balancing content generation and style transfer. To address content leakage in visual style prompting tasks, we propose Negative visual guidance, a simple method that ensures the reference image's content does not interfere with the text-specified content. Additionally, Stochastic encoding maps real images to suitable latents, improving the overall accuracy of the generated style, while **Color** calibration aligns the final output to the reference's color statistics.

Our method demonstrates qualitative and quantitative improvements over existing approaches, offering a robust solution for visual style prompting without complex training. We also provide a detailed analysis on where to apply swapping self-attention, identifying the optimal layers to balance style transfer and content fidelity. StyleGuide outperforms existing methods both qualitatively and quantitatively, and can be easily combined with algorithms such as ControlNet (Zhang et al., 2023) and Dreambooth-LoRA (Ryu, 2023), as shown in Figure A20.

- However, StyleGuide is limited by the pretrained diffusion models' capabilities, unable to generate elements beyond the original model's scope (e.g., "stone golem" in Figure A24a). In adition the visually specified style overrides the textually specified style if they disagree as shown in Figure A24b.
- 539 For future work, expanding our method to other domains, such as video content, could broaden its applicability and open new research directions.

540 REFERENCES 541

547

551

553

559

563

565

571

576

577

582

583

584

542	Yuval Alaluf, Daniel Garibi, Or Patashnik, Hadar Averbuch-Elor, and Daniel Cohen-Or. Cross-image
543	attention for zero-shot appearance transfer. In ACM SIGGRAPH 2024 Conference Papers, pp.
544	1–12, 2024.

- Omri Avrahami, Kfir Aberman, Ohad Fried, Daniel Cohen-Or, and Dani Lischinski. Break-a-scene: 546 Extracting multiple concepts from a single image. In SIGGRAPH Asia, 2023.
- Mingdeng Cao, Xintao Wang, Zhongang Qi, Ying Shan, Xiaohu Qie, and Yinqiang Zheng. Masactrl: 548 Tuning-free mutual self-attention control for consistent image synthesis and editing. In ICCV), 549 2023. 550
- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and 552 Armand Joulin. Emerging properties in self-supervised vision transformers. In ICCV, 2021.
- Hyungjin Chung, Jeongsol Kim, Geon Yeong Park, Hyelin Nam, and Jong Chul Ye. Cfg++: Manifold-554 constrained classifier free guidance for diffusion models. arXiv preprint arXiv:2406.08070, 2024a.
- 556 Jiwoo Chung, Sangeek Hyun, and Jae-Pil Heo. Style injection in diffusion: A training-free approach for adapting large-scale diffusion models for style transfer. In Proceedings of the IEEE/CVF 558 Conference on Computer Vision and Pattern Recognition, pp. 8795-8805, 2024b.
- Guillaume Couairon, Jakob Verbeek, Holger Schwenk, and Matthieu Cord. Diffedit: Diffusion-based 560 semantic image editing with mask guidance, 2022. 561
 - Yingying Deng, Fan Tang, Weiming Dong, Chongyang Ma, Xingjia Pan, Lei Wang, and Changsheng Xu. Stytr2: Image style transfer with transformers. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 11326–11336, 2022.
- Martin Nicolas Everaert, Marco Bocchio, Sami Arpa, Sabine Süsstrunk, and Radhakrishna Achanta. 566 Diffusion in style. In ICCV, 2023. 567
- 568 Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit Haim Bermano, Gal Chechik, and 569 Daniel Cohen-or. An image is worth one word: Personalizing text-to-image generation using 570 textual inversion. In ICLR, 2022.
- Daniel Garibi, Or Patashnik, Andrey Voynov, Hadar Averbuch-Elor, and Daniel Cohen-Or. Renoise: 572 Real image inversion through iterative noising, 2024. 573
- 574 Leon A Gatys. A neural algorithm of artistic style. arXiv preprint arXiv:1508.06576, 2015. 575
 - Inhwa Han, Serin Yang, Taesung Kwon, and Jong Chul Ye. Highly personalized text embedding for image manipulation by stable diffusion. arXiv, 2023a.
- 578 Ligong Han, Song Wen, Qi Chen, Zhixing Zhang, Kunpeng Song, Mengwei Ren, Ruijiang Gao, 579 Anastasis Stathopoulos, Xiaoxiao He, Yuxiao Chen, Di Liu, Qilong Zhangli, Jindong Jiang, 580 Zhaoyang Xia, Akash Srivastava, and Dimitris Metaxas. Improving tuning-free real image editing 581 with proximal guidance, 2023b.
 - Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Promptto-prompt image editing with cross attention control, 2022.
- 585 Amir Hertz, Andrey Voynov, Shlomi Fruchter, and Daniel Cohen-Or. Style aligned image generation 586 via shared attention. arXiv, 2023.
- Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. arXiv preprint arXiv:2207.12598, 588 2022.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In NeurIPS, 591 2020. 592
- Inbar Huberman-Spiegelglas, Vladimir Kulikov, and Tomer Michaeli. An edit friendly ddpm noise space: Inversion and manipulations, 2024.

- 594 Jaeseok Jeong, Mingi Kwon, and Youngjung Uh. Training-free content injection using h-space 595 in diffusion models. In Proceedings of the IEEE/CVF Winter Conference on Applications of 596 Computer Vision, pp. 5151–5161, 2024. 597 Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept 598 customization of text-to-image diffusion. In CVPR, 2023. 600 Mingi Kwon, Jaeseok Jeong, and Youngjung Uh. Diffusion models already have a semantic latent 601 space. In ICLR, 2023. 602 Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chunyuan Li, 603 and Yong Jae Lee. Gligen: Open-set grounded text-to-image generation. In CVPR, 2023. 604 605 Nan Liu, Shuang Li, Yilun Du, Antonio Torralba, and Joshua B Tenenbaum. Compositional visual 606 generation with composable diffusion models. In European Conference on Computer Vision, pp. 607 423-439. Springer, 2022. 608 Songhua Liu, Tianwei Lin, Dongliang He, Fu Li, Meiling Wang, Xin Li, Zhengxing Sun, Qian Li, 609 and Errui Ding. Adaattn: Revisit attention mechanism in arbitrary neural style transfer. In ICCV, 610 2021. 611 612 Frank J. Massey. The kolmogorov-smirnov test for goodness of fit. Journal of the American 613 Statistical Association, 46(253):68-78, 1951. ISSN 01621459. URL http://www.jstor. 614 org/stable/2280095. 615 Benyuan Meng, Qianqian Xu, Zitai Wang, Xiaochun Cao, and Qingming Huang. Not all diffusion 616 model activations have been evaluated as discriminative features. arXiv preprint arXiv:2410.03558, 617 2024. 618 619 Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. 620 SDEdit: Guided image synthesis and editing with stochastic differential equations. In *ICLR*, 2022. 621 Daiki Miyake, Akihiro Iohara, Yu Saito, and Toshiyuki Tanaka. Negative-prompt inversion: Fast 622 image inversion for editing with text-guided diffusion models, 2023. 623 624 Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for 625 editing real images using guided diffusion models. In CVPR, 2023. 626 Thao Nguyen, Yuheng Li, Utkarsh Ojha, and Yong Jae Lee. Visual instruction inversion: Image 627 editing via visual prompting. NeurIPS, 2023. 628 629 Dae Young Park and Kwang Hee Lee. Arbitrary style transfer with style-attentional networks. In 630 CVPR, 2019. 631 632 Yong-Hyun Park, Mingi Kwon, Jaewoong Choi, Junghyo Jo, and Youngjung Uh. Understanding the latent space of diffusion models through the lens of riemannian geometry. Advances in Neural 633 Information Processing Systems, 36:24129–24142, 2023. 634 635 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe 636 Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image 637 synthesis. arXiv, 2023. 638 Tianhao Qi, Shancheng Fang, Yanze Wu, Hongtao Xie, Jiawei Liu, Lang Chen, Qian He, and Yong-639 dong Zhang. Deadiff: An efficient stylization diffusion model with disentangled representations. 640 In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 641 8693-8702, 2024. 642 643 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, 644 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual 645 models from natural language supervision. In *ICML*, 2021. 646 Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, 647
 - and Ilya Sutskever. Zero-shot text-to-image generation. In *ICML*, pp. 8821–8831. PMLR, 2021.

648 649 650	Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In <i>CVPR</i> , 2022.
651 652 653	Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In <i>CVPR</i> , 2023.
654 655 656	Simo Ryu. Low-rank adaptation for fast text-to-image diffusion fine-tuning. https://github.com/cloneofsimo/lora, 2023.
657 658	Lu Sheng, Ziyi Lin, Jing Shao, and Xiaogang Wang. Avatar-net: Multi-scale zero-shot style transfer by feature decoration. In <i>CVPR</i> , 2018.
659 660 661 662	Kihyuk Sohn, Lu Jiang, Jarred Barber, Kimin Lee, Nataniel Ruiz, Dilip Krishnan, Huiwen Chang, Yuanzhen Li, Irfan Essa, Michael Rubinstein, Yuan Hao, Glenn Entis, Irina Blok, and Daniel Castro Chin. Styledrop: Text-to-image synthesis of any style. In <i>NeurIPS</i> , 2023.
663 664	Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In <i>ICLR</i> , 2020.
665 666 667	Adéla Šubrtová, Michal Lukáč, Jan Čech, David Futschik, Eli Shechtman, and Daniel Sỳkora. Diffusion image analogies. In <i>SIGGRAPH</i> , 2023.
668 669 670	Yasheng Sun, Yifan Yang, Houwen Peng, Yifei Shen, Yuqing Yang, Han Hu, Lili Qiu, and Hideki Koike. Imagebrush: Learning visual in-context instructions for exemplar-based image manipulation. <i>NeurIPS</i> , 2023.
671 672 673	Narek Tumanyan, Michal Geyer, Shai Bagon, and Tali Dekel. Plug-and-play diffusion features for text-driven image-to-image translation. In <i>CVPR</i> , 2023.
674 675	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In <i>NeurIPS</i> , 2017.
676 677 678	Andrey Voynov, Qinghao Chu, Daniel Cohen-Or, and Kfir Aberman. $p+$: Extended textual conditioning in text-to-image generation. <i>arXiv</i> , 2023.
679 680	Bram Wallace, Akash Gokul, and Nikhil Naik. Edict: Exact diffusion inversion via coupled transfor- mations, 2022.
682 683 684	Haofan Wang, Matteo Spinelli, Qixun Wang, Xu Bai, Zekui Qin, and Anthony Chen. Instantstyle: Free lunch towards style-preserving in text-to-image generation. <i>arXiv preprint arXiv:2404.02733</i> , 2024a.
685 686 687	Haofan Wang, Peng Xing, Renyuan Huang, Hao Ai, Qixun Wang, and Xu Bai. Instantstyle-plus: Style transfer with content-preserving in text-to-image generation. <i>arXiv preprint arXiv:2407.00788</i> , 2024b.
689 690	Zhouxia Wang, Xintao Wang, Liangbin Xie, Zhongang Qi, Ying Shan, Wenping Wang, and Ping Luo. Styleadapter: A single-pass lora-free model for stylized image generation. <i>arXiv</i> , 2023.
691 692 693	Chen Henry Wu and Fernando De la Torre. A latent space of stochastic diffusion models for zero-shot image editing and guidance. In <i>ICCV</i> , 2023.
694 695 696	Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Weixian Lei, Yuchao Gu, Yufei Shi, Wynne Hsu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation. In <i>ICCV</i> , 2023.
697 698 699 700	Peng Xing, Haofan Wang, Yanpeng Sun, Qixun Wang, Xu Bai, Hao Ai, Renyuan Huang, and Zechao Li. Csgo: Content-style composition in text-to-image generation. arXiv preprint arXiv:2408.16766, 2024.
701	Sihan Xu, Yidong Huang, Jiayi Pan, Ziqiao Ma, and Joyce Chai. Inversion-free image editing with natural language. <i>arXiv</i> , 2023.

702 703	Shuai Yang, Yifan Zhou, Ziwei Liu, and Chen Change Loy. Rerender a video: Zero-shot text-guided video-to-video translation. In <i>SIGGRAPH Asia</i> , 2023.
704 705 706	Yuan Yao, Jianqiang Ren, Xuansong Xie, Weidong Liu, Yong-Jin Liu, and Jun Wang. Attention-aware multi-stroke style transfer, 2019.
707 708	Hu Ye, Jun Zhang, Sibo Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models. <i>arXiv</i> , 2023.
709 710 711	Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In <i>ICCV</i> , 2023.
712 713	Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In <i>CVPR</i> , 2018.
714	
715	
717	
718	
719	
720	
721	
722	
723	
724	
725	
726	
727	
728	
729	
730	
731	
732	
733	
734	
735	
730	
737	
730	
740	
741	
742	
743	
744	
745	
746	
747	
748	
749	
750	
751	
752	
753	
754	
755	