MAGIC INSERT: STYLE-AWARE DRAG-AND-DROP

Anonymous authors

000

001 002 003

004

021

023 024 025

026

027 028

029

031

032

034

Paper under double-blind review



Figure 1: Using *Magic Insert* we are able to, for the first time, drag-and-drop a subject from an image with an arbitrary style onto another target image with a vastly different style and achieve a style-aware and realistic insertion of the subject into the target image.

ABSTRACT

We present **Magic Insert**, a method for dragging-and-dropping subjects from a user-provided image into a target image of a different style in a physically plausible manner while matching the style of the target image. This work formalizes the problem of style-aware drag-and-drop and presents a method for tackling it by addressing two sub-problems: *style-aware personalization* and *realistic object insertion in stylized images*. For style-aware personalization, our method first fine-tunes a pretrained text-to-image diffusion model using LoRA and learned text tokens on the subject image, and then infuses it with a CLIP representation of the target style. For object insertion, we use *Bootstrapped Domain Adaption* to adapt a domain-specific photorealistic object insertion model to the domain of diverse artistic styles. Overall, the method significantly outperforms traditional approaches such as inpainting. Finally, we present a dataset, SubjectPlop, to facilitate evaluation and future progress in this area.

041 042 043

044

038

039

040

1 INTRODUCTION

Large text-to-image models have recently made significant progress in generating high-quality images. However, to make these models truly useful, controllability is essential. Users have diverse needs and want to interact with these models in different ways depending on their specific use case. Influential work has been done to enable controllability in these networks, robustly addressing foundational applications and controls such as subject personalization, style learning, layout controls, and semantic controls. Despite this progress, the full potential of these powerful large models has not been fully realized. Some applications that seemed clearly out of reach just a couple of years ago are now possible with careful approaches.

053 We present one such application: *style-aware drag-and-drop*. We formalize this problem and introduce *Magic Insert*, our method to tackle it, which shows strong performance compared to current

baselines. One might initially consider addressing style-aware drag-and-drop by trying to inpaint using a stylized subject, for example by combining Dreambooth Ruiz et al. (2023a), StyleDrop Sohn et al. (2023), and inpainting. We find that approaches of this type are very expensive and achieve subpar results.

058 In developing Magic Insert, we address two interesting sub-problems: *style-aware personalization* 059 and *realistic object insertion in stylized images*. For style-aware personalization, there have been 060 attempts on adjacent problems, such as learning a style and then representing a specific subject in 061 that style Sohn et al. (2023); Hertz et al. (2023), or combining pre-trained custom style and subject 062 models Shah et al. (2023); Frenkel et al. (2024). Recent style work suggests that fast style learning 063 is possible, but fast learning of a subject, including all the intricacies of identity, is a much harder 064 problem that has arguably not been solved yet Ye et al. (2023); Ruiz et al. (2023b); Gal et al. (2023); Wang et al. (2024b). We propose leveraging learnings from both domains and settle on a solution that 065 uses adapter injection of style paired with subject-learning in the embedding and weight space of a 066 diffusion model. 067

One key idea we propose is to not attempt inpainting directly into an image after achieving style-aware personalization. Instead, for best results, we first generate a high-quality subject and then insert that subject into the target image. To achieve our results, we introduce an innovation called *Bootstrap Domain Adaptation*, that allows progressive retargeting of a model's initial distribution to a target distribution. We apply this idea to adapt a subject insertion network that has been trained on real images to perform well on the stylized image domain, enabling the insertion of our generated stylized subject into the background image.

Our method allows the generated output to exhibit strong adherence to the target style while preserving the essence and identity of the subject, and for realistic insertion of the stylized subject into the generated image. The method also provides flexibility in terms of the degree of stylization desired and how closely to adhere to the original subject's specific details and pose (or allow more novelty in the generation).

- ⁰⁸⁰ In summary, we propose the following contributions:
 - We propose and formalize the problem of **style-aware drag-and-drop**, where a subject (a character or object) is dragged from one image into another. Specifically, in our problem formulation the subject reference image and the target image may be in vastly different styles, and the plausibility and realism of the subject insertion is important.
 - In order to encourage exploration into this new problem, we present **SubjectPlop**, a dataset of subjects and backgrounds that span widely different styles and overall semantics. We will release this dataset for public use, as well as our evaluation suite.
 - We propose **Magic Insert**, a method to tackle the style-aware drag-and-drop problem. Our method is composed of a style-aware personalization component and a style-consistent drag-and-drop component.
 - For **style-aware personalization**, we demonstrate strong and consistent results using subjectlearning in the embedding and weight space of a pre-trained diffusion models, along with adapter injection of style.
 - For **drag-and-drop**, we propose *Bootstrapped Domain Adaptation*, a method that allows for progressive retargeting of a model's initial distribution unto a target distribution. We use this to adapt an object insertion network trained on real images to perform well on the stylized image domain.
- 099 100 101

102

081 082

084

085

090

092

093

095 096

097

098

2 RELATED WORK

Text-to-Image Models Recent text-to-image models such as Imagen Saharia et al. (2022b), DALLE 2 Ramesh et al. (2022), Stable Diffusion (SD) Rombach et al. (2022), Muse Chang et al. (2023)
and Parti Yu et al. (2022) have demonstrated remarkable capabilities in generating high-quality
images from text descriptions. They leverage advancements in diffusion models Sohl-Dickstein et al.
(2015); Ho et al. (2020); Song et al. (2022a) and generative transformers. Our work builds on top of
SDXL Podell et al. (2023) and the LDM architecture Rombach et al. (2022).



Figure 2: **Style-Aware Personalization:** To generate a subject that fully respects the style of the target image while also conserving the subject's essence and identity, we (1) personalize a diffusion model in both weight and embedding space, by training LoRA deltas on top of the pre-trained diffusion model and simultaneously training the embedding of two text tokens using the diffusion denoising loss (2) use this personalized diffusion model to generate the style-aware subject by embedding the style of the target image and conducting adapter style-injection into select upsampling layers of the model during denoising.

Image Inpainting The task of filling masked pixels of a target image has been explored using a
wide range of approaches: Generative adversarial networks Goodfellow et al. (2014) e.g. Pathak
et al. (2016); Hui et al. (2020); Liu et al. (2020); Ntavelis et al. (2020); Ren et al. (2019); Zeng et al.
(2019) and end-to-end learning methods Iizuka et al. (2017); Liu et al. (2018); Suvorov et al. (2022);
Wu et al. (2022). More recently, diffusion models enabled significant progress Meng et al. (2021b);
Lugmayr et al. (2022); Wang et al. (2023); Saharia et al. (2022a); Avrahami et al. (2022). Such inpainting methods are a precursor to many object insertion approaches.

138 139

Generative Object Insertion The problem of inserting an object into an existing scene has been 140 originally explored using Generative Adversarial Networks (GANs) Goodfellow et al. (2014). Lee 141 et al. (2018) breaks down the task into two generative modules, one determines where the inserted 142 object mask should be and the other determines what the mask shape and pose. ShadowGAN Zhang 143 et al. (2019) addresses the need to add a shadow cast by the inserted object, leveraging 3D rendering 144 for training data. More recent works use diffusion models. Paint-By-Example Yang et al. (2023) 145 allows inpainting a masked area of the target image with reference to the object source image, but it only preserves semantic information and has low fidelity to the original object's identity. Recent 146 work also explores swapping objects in a scene while harmonizing, but focuses on swapping areas 147 of the image which were previously populated Gu et al. (2024). There also exists an array of work 148 that focuses on inserting subjects or concepts in a scene either by inpainting Safaee et al. (2024); 149 Lu et al. (2023) or by other means Song et al. (2022b); Sarukkai et al. (2024) - these do not handle 150 large style adaptation and inpainting methods usually suffer from problems with insertion such as 151 background removal, incomplete insertion and low quality results. ObjectDrop Winter et al. (2024) 152 trains a diffusion model for object removal/insertion using a counterfactual dataset captured in the 153 real world. The trained model can insert segmented objects into real images with contextual cues 154 such as shadows and reflections. We build upon this novel and incredibly useful paradigm by tackling 155 the challenging domain of stylized images instead.

156

Personalization, Style Learning and Controllability Text-to-image models enable users to provide text prompts and sometimes input images as conditioning input, but do not allow for fine-grained control over subject, style, layout, etc. Textual Inversion Gal et al. (2022) and DreamBooth Ruiz et al. (2023a) are pioneering works that demonstrated personalization of such models to generate images of specific subjects, given few casual images as input. Textual Inversion Gal et al. (2022) and follow-up techniques such as P+ Voynov et al. (2023) optimize text embeddings, while DreamBooth

2 **Subject** Insertion Model

169 Figure 3: Subject Insertion: In order to insert the style-aware personalized generation, we (1) copy-paste a 170 segmented version of the subject onto the target image (2) run our subject insertion model on the deshadowed 171 image - this creates context cues and realistically embeds the subject into the image including shadows and reflections. 172

173 174

185

186 187

188

191

193

194

196

197

199 200 201

162

163 164

165

166

167 168

optimizes the model weights. This type of work has also been extended to 3D models Raj et al. 175 (2023), scene completion Tang et al. (2023) and others. There also exists work on fast subject-driven 176 generation Chen et al. (2024); Ruiz et al. (2023b); Gal et al. (2023); Wang et al. (2024b); Arar et al. 177 (2023). Other work allows for conditioning on new modalities such as ControlNet Zhang et al. (2023) 178 and on image features (IP-Adapter Ye et al. (2023)). There is a body of work that dives more deeply 179 into style learning and generating consistent style as well with StyleDrop Sohn et al. (2023) as a pioneer, with newer work that achieves fast stylization Shah et al. (2023); Wang et al. (2024a); Hertz 181 et al. (2023); Rout et al. (2024), or combines subject models with style models like ZipLoRA Shah et al. (2023) and others Frenkel et al. (2024). Our work leverages ideas from Textual Inversion, 182 DreamBooth and IP-Adapter to unlock style-aware personalization prior and combine it with subject 183 insertion. 184

3 METHOD

STYLE-AWARE DRAG-AND-DROP PROBLEM FORMULATION 3.1

189 We formalize the style-aware drag-and-drop problem as follows. Let \mathcal{I}_s and \mathcal{I}_t denote the space of 190 subject and target images, respectively. The space of subject images consists of images of solely the subject in front of plain backgrounds. Given a subject image $x_s \in \mathcal{I}_s$ and a target image $x_t \in \mathcal{I}_t$, our goal is to generate a new image $\hat{x}_t \in \mathcal{I}_t$ such that: 192

- 1. The subject from x_s is inserted into \hat{x}_t in a semantically consistent and realistic manner, accounting for factors such as occlusion, shadows, and reflections.
- 2. The inserted subject in \hat{x}_t adopts the style characteristics of the target image x_t while preserving its essential identity and attributes from x_s .

Formally, we aim to learn a function $h : \mathcal{I}_s \times \mathcal{I}_t \to \mathcal{I}_t$ that satisfies:

$$h(x_s, x_t) = \hat{x}_t \quad \text{s.t.} \quad \hat{x}_t \sim p(\hat{x}_t | x_t, x_s) \tag{1}$$

where $p(\hat{x}_t|x_t, x_s)$ represents the conditional distribution of the output image given the subject 202 and target images. This distribution encapsulates the desired properties of semantic consistency, 203 realistic insertion, and style adaptation. To learn the function h, we decompose the problem into 204 two sub-tasks: style-aware personalization and realistic object insertion in stylized images. Style-205 aware personalization focuses on generating a subject that adheres to the target image's style while 206 maintaining its identity. Realistic object insertion aims to seamlessly integrate the stylized subject 207 into the target image, accounting for the scene's geometry and lighting conditions. By addressing 208 these sub-tasks, we can effectively solve the style-aware drag-and-drop problem and generate visually 209 coherent and compelling results. In the following sections, we present our dataset and the components 210 of our proposed method.

211

212 3.2 SUBJECTPLOP DATASET

213

To facilitate the evaluation of the style-aware drag-and-drop problem, we introduce the SubjectPlop 214 dataset and make it publicly available. As this is a novel problem, a dedicated dataset is crucial for 215 enabling the research community to make progress in this area.



231 Figure 4: Bootstrapped Domain Adaptation: Surprisingly, a diffusion model trained for subject inser-232 tion/removal on data captured in the real world can generalize to images in the wider stylistic domain in a limited 233 fashion. We introduce *bootstrapped domain adaptation*, where a model's effective domain can be adapted by 234 using a subset of its own outputs. (left) Specifically, we use a subject removal/insertion model to first remove subjects and shadows from a dataset from our target domain. Then, we filter flawed outputs, and use the filtered 235 set of images to retrain the subject removal/insertion model. (right) We observe that, the initial distribution 236 (blue) changes after training (purple) and initially incorrectly treated images (red samples) are subsequently 237 correctly treated (green). When doing bootstrapped domain adaptation, we train on only the initially correct 238 samples (green).

SubjectPlop consists of a diverse collection of subjects generated using DALL-E3 Ramesh et al. (2022) and backgrounds generated using the open-source SDXL model Podell et al. (2023). The dataset includes various subject types, such as animals and fantasy characters, and both subjects and backgrounds exhibit a wide range of styles, including 3D, cartoon, anime, realistic, and photographic. The diversity in color hues and lighting conditions ensures comprehensive coverage of different scenarios for evaluation. No real people are represented in the dataset.

The dataset comprises 20 distinct backgrounds and 35 unique subjects, allowing for a total of 700 possible subject-background pairs. The entire dataset is meant for evaluation of the task. This rich set of test cases enables the assessment of performance and generalization capabilities of style-aware drag-and-drop techniques. By introducing SubjectPlop, we aim to provide a standardized benchmark for evaluating and comparing different approaches to the style-aware drag-and-drop problem. We believe this dataset will serve as a valuable resource for researchers and practitioners working in image manipulation and generation, fostering further advancements in this area.

254 255

266 267

3.3 STYLE-AWARE PERSONALIZATION

Our style-aware personalization approach is illustrated in Figure 2. Let f_{θ} denote a pre-trained diffusion model with parameters θ . Given a subject image $x_s \in \mathcal{I}_s$, our method personalizes f_{θ} on x_s in both the weight and embedding space, similar to DreamBooth Ruiz et al. (2023a) and Textual Inversion Gal et al. (2022).

In the first step, we train LoRA Hu et al. (2021) (Low-Rank Adaptation) deltas Δ_{θ} to produce an efficiently fine-tuned adapted model $f_{\theta'}$ where $\theta' = \theta + \Delta_{\theta}$, while preserving the model's original capabilities. Simultaneously, we learn embeddings $e_1, e_2 \in \mathbb{R}^d$ for two personalized text tokens, where *d* is the embedding dimensionality. We use two learned embeddings since we found better performance for both subject preservation and editability in this configuration. The LoRA deltas and and embeddings are jointly trained using the diffusion denoising loss:

$$\mathcal{L}_{\text{joint}} = \mathbb{E}_{t,\epsilon} \left[\|\epsilon - \epsilon_{\theta'}(x_s^t, t, [e_1; e_2])\|_2^2 \right]$$
(2)

where $t \sim \mathcal{U}(0,1)$, $\epsilon \sim \mathcal{N}(0,\mathbf{I})$, $x_s^t = \sqrt{\bar{\alpha}_t}x_s + \sqrt{1-\bar{\alpha}_t}\epsilon$, and $\epsilon_{\theta'}$ is the noise prediction of the adapted model $f_{\theta'}$. The joint optimization of Δ_{θ} , e_1 , and e_2 is performed using the loss $\mathcal{L}_{\text{joint}}$. These personalized text tokens $[e_1; e_2]$ serve as a compact representation of the subject's identity. By performing embedding and weight-space learning simultaneously, We find that performing embedding
 and weight-space learning simultaneously, with two text tokens, captures the subject's identity more
 strongly while allowing sufficient editability to introduce the target style.

In the second step, we leverage the personalized diffusion model $f_{\theta'}$ to generate the style-aware subject \hat{x}_s . To infuse the target image x_t 's style into \hat{x}_s , we employ style injection. Specifically, we generate a style embedding $e_t = \text{CLIP}(x_t)$ of x_t using a frozen CLIP encoder CLIP. We then use a frozen IP-Adapter model v to inject e_t into a subset of the UNet blocks of $f_{\theta'}$ during inference:

$$\hat{x}_s = f_{\theta'}([e_1; e_2], v(e_t)) \tag{3}$$

This approach is similar to InstantStyle Wang et al. (2024a), with injection into the upsample block that is adjacent to the midblock, with some key differences being omitting content/style embedding separation, and injecting into a personalized model. To the best of our knowledge, our central idea of combining adapter injection and personalized models remains unexplored in the published literature. This ensures that \hat{x}_s maintains the subject's identity while adopting x_t 's style characteristics.

By combining style-aware personalization with style injection, our method generates subjects that
harmoniously blend into the target image while retaining their essential identity, effectively tackling
the first challenge of style-aware drag-and-drop and enabling the creation of visually coherent and
style-consistent results.

289 3.4 BOOTSTRAPPED DOMAIN ADAPTATION FOR SUBJECT INSERTION

In this section, we address the problem of subject insertion and propose a novel solution using 291 bootstrapped domain adaptation. We formalize the concept of bootstrapped domain adaptation and 292 describe the dataset used for this purpose. Subject insertion is a crucial component of the style-293 aware drag-and-drop problem, as it involves seamlessly integrating a stylized subject into a target 294 background image. While diffusion-based inpainting approaches Meng et al. (2021a); Saharia et al. 295 (2022b); Rombach et al. (2022) can be used for this, they still face challenges such as generating 296 content in smooth regions, producing incomplete figures, erasing objects behind inserted subjects, and 297 having problems with boundary harmonization. We take a simpler and stronger approach, which is to 298 insert the subject by copying and pasting it into the target image, and then subsequently generating 299 contextual cues such as shadows and reflections Winter et al. (2024) in a second step. Unfortunately, existing subject insertion models are trained on data captured in the real world, severely limiting their 300 ability to generalize to images with diverse artistic styles. 301

302 Let \mathcal{D}_r denote the distribution of real-world images and \mathcal{D}_s denote the distribution of stylized images. 303 Existing subject insertion models are trained on samples from \mathcal{D}_r , but our goal is to adapt them to 304 perform well on samples from \mathcal{D}_s . To overcome this limitation, we introduce bootstrapped domain adaptation, a technique that enables a model to adapt its effective domain by leveraging a subset of its 305 own outputs. As illustrated in Figure 4 (left), we employ a subject removal/insertion model q_{θ} trained 306 on real-data (Winter et al. (2024) in our case) to first remove subjects and shadows from a dataset 307 $S \sim D_s$ belonging to our target domain. Subsequently, we filter out flawed outputs and obtain a 308 filtered set of images $\mathcal{S}' \subset \mathcal{S}$, which we use to retrain the subject removal/insertion model. Filtering 309 can be done using human feedback or automatically given a quality evaluation module. 310

311 The bootstrapped domain adaptation process can be formalized as follows:

312

278

288

290

$$\omega = \arg\min_{\omega} \mathbb{E}_{(x,y)\sim\mathcal{S}'} \mathcal{L}(g_{\omega}(x), y) \tag{4}$$

where ω denotes the adapted model parameters, \mathcal{L} is the diffusion denoising loss, and (x, y) are pairs of input images and corresponding subject removal/insertion ground truths from the filtered set \mathcal{S}_f . The concept of bootstrapped domain adaptation is based on the surprising observation that a diffusion model trained for subject insertion/removal on real-world data can generalize to a wider stylistic domain to a limited extent. By retraining the model on its own filtered outputs, we can effectively adapt its domain to better handle stylized images.

Figure 4 (right) demonstrates the effect of bootstrapped domain adaptation on the model's distribution. The initial distribution, represented as $p_{\omega}(x)$, evolves after training, becoming $p_{\omega^*}(x)$. Images that were initially treated incorrectly, shown as samples from $\mathcal{D}_s \setminus \mathcal{S}'$, are subsequently handled correctly, as indicated by their inclusion in \mathcal{S}' . During the bootstrapped domain adaptation process, we train the model only on the initially correct samples from \mathcal{S}' to further refine its performance on the target



Figure 5: **Result Gallery:** Examples of our Magic Insert method for different subjects and backgrounds with vastly different styles.

Table 1: **Subject Fidelity Comparisons.** We compare our method for subject fidelity (DINO, CLIP-I, CLIP-T Simple, CLIP-T Detailed) across different methods. Our method variants show high subject fidelity.

Method	DINO \uparrow	CLIP-I↑	CLIP-T Simple \uparrow	CLIP-T Detailed \uparrow	Overall Mean \uparrow
StyleAlign Prompt	0.223	0.743	0.266	0.299	0.383
StyleAlign ControlNet	0.414	0.808	0.289	0.294	0.451
InstantStyle Prompt	0.231	0.778	0.283	0.300	0.398
InstantStyle ControlNet	0.446	0.806	0.281	0.283	0.454
Ours	0.295	0.829	0.276	0.293	0.423
Ours ControlNet	0.514	0.869	0.289	0.308	0.495

domain. Several steps of bootstrapped domain adaptation can be performed, further enhancing the model's performance. In our work we find that one step suffices, with a small set of samples (around 50). Figure 7 shows results with and without bootstrap domain adaptation.

To facilitate the bootstrapped domain adaptation process, we curate a dataset S specifically tailored to this task. The dataset comprises a diverse range of stylized images, selected to represent the target domain \mathcal{D}_s . In our case, this dataset is constructed by sampling from different text-to-image generative models with diverse prompts that elicit prominent subjects with shadows and reflections in a variety of global styles. By finetuning the subject removal/insertion model on this dataset using the bootstrapped domain adaptation technique, we enable it to effectively handle subject insertion in the context of style-aware drag-and-drop.

4 EXPERIMENTS

In this section, we show experiments and applications. Our full method enables insertion of arbitrary
subjects into images with diverse styles, with a large expanse of text-guided semantic modifications.
Specifically, not only does the subject retain its identity and essence while inheriting the style of the
target image, but we can modify key subject characteristics such as the pose and other core attributes
such as adding accessories, changing appearance, changing shapes, or even species hybrids (see
appendix). These changes can be integrated with components such as LLMs that allow for automatic affordances and environment interactions (Figure 6).



Figure 6: **LLM-Guided Affordances:** Examples of an LLM-guided pose modification for Magic Insert, with the LLM suggesting plausible poses and environment interactions for areas of the image and Magic Insert generating and inserting the stylized subject with the corresponding pose into the image.

Table 2: **Style Fidelity Comparisons.** We compare our method for style fidelity (CLIP-I, CSD, CLIP-T). Our method variants show strong style-following.

Method	CLIP-I \uparrow	$\mathbf{CSD}\uparrow$	$\text{CLIP-T} \uparrow$	$Overall Mean \uparrow$
StyleAlign Prompt	0.570	0.150	0.248	0.323
StyleAlign ControlNet	0.575	0.188	0.274	0.345
InstantStyle Prompt	0.583	0.312	0.276	0.390
InstantStyle ControlNet	0.588	0.334	0.279	0.400
Ours	0.560	0.243	0.268	0.357
Ours ControlNet	0.575	0.294	0.274	0.381

4.1 STYLE-AWARE DRAG-AND-DROP RESULTS

Magic Insert Results We present a gallery of qualitative results in Figure 5 to highlight the effectiveness and versatility of our method. The examples span a wide range of subjects and target backgrounds with vastly different artistic styles, from photorealistic scenes to cartoons, and paintings. For style-aware personalization we use the SDXL model Podell et al. (2023), and for subject insertion we use our trained subject insertion model based on a latent diffusion model architecture.

In each case, our method successfully extracts the subject from the source image and blends it into
 the target background, adapting the subject's appearance to match the background's style. Notice
 how the inserted subjects take on the colors, textures, and stylistic elements of the target images. The
 coherent shadows and reflections enhance the plausibility of the results.

LLM-Guided Affordances Our proposed style-aware personalization method allows for large changes in character pose, with support from the diffusion model prior. Using and LLM (ChatGPT 40) we are able to generate LLM-guided affordances for different subjects, by feeding an instruction

Table 3: ImageReward Metric Comparisons. We compare different methods using the ImageReward metric,
 which correlates with human preference for aesthetic evaluation. Higher scores indicate better performance. Our
 variants outperform all benchmarks

Method	ImageReward Score \uparrow
StyleAlign Prompt	-1.1942
StyleAlign ControlNet	-0.5180
InstantStyle Prompt	-0.4638
InstantStyle ControlNet	-0.2759
Ours	-0.2108
Ours ControlNet	-0.1470

Table 4: User Study. This study evaluates our method against two different baselines (StyleAlign ControlNet and InstantStyle ControlNet) based on subject identity, style fidelity, and realistic insertion. Participants ranked each method by preference.

Method	User Preference \uparrow
Ours over StyleAlign ControlNet	85%
Ours over InstantStyle ControlNet	80%

prompt, the full background image, and the section of the background image in which the character
will be positioned. Using these LLM suggestions, we can generate the character following these poses
and environment interactions and insert it in the appropriate space. With this, we show in Figure 6 a
first attempt at the previously unassailable task of inserting subjects into images realistically with
automatic interactions with the scene.

Bootstrap Domain Adaptation We show in Figure 4 a sample case of subject insertion with an insertion model that is trained on real images without adaptation, and on the same model that uses our proposed bootstrap domain adaptation on a small set of 50 samples. Insertion without bootstrap domain adaptation generates subpar results, with problems such as missing shadows, reflections and even added distortions.

4.2 Comparisons

463 Here we introduce baselines, as well as quantitative and qualitative comparisons, as well as a 464 user study. Specifically, our proposed baselines utilize the StyleAlign Hertz et al. (2023) and 465 InstantStyle Wang et al. (2024a) stylization methods, which can generate images in reference styles 466 given either inversion or embedding of the reference image. We combine these methods with either 467 sufficiently detailed prompting guided by a VLM (ChatGPT 4) or edge-conditioned ControlNet. 468 For prompting we use the VLM to describe the subjects while eliminating style cues, and for edgeconditioning we use Canny edges extracted from the subject reference images to guide the stylized 469 outputs using ControlNet. 470

471

455

461

462

Baseline Comparisons We run studies in order to compare the performance of subject stylization
 for different baselines and our style-aware personalization method. We study the performance of
 these methods on subject fidelity, style fidelity, and human preference.

For subject fidelity (Table 1), our proposed variants achieve high scores across various subject fidelity
metrics (DINO, CLIP-I, CLIP-T Simple, CLIP-T Detailed). DINO and CLIP-I metrics are identical
to those presented in DreamBooth Ruiz et al. (2023a) and CLIP-T Simple / Detailed denotes the
CLIP similarity between the output image CLIP embedding and the CLIP embedding of simple and
detailed text prompts describing the subject, which are in turn generated by ChatGPT 4.

Regarding style fidelity (Table 2), our proposed variants demonstrate strong style-following performance using CLIP-I Ruiz et al. (2023a); Sohn et al. (2023), CSD Somepalli et al. (2024), CLIP-T Ruiz et al. (2023a); Sohn et al. (2023) metrics. For style fidelity, InstantStyle ControlNet outperforms our variants using these automatic metrics, although we observe that subject details and contrast is lost in many of these samples as shown in Figure 8. For this, we also compute ImageReward Xu et al. (2024) scores in Table 3, which correlate strongly with human preference in aesthetic evaluations. We observe that our variants strongly outperform the benchmarks.



Figure 8: Style-Aware Personalization Baseline Comparison: We show some comparisons of our style-aware 509 personalization method with respect to the top performing baselines StyleAlign + ControlNet and InstantStyle + ControlNet. We can see that the baselines can yield decent outputs, but lag behind our style-aware personalization 510 method in overall quality. In particular InstantStyle + ControlNet outputs often appear slightly blurry and don't 511 capture subject features with good contrast. 512

Moreover, finding strong quantitative metrics for subject fidelity and for style fidelity is an open 514 problem in the field, and metrics can have strong biases that can make them suboptimal. Again, we 515 show some examples for our proposed style-aware personalization, along with top baseline contenders 516 StyleAlign ControlNet and InstantStyle ControlNet in Figure 8. We observe that the generation quality 517 of our variants is stronger than the benchmarks, especially with both strong stylization performance 518 while still retaining the essence of the subjects. Our Magic Insert + ControlNet variant is powerful 519 given that it exactly follows the outline of the character, and thus has the strongest subject fidelity 520 over all approaches, although it does not have the desirable properties of our method w/o ControlNet 521 which include pose, form and attribute modification of the subject. We further the discussion on subject fidelity vs. editability tradeoff in the appendix. 522

523 524

525

527

529

530

531

User Study Following previous work Ruiz et al. (2023a); Sohn et al. (2023); Ruiz et al. (2023b); Tang et al. (2023) we perform a robust user study to compare our full method (w/ ControlNet) with the strongest baselines: StyleAlign ControlNet and InstantStyle ControlNet. We recruit a total of 526 60 users (4 sets of 15 users) to answer 40 evaluation tasks (2 sets of 20 tasks) for each baseline comparison (2 baseline comparisons). We collect a total of 1200 user evaluations. We ask users to 528 rank their preferred methods with respect to subject identity preservation, style fidelity with respect to the background image, and realistic insertion of the subject into the background image. We show the results in Table 4. We observe a strong preference of users for our generated outputs compared to baselines.

- 532
- 534

CONCLUSION 5

In this work, we introduced the problem of style-aware drag-and-drop, a new challenge in image 536 generation that enables intuitive subject insertion while maintaining style consistency. We proposed Magic Insert, a method combining style-aware personalization and insertion through bootstrapped domain adaptation, which outperforms baselines in style adherence and insertion realism. To support 538 further research, we presented the SubjectPlop dataset, featuring subjects and backgrounds across diverse styles and semantics.

540 REFERENCES 541

542 543 544	Moab Arar, Rinon Gal, Yuval Atzmon, Gal Chechik, Daniel Cohen-Or, Ariel Shamir, and Amit H. Bermano. Domain-agnostic tuning-encoder for fast personalization of text-to-image models. In <i>SIGGRAPH Asia 2023 Conference Papers</i> , pp. 1–10, 2023.
545 546 547	Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural images. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 18208–18218, 2022.
548 549 550 551	Huiwen Chang, Han Zhang, Jarred Barber, AJ Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Murphy, William T Freeman, Michael Rubinstein, et al. Muse: Text-to-image generation via masked generative transformers. arXiv preprint arXiv:2301.00704, 2023.
552 553 554	Wenhu Chen, Hexiang Hu, Yandong Li, Nataniel Ruiz, Xuhui Jia, Ming-Wei Chang, and William W Cohen. Subject-driven text-to-image generation via apprenticeship learning. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
555 556 557	Yarden Frenkel, Yael Vinker, Ariel Shamir, and Daniel Cohen-Or. Implicit style-content separation using b-lora. <i>arXiv preprint arXiv:2403.14572</i> , 2024.
558 559 560	Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. <i>arXiv preprint arXiv:2208.01618</i> , 2022.
561 562 563	Rinon Gal, Moab Arar, Yuval Atzmon, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. Encoder-based domain tuning for fast personalization of text-to-image models. <i>ACM Transactions</i> <i>on Graphics (TOG)</i> , 42(4):1–13, 2023.
564 565 566 567	Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. <i>Advances in neural information processing systems</i> , 27, 2014.
568 569 570	Jing Gu, Yilin Wang, Nanxuan Zhao, Wei Xiong, Qing Liu, Zhifei Zhang, He Zhang, Jianming Zhang, HyunJoon Jung, and Xin Eric Wang. Swapanything: Enabling arbitrary object swapping in personalized visual editing. <i>arXiv preprint arXiv:2404.05717</i> , 2024.
571 572 573	Amir Hertz, Andrey Voynov, Shlomi Fruchter, and Daniel Cohen-Or. Style aligned image generation via shared attention. <i>arXiv preprint arXiv:2312.02133</i> , 2023.
574	Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. 2020.
575 576 577 578	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. <i>arXiv preprint arXiv:2106.09685</i> , 2021.
579 580	Zheng Hui, Jie Li, Xiumei Wang, and Xinbo Gao. Image fine-grained inpainting. <i>arXiv preprint arXiv:2002.02609</i> , 2020.
581 582 583	Satoshi Iizuka, Edgar Simo-Serra, and Hiroshi Ishikawa. Globally and locally consistent image completion. <i>ACM Transactions on Graphics (ToG)</i> , 36(4):1–14, 2017.
584 585 586	Donghoon Lee, Sifei Liu, Jinwei Gu, Ming-Yu Liu, Ming-Hsuan Yang, and Jan Kautz. Context- aware synthesis and placement of object instances. <i>ArXiv</i> , abs/1812.02350, 2018. URL https: //api.semanticscholar.org/CorpusID:53973622.
587 588 589 590	Guilin Liu, Fitsum A Reda, Kevin J Shih, Ting-Chun Wang, Andrew Tao, and Bryan Catanzaro. Image inpainting for irregular holes using partial convolutions. In <i>Proceedings of the European</i> <i>conference on computer vision (ECCV)</i> , pp. 85–100, 2018.
591 592 593	Hongyu Liu, Bin Jiang, Yibing Song, Wei Huang, and Chao Yang. Rethinking image inpainting via a mutual encoder-decoder with feature equalizations. In <i>Computer Vision–ECCV 2020: 16th</i> <i>European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16</i> , pp. 725–741. Springer, 2020.

594 Lingxiao Lu, Bo Zhang, and Li Niu. Dreamcom: Finetuning text-guided inpainting model for image 595 composition. arXiv preprint arXiv:2309.15508, 2023. 596 Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. 597 Repaint: Inpainting using denoising diffusion probabilistic models. In Proceedings of the 598 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 11461–11471, 2022. 600 Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. 601 Sdedit: Guided image synthesis and editing with stochastic differential equations. arXiv preprint arXiv:2108.01073, 2021a. 602 603 Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. 604 Sdedit: Guided image synthesis and editing with stochastic differential equations. arXiv preprint 605 arXiv:2108.01073, 2021b. 606 Evangelos Ntavelis, Andrés Romero, Siavash Bigdeli, Radu Timofte, Zheng Hui, Xiumei Wang, 607 Xinbo Gao, Chajin Shin, Taeoh Kim, Hanbin Son, et al. Aim 2020 challenge on image extreme 608 inpainting. In Computer Vision–ECCV 2020 Workshops: Glasgow, UK, August 23–28, 2020, 609 Proceedings, Part III 16, pp. 716–741. Springer, 2020. 610 611 Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context 612 encoders: Feature learning by inpainting. In Proceedings of the IEEE conference on computer 613 vision and pattern recognition, pp. 2536-2544, 2016. 614 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe 615 Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image 616 synthesis. arXiv preprint arXiv:2307.01952, 2023. 617 Amit Raj, Srinivas Kaza, Ben Poole, Michael Niemeyer, Nataniel Ruiz, Ben Mildenhall, Shiran 618 Zada, Kfir Aberman, Michael Rubinstein, Jonathan Barron, et al. Dreambooth3d: Subject-driven 619 text-to-3d generation. In Proceedings of the IEEE/CVF international conference on computer 620 vision, pp. 2349–2359, 2023. 621 622 Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-623 conditional image generation with clip latents. arXiv preprint arXiv:2204.06125, 2022. 624 Yurui Ren, Xiaoming Yu, Ruonan Zhang, Thomas H Li, Shan Liu, and Ge Li. Structureflow: Image 625 inpainting via structure-aware appearance flow. In Proceedings of the IEEE/CVF international 626 conference on computer vision, pp. 181–190, 2019. 627 628 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-629 resolution image synthesis with latent diffusion models. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 10684–10695, 2022. 630 631 Litu Rout, Yujia Chen, Nataniel Ruiz, Abhishek Kumar, Constantine Caramanis, Sanjay Shakkottai, 632 and Wen-Sheng Chu. Rb-modulation: Training-free personalization of diffusion models using 633 stochastic optimal control. arXiv preprint arXiv:2405.17401, 2024. 634 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. 635 Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In 2023 636 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 22500–22510. 637 IEEE, 2023a. 638 639 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Wei Wei, Tingbo Hou, Yael Pritch, Neal Wadhwa, 640 Michael Rubinstein, and Kfir Aberman. Hyperdreambooth: Hypernetworks for fast personalization of text-to-image models. arXiv preprint arXiv:2307.06949, 2023b. 641 642 Mehdi Safaee, Aryan Mikaeili, Or Patashnik, Daniel Cohen-Or, and Ali Mahdavi-Amiri. Clic: 643 Concept learning in context. In Proceedings of the IEEE/CVF Conference on Computer Vision and 644 Pattern Recognition, pp. 6924–6933, 2024. 645 Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, 646 and Mohammad Norouzi. Palette: Image-to-image diffusion models. In ACM SIGGRAPH 2022 647 Conference Proceedings, pp. 1–10, 2022a.

648 649 650 651	Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. <i>Advances in Neural Information Processing Systems</i> , 35:36479–36494, 2022b.
652 653 654 655	Vishnu Sarukkai, Linden Li, Arden Ma, Christopher Ré, and Kayvon Fatahalian. Collage diffusion. In <i>Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)</i> , pp. 4208–4217, January 2024.
656 657	Viraj Shah, Nataniel Ruiz, Forrester Cole, Erika Lu, Svetlana Lazebnik, Yuanzhen Li, and Varun Jampani. Ziplora: Any subject in any style by effectively merging loras. 2023.
658 659 660	Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. 2015.
661 662 663 664	Kihyuk Sohn, Nataniel Ruiz, Kimin Lee, Daniel Castro Chin, Irina Blok, Huiwen Chang, Jarred Barber, Lu Jiang, Glenn Entis, Yuanzhen Li, et al. Styledrop: Text-to-image generation in any style. In <i>37th Conference on Neural Information Processing Systems (NeurIPS)</i> . Neural Information Processing Systems Foundation, 2023.
665 666 667	Gowthami Somepalli, Anubhav Gupta, Kamal Gupta, Shramay Palta, Micah Goldblum, Jonas Geiping, Abhinav Shrivastava, and Tom Goldstein. Measuring style similarity in diffusion models. <i>arXiv preprint arXiv:2404.01292</i> , 2024.
668 669	Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. 2022a.
670 671 672	Yizhi Song, Zhifei Zhang, Zhe Lin, Scott Cohen, Brian Price, Jianming Zhang, Soo Ye Kim, and Daniel Aliaga. Objectstitch: Generative object compositing. arXiv preprint arXiv:2212.00932, 2022b.
673 674 675 676 677	Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin, Anastasia Remizova, Arsenii Ashukha, Aleksei Silvestrov, Naejin Kong, Harshith Goka, Kiwoong Park, and Victor Lempitsky. Resolution-robust large mask inpainting with fourier convolutions. In <i>Proceedings of the IEEE/CVF winter conference on applications of computer vision</i> , pp. 2149–2159, 2022.
678 679 680	Luming Tang, Nataniel Ruiz, Qinghao Chu, Yuanzhen Li, Aleksander Holynski, David E Jacobs, Bharath Hariharan, Yael Pritch, Neal Wadhwa, Kfir Aberman, et al. Realfill: Reference-driven generation for authentic image completion. <i>arXiv preprint arXiv:2309.16668</i> , 2023.
681 682	Andrey Voynov, Qinghao Chu, Daniel Cohen-Or, and Kfir Aberman. <i>p</i> +: Extended textual conditioning in text-to-image generation. <i>arXiv preprint arXiv:2303.09522</i> , 2023.
684 685	Haofan Wang, 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.
686 687	Qixun Wang, Xu Bai, Haofan Wang, Zekui Qin, and Anthony Chen. Instantid: Zero-shot identity- preserving generation in seconds. <i>arXiv preprint arXiv:2401.07519</i> , 2024b.
689 690 691 692	Su Wang, Chitwan Saharia, Ceslee Montgomery, Jordi Pont-Tuset, Shai Noy, Stefano Pellegrini, Yasumasa Onoe, Sarah Laszlo, David J Fleet, Radu Soricut, et al. Imagen editor and editbench: Advancing and evaluating text-guided image inpainting. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 18359–18369, 2023.
693 694 695	Daniel Winter, Matan Cohen, Shlomi Fruchter, Yael Pritch, Alex Rav-Acha, and Yedid Hoshen. Objectdrop: Bootstrapping counterfactuals for photorealistic object removal and insertion. <i>arXiv</i> preprint arXiv:2403.18818, 2024.
696 697 698 699	Chenfei Wu, Jian Liang, Xiaowei Hu, Zhe Gan, Jianfeng Wang, Lijuan Wang, Zicheng Liu, Yuejian Fang, and Nan Duan. Nuwa-infinity: Autoregressive over autoregressive generation for infinite visual synthesis. <i>arXiv preprint arXiv:2207.09814</i> , 2022.
700 701	Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation. <i>Advances</i> <i>in Neural Information Processing Systems</i> , 36, 2024.

702 703 704 705	Binxin Yang, Shuyang Gu, Bo Zhang, Ting Zhang, Xuejin Chen, Xiaoyan Sun, Dong Chen, and Fang Wen. Paint by example: Exemplar-based image editing with diffusion models. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 18381–18391, 2023.
705 706 707	Hu Ye, Jun Zhang, Sibo Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models. 2023.
708 709 710	Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content- rich text-to-image generation. <i>arXiv preprint arXiv:2206.10789.2022</i>
711 712 713 714	 Yanhong Zeng, Jianlong Fu, Hongyang Chao, and Baining Guo. Learning pyramid-context encoder network for high-quality image inpainting. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i>, pp. 1486–1494, 2019.
715 716 717	Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pp. 3836–3847, 2023.
718 719 720	Shuyang Zhang, Runze Liang, and Miao Wang. Shadowgan: Shadow synthesis for virtual objects with conditional adversarial networks. <i>Computational Visual Media</i> , 5:105–115, 2019.
721 722 723	
724 725	
726 727	
728 729 730	
731 732	
733 734	
735 736 737	
738 739	
740 741 742	
743 744	
745 746	
747 748 749	
750 751	
752 753	
754 755	

756 A APPENDIX

758 Semantic Modifications of Subject Our method inherits all benefits of DreamBooth Ruiz et al.
 (2023a) and thus allows for modification of subject characteristics such as pose, adding accessories, changing appearance, shapeshifting and hybrids. We show some examples in Figure 9. The generated subjects can then be inserted into the background image.

Editability / Fidelity Tradeoff Our method (w/o ControlNet) also inherits DreamBooth's editability
/ fidelity tradeoff. Specifically, the longer the personalization training, the stronger the subject fidelity
but the lesser the editability. This phenomenon is shown in Figure 10. In most cases a sweet spot can
be found for different applications. For our work we use 600 iterations with batch size 1, a learning
rate of 1e-5 and weight decay of 0.3 for the UNet. We also train the text encoder with a learning rate
of 1e-3 and weight decay of 0.1.



Figure 9: **Style-Aware Personalization with Attribute Modification:** Our method allows us to modify key attributes for the subject, such as the ones reflected in this figure, while consistently applying our target style over the generations. This allows us to reinvent the character, or add accessories, which gives large flexibility for creative uses. Note that when using ControlNet this capability disappears.



Figure 10: **Editability / Fidelity Tradeoff:** We show the phenomenon of editability / fidelity tradeoff by showing generations for different finetuning iterations of the space marine (shown above the images) with the "green ship" stylization and additional text prompting "sitting down on the floor". When the style-aware personalized model is finetuned for longer on the subject, we get stronger fidelity to the subject but have less flexibility on editing the pose or other semantic properties of the subject. This can also translate to style editability.