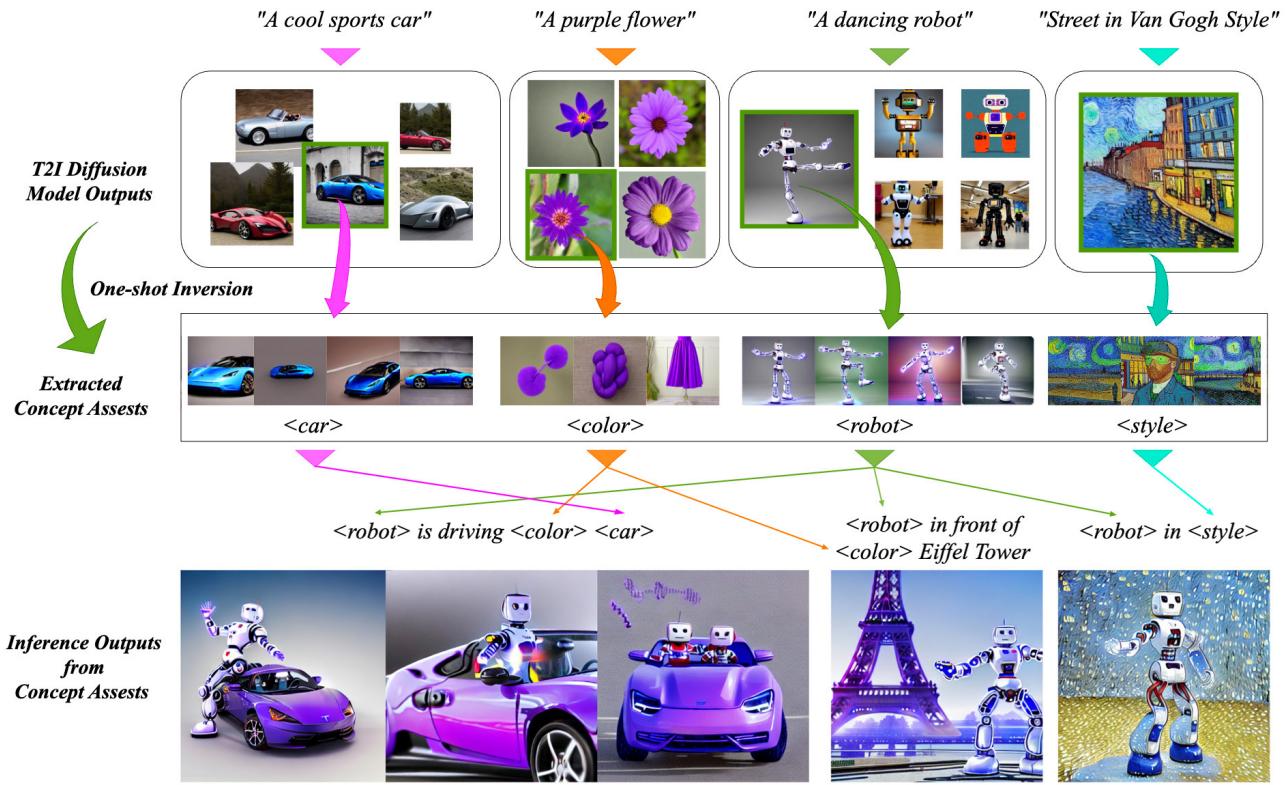


# 1 Concept Assets: Naming Custom Concepts in Diffusion Network Outputs

2 ANONYMOUS AUTHOR(S)

3 SUBMISSION ID: 1055



35 Fig. 1. Concept Assets: The diffusion model takes user-provided text prompts and produces a range of outputs. The user selects their favorite from these (top).  
36 Subsequently, we apply our one-shot inversion technique to extract the selected concept(Middle). Finally, these extracted concepts are combined to craft the  
37 story of a 'traveling robot' (Bottom).

38 Recent breakthroughs in text-to-image generation have greatly expanded  
39 the possibilities for creating visual content. Despite this progress, these models  
40 often struggle to consistently extract not only visual characters but also  
41 general concepts, such as style, pose, and color, from pre-trained models.  
42 Common approaches typically require multiple pre-existing images of the  
43 target subject or fine-tuning with large datasets of target class for one-shot  
44 inversion. Our research presents an innovative interactive method for text-  
45 driven diffusion models that allows for the extraction and preservation of  
46 custom concepts, which we refer to as "concept assets." This mechanism  
47 is an augmentation of the textual inversion process, incorporating inter-  
48 activity that empowers users to edit, define and preserve distinct concepts

49 in generated images for future use, all within a one-shot framework. With  
50 this approach, a concept can be captured within a pseudo-token, making it  
51 readily available for future use in various applications like character design,  
52 story telling and personalized image generation. We quantitatively evaluate  
53 our method against existing baselines using multiple metrics and validate  
54 our results with a user study. We also provide examples of how our method  
55 can be applied in practical scenarios.

56 CCS Concepts: • Computing methodologies → Image processing.

57 Additional Key Words and Phrases: Diffusion customization, concept extraction,  
58 generative AI

## 59 ACM Reference Format:

60 Anonymous Author(s). 2024. Concept Assets: Naming Custom Concepts  
61 in Diffusion Network Outputs. *ACM Trans. Graph.* 1, 1 (February 2024),  
62 12 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

63 Permission to make digital or hard copies of all or part of this work for personal or  
64 classroom use is granted without fee provided that copies are not made or distributed  
65 for profit or commercial advantage and that copies bear this notice and the full citation  
66 on the first page. Copyrights for components of this work owned by others than the  
67 author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or  
68 republish, to post on servers or to redistribute to lists, requires prior specific permission  
69 and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

70 © 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
71 ACM 0730-0301/2024/2-ART  
72 <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 73 1 INTRODUCTION

74 State of art text to image models exhibit remarkable ability to generate  
75 convincing images from natural language descriptions; still

there remains an inherent gap between the two domains, which complicates practical usage. Because of the structure and finite quantity of available training data, there is always a degree of ambiguity between a verbal description and appearance, making it difficult to request a particular output in most cases. For example, the phrase *sports car* may correspond to many different images in which the referenced objects are quite diverse in appearance. When generating images for a particular purpose however, the user may wish to retain a particular appearance between iterations to change the context in which a particular object appears, or fine-tune its properties by prompt engineering. But with the exception of very particular terms strongly associated with an appearance in the training data, there is no guarantee that another invocation of the network will yield the same appearance.

A variety of personalization techniques allows the user to associate a special token with the particular appearance of a concept to ensure consistency, but to achieve this they require either multiple pre-existing examples or involved image manipulation such as selection. This limits their utility to concepts of which there are either sufficient pre-existing real world examples, or which are sufficiently concrete to be easily selected.

In this paper, we introduce a new method for personalization which iteratively incorporates the user’s feedback to specify a concept that may then be named, saved, and reused in further generation. We extract these concepts from a root image which is either itself generated or provided by the user. This image is automatically modified and the user is then prompted to accept or reject these modifications in order to identify attributes of interest. These attributes can be arbitrary, and include identity (for instance, ‘robot’), appearance (e.g. ‘purple’), and action (such as ‘dancing’).

A key aspect of our approach is its ability to disentangle multiple attributes within an image through a combination of pseudoword parameters driven by textual embeddings and LoRA (Low-Rank Adaptation) parameter space. Thus disentangled, attributes from a single root image may then be used and edited independently in future synthesis. For example, we can create variations of a ‘dancing robot’ with different appearances but the same dance pose. Each resulting variant is then labeled with a descriptive caption, such as ‘a <dance1> <robotn>’, to highlight the altered attribute. This method of pairing not only aids in distinguishing between attributes but also acts as a self-regulating mechanism in our training process. Through this disentanglement, our approach allows for the efficient extraction of a broad spectrum of concepts in a single training session, each of which is self-consistent. These named attributes, which we dub *concept assets*, then represent concepts of consistent, specific and custom appearance which bridge the domain gap.

We then show that our method can effectively combine both visual and abstract concepts extracted directly from diffusion outputs. This results in output images that feature a natural semantic blend and novel layouts, demonstrating versatility and generalization capability of the extracted assets. Such properties make our method particularly suited for applications like comic creation, storytelling, and poster design.

A key aspect of our approach is that post-extraction, the concept assets are encapsulated entirely within the word embedding of the pseudotokens, and so our technique is not dependent on any

particular model architecture or pre-trained weights, and may be transparently combined with any number of techniques which e.g. modify attention weights or the inference process.

In summary, our primary contributions are as follows: (1) We establish a framework for the extraction of multiple general concepts, moving beyond merely generating visually consistent objects. (2) We introduce a lightweight, innovative solution that is non-intrusive to existing architectures and tailored to this expanded scope. It uniquely incorporates user feedback for customization. (3) We conduct comprehensive evaluations, both quantitative and qualitative, including a user study, to demonstrate the robustness and effectiveness of our method in handling a broader range of concepts.

## 2 RELATED WORK

**Text-to-Image Generation** Recently diffusion models have driven input-conditioned image generation has undergone a revolution in quality and flexibility. Previously, GAN based methods [6, 17, 32] set a foundation for the field, offering impressive diversity and quality in image generation, yet they often fell short in accurately translating complex textual descriptions into images. Diffusion-based models [8, 11, 13, 15, 26, 28–31] have since made significant progress overcoming this challenge. An increasing number of large-scale text-to-image models [7, 20, 21, 25] are being trained using billions of data sources. These foundation models [4] excel at synthesizing photo-realistic images that correspond to user prompts.

**Personalization** One of the most impressive aspects of foundation text-to-image models is their capacity for personalization. In this context personalization focuses on consistently generating a user-defined visual concept, bound to a text token, based on a handful of reference images. Two foundational approaches in this area are Textual Inversion (TI) [9] and DreamBooth (DB) [23].

Textual Inversion introduces a unique, trainable text token. This token is optimized to reconstruct the target concept using standard diffusion loss techniques, while the overall model weights remain static. DreamBooth, in contrast, utilizes an existing, but rare, text token. It then fine-tunes the entire model, adapting its weights to better reconstruct the target concept.

Custom Diffusion [16] represents a slightly different approach, fine-tuning only specific layers of the model, whereas LoRA [14, 24] confines its updates to low-rank adjustments. In our research, we employ a two-phase strategy. First, optimize only the textual embeddings of the target concepts. Second, we jointly train both the embeddings and the model weights simultaneously.

Personalization typically involves using target images that contain the same object, to ensure a diversity of outputs. However, extracting concepts directly from diffusion model outputs is more challenging, especially when the desired concept appears only sporadically. Recent developments in this area include one-shot textual inversion techniques like E4T [10], which leverages extensive pre-training across a broad spectrum of concepts within a given domain for quick inversion. Break-a-Scene [1] introduces an innovative approach, utilizing additional mask information and an attention control mechanism for the target visual concept, allowing for more

229 detailed control over the generated scenes. Like our approach Break-  
 230 a-Scene [1] uses a similar two-stage training approach.

231 Despite these advances, methods like E4T and Break-a-Scene face  
 232 limitations in generalizing to abstract and broad concepts due to  
 233 their heavy reliance on visual cues.

234 Another research direction focuses on using text prompts alone  
 235 to generate consistent concepts, as seen in studies like Elodin [18]  
 236 and The Chosen One [2]. While these methods achieve visual con-  
 237 sistency, their training methodology—expectation and clustering  
 238 aligned with text prompts—often results in outputs that lack de-  
 239 tail and appear overly smoothed. In addition, these methods can  
 240 only learn a single visually consistent character per text-prompt. In  
 241 contrast our work enables extracting a consistent character from a  
 242 single image, giving users much more flexibility.

### 244 3 METHOD

245 The objective of our method is to extract concepts from the diffusion  
 246 prior by starting from a single seed image, and then augmenting  
 247 this to a handful of additional images that facilitate learning dis-  
 248 entangled concepts. This supports objects and characters, but also  
 249 the identification of more abstract elements such as actions, colors,  
 250 styles, and camera views.

251 In contrast to previous work that relies on training sets of multiple  
 252 images of the same object or style combined with a fixed user-  
 253 provided text template, we start from a single root image  $I$  and  
 254 no fixed text template. This image can either be a natural image  
 255 provided by the user, or one generated by the diffusion model.

256 Our key insight is to view an image as a mix of different attributes,  
 257 such as ‘identity’ (like ‘robot’) and ‘action’ (such as ‘dancing’). By  
 258 applying the SDEdit technique with specific text prompts, we are  
 259 able to alter one attribute while keeping the others constant, creating  
 260 a variety of images related to the original.

261 We pair each modified image with a distinct caption that high-  
 262 lights the altered attribute. For instance, if we modify the ‘robot’  
 263 aspect while maintaining the ‘dancing’ action, the resulting image,  
 264 denoted as  $I_n$ , is then labeled with a caption such as ‘a <dance1>  
 265 <robot N>’. This augmentation process acts as a regularization mech-  
 266 anism which helps correctly disentangled visual attributes. Conse-  
 267 quently, a single training yields a variety of pseudo words that can  
 268 be mixed and matched at inference.

#### 270 3.1 Pretrained Text-guided Latent Diffusion Models:

271 Consider a text-guided diffusion model that initiates with a textual  
 272 embedding  $C$  and a random Gaussian noisy image, denoted as  $I_T$ .  
 273 Note that  $C = \mathcal{T}(P)$  is the embedding of text prompt  $P$  in natural  
 274 language projected by the text encoder  $\mathcal{T}$  and  $I_T \in \mathbb{R}^{H \times W \times C}$  is  
 275 characterized by Gaussian i.i.d pixel values.

276 The goal of the diffusion model is to progressively denoises the  
 277 image, resulting in a sequence  $I_T, I_{T-1}, \dots, I_0$  such that  $I_0$  corre-  
 278 sponds to the text prompt  $P$ . But to allow the diffusion model to  
 279 operate on a more compact representation, which reduces both time  
 280 complexity and memory usage, Latent Diffusion Models (LDMs) [22]  
 281 uses an encoder  $\mathcal{E}$  to map a given image  $I$  into a latent embedding  
 282  $z$ . Subsequently, a decoder  $\mathcal{D}$ , is employed to reconstruct the input  
 283 image from  $z$ , such that  $\mathcal{D}(\mathcal{E}(I)) \approx I$ . Therefore, we only need to

284 replace the image  $I$  with its latent embedding  $z$  in the following  
 285 algorithm.

286  $z_t$  is a sample with added standard Gaussian noise  $\epsilon$ , where the  
 287 noise is introduced according to a time-dependent schedule  $\alpha_t$ . The  
 288 relation is given by  $z_t = \sqrt{\alpha_t} z_0 + \sqrt{1 - \alpha_t} \epsilon$ , corresponding to a  
 289 forward diffusion process:

$$290 q(x_t | x_{t-1}) := \mathcal{N}(x_t; \sqrt{\alpha_t} x_{t-1}, (1 - \alpha_t)I) \quad (1) \\ 291$$

292 To learn the reverse process, the mean value  $\mu_\theta(x_t, t)$ , parame-  
 293 terized by  $\theta$ , is predicted:

$$294 p_\theta(x_{t-1} | x_t) := \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta) \quad (2) \\ 295$$

296 A denoiser module, parameterized by the network  $\epsilon^\theta$ , is utilized  
 297 for fitting the objective, which is equivalent to fitting the mean value  
 298 prediction in the reverse process:

$$299 \min_\theta \mathbb{E}_{z_0, \epsilon \sim \mathcal{N}(0, I), t \sim \text{Uniform}(1, T)} \|\epsilon - \epsilon_\theta(z_t, t, C)\|^2. \quad (3) \\ 300$$

301 During the inference phase, the model can employ either sto-  
 302 chastic sampling as in DDPM [12], which introduces noise at each  
 303 sampling step, or deterministic sampling as in DDIM [27], which  
 304 follows an ODE-like deterministic trajectory. Specifically, we utilize  
 305 the deterministic DDIM sampling approach to leverage the inverse  
 306 characteristics of the ODE path:

$$307 \frac{z_{t-1}}{\sqrt{\alpha_{t-1}}} = \frac{z_t}{\sqrt{\alpha_t}} + \left( \frac{\sqrt{1 - \alpha_{t-1}}}{\sqrt{\alpha_{t-1}}} - \frac{\sqrt{1 - \alpha_t}}{\sqrt{\alpha_t}} \right) \epsilon^\theta(z_t, t, C) \quad (4) \\ 308$$

#### 313 3.2 Text Guided SDEdit for Data Augmentation:

314 In our approach, we adapt the Stochastic Differential Equation (SDE)  
 315 editing technique, which introduces noise into an existing image  
 316 and then partially reverses this process to produce a new image.  
 317 Originally, this technique was used for editing real-world images,  
 318 using manual strokes for edits, and the subsequent denoising step  
 319 did not utilize textual prompts. However, we have modified the  
 320 method to better suit our objectives. Our adapted procedure in-  
 321 volves introducing noise into the image and then employing text-  
 322 guided, classifier-free guidance  $\epsilon^\theta(z_t, t, C)$  during the DDIM denoising  
 323 phase to achieve the desired outcome.

324 Initially, the input image is converted into a latent space rep-  
 325 resentation through an encoder, denoted as  $z = \mathcal{E}(I)$ . This latent  
 326 representation,  $z$ , is then subjected to the addition of independent  
 327 standard Gaussian noise, expressed as  $z_t = \sqrt{\alpha_t} z + \sqrt{1 - \alpha_t} \epsilon$ . The  
 328 denoising process commences not from the initial step  $T$  but from  
 329 an intermediate step  $t$ . Note that if noise is added too early in the  
 330 sequence, specifically at a timestep near  $T$ , this can lead to newly  
 331 generated images with no relation to the root image  $I$ . If the noise is  
 332 introduced at a timestep close to 0, which is too late in the sequence.  
 333 Such a timestep can result in insufficient variations in the content of  
 334 the newly generated images. Therefore we choose to introduce noise  
 335 at timestep  $T = 0.86$  to balance variety and identity preservation.

#### 336 3.3 Interactive User Selection for Dual images and 337 Captions:

338 Based on their goals users can use text-guided SDEdit to generate a  
 339 wide range of variations from the initial seed image; then select their  
 340 preferred variations to use when learning disentangled concepts.

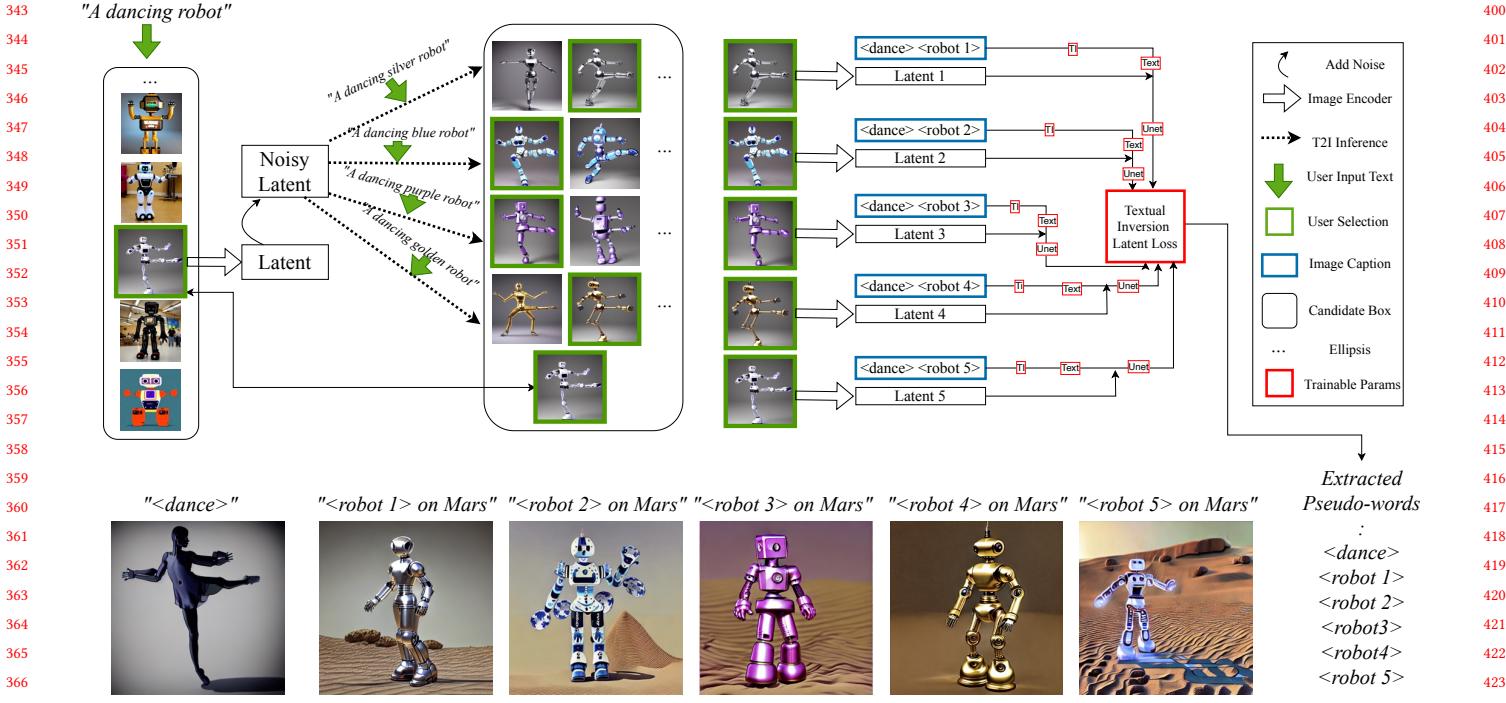


Fig. 2. Overview of the Method: Our approach is divided into three primary phases. (1) From the user’s text prompt, an initial image is chosen as the starting point of the process from among various outputs, ensuring it closely aligns with the user’s requirements. (2) We then employ text-guided image perturbation methods to create a range of potential images. The user selects a group of images that exhibit a common characteristic, such as a similar ‘dance’ pose. (3) Each chosen image is then labeled with a simple caption that highlights its distinct, intrinsic attribute. Following this, we use textual inversion and LoRA training to simultaneously extract multiple pseudo-concepts (or concepts).

This selection process can be iteratively refined by user feedback, maintaining image quality without introducing degradation or artifacts. This interactive step is crucial for users to tailor the semantic details they desire from the pre-trained model, as shown in Figure 2.

For instance, if the original image  $I$  is categorized as a ‘dancing robot’, it represents a combination of the action ‘dance’ and the identity ‘robot’. To isolate the ‘dance’ action and identify the ‘robot’, we can use SDEdit with prompts like ‘a dancing silver robot’, ‘a dancing purple robot’, ‘a dancing single-eyed robot’, etc. Each prompt generates a batch of images from which users can choose. The chosen image can either be added to the training dataset or returned to the SDEdit process for further refinement. For our training set, we select eight images based on empirical findings. The effects of using fewer or more images in the training set are shown in Section 4.3. Each training image is labeled with a unique caption, such as ‘<dance 1> <robot 1>’, through to ‘<dance 1> <robot N>’. While we vary the robot identities, we keep the dance action constant, resulting in a single action pseudoword ‘<dance 1>’ to represent the common dancing pose, and multiple identity pseudowords to distinguish the different robot appearances.

Figure 3 showcases various use cases, demonstrating how our method allows for the creation of desired image variations based on textual guidance provided by users. Utilizing our approach, these

varied images can be converted into pseudo-words which can be invoked for subsequent text-to-image inference tasks. Thus, the variation generation stage is a valuable tool for users to directly design and refine concepts. In Figure 8, we display the results of inference using these extracted and combined assets.

### 3.4 Multiple Words Textual Inversion:

Our methodology is anchored in the application of the pre-trained Stable Diffusion [21] model, adapted for our specific purposes. Consistent with the textual inversion framework, we extend the tokenizer’s lexicon by incorporating each pseudo word  $< S_i >$ , represented as a novel word embedding  $v_i$  within the text encoder. Our approach diverges from prior formulations by learning a set of parameters  $V = \{v_i\}_{i=1}^n$  concurrently. Moreover, we enhance the consistency and expressive capacity of the pseudo words by updating the residual weights of the model  $\Delta W$ . This update is effected through a low-rank adaptation (LoRA) applied to both the text-encoder and the U-Net component of the architecture.

The optimization objective is formulated thus:

$$\Theta = \arg \min_{\Theta=\{V, \Delta W\}} \mathbb{E}_{z \sim \mathcal{E}(x), y \sim \mathcal{N}(0,1), t} \left[ \|\epsilon - \epsilon_\Theta(z_t^k, t, c_\Theta^k(y))\|_2^2 \right], \quad (5)$$

$$k = 1, 2, \dots, n \quad (6)$$

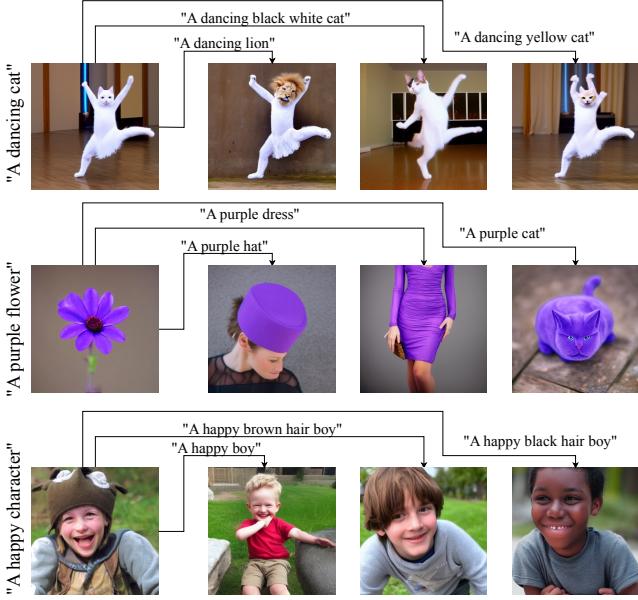


Fig. 3. User-Driven Conceptual Editing Phase: This stage begins with the images selected by the user. Users provide text prompts for various alterations. As this editing process relies solely on text, it enables the creation of more diverse concepts from the original image with ease. The provided examples demonstrate that the edits can range from subtle (altering only the head style in the ‘cat’ example) to imaginative (creating a ‘Ballet Lion’). The results effectively retain key characteristics from the original image, such as the consistent purple color seen in the second row and consistent happy mood seen in the third row. Additionally, the edits blend seamlessly into the natural image distribution, as seen in the perturbed ‘boy face’ images, which appear natural and undistorted.

In this expression,  $y^k$  corresponds to the caption of the  $k$ -th image, and  $z_t^k$  is the perturbed latent variable of the  $k$ -th image, injected with noise according to the timestep  $t$ . It is important to emphasize that the pairing of images and captions is fixed, as delineated in the previous section. Each selected image  $I_k$  is paired exclusively with a single caption ' $< A_1 >< B_k >$ ', where  $A$  and  $B$  represent two distinct characteristics that inherently differentiate the various elements within the image’s composition.

Note that in the second stage of our approach, we fine-tune not only the textual embedding but also text-encoder and the U-Net component of the stable-diffusion model. During inference, especially when incorporating multiple custom pseudo words into a single text prompt, it’s necessary to integrate more than just the textual embeddings of these pseudo words. Specifically, we also need to combine their LoRA matrices. When adding the delta parameterized by a LoRA matrix, apply a scaling factor of 0.8. In cases where multiple LoRA meetings (from multiple concepts) are invoked in a single training, we scale the first by 0.8 and subsequent ones by 0.5 (order determined by pseudo-word order in text prompt). For a more comprehensive understanding of the LoRA fine-tuning method, readers are encouraged to consult the referenced citation[14][24]."

## 4 EXPERIMENTS

In Section 4.1, we evaluate our method against various baselines through qualitative and quantitative comparisons. Section 4.2 describes the user study along with its findings. In Section 4.3, we present the results of our ablation study. Lastly, Section 4.4 showcases several practical applications of our method.

### 4.1 Quantitative and Qualitative Comparison

For a quantitative evaluation of our method, we compute the text-image and image-image feature similarity to numerically measure the prompt consistency and identity consistency, following [3]. First, we automatically generate 10 candidates using ChatGPT for the image character and the text description, similar to [3]. Then, we compute the text-image feature similarity as a prompt consistency metric by using the CLIP [19] text and image encoder to calculate the feature distance between the corresponding prompt and the generated image. For the image-image similarity, CLIP image feature embedding similarities are measured for the generated images across the same prompt concept as the identity consistency.

The left plot in Figure 5 visualizes the result of the quantitative evaluation of baselines and our models in a 2D plot where the x-axis is prompt consistency and y-axis is identity consistency. Our model achieves the best prompt consistency with relatively good identity consistency compared to the baselines, LoRA DB [24], Textual Inversion (TI) [9], The Chosen One [2] and Break-A-Scene [1]. We have observed that baselines with a relatively high identity consistency score, such as LORA DB, actually produce candidates with limited diversity. In contrast, our method maintains a sweet spot between identity consistency and generative diversity.

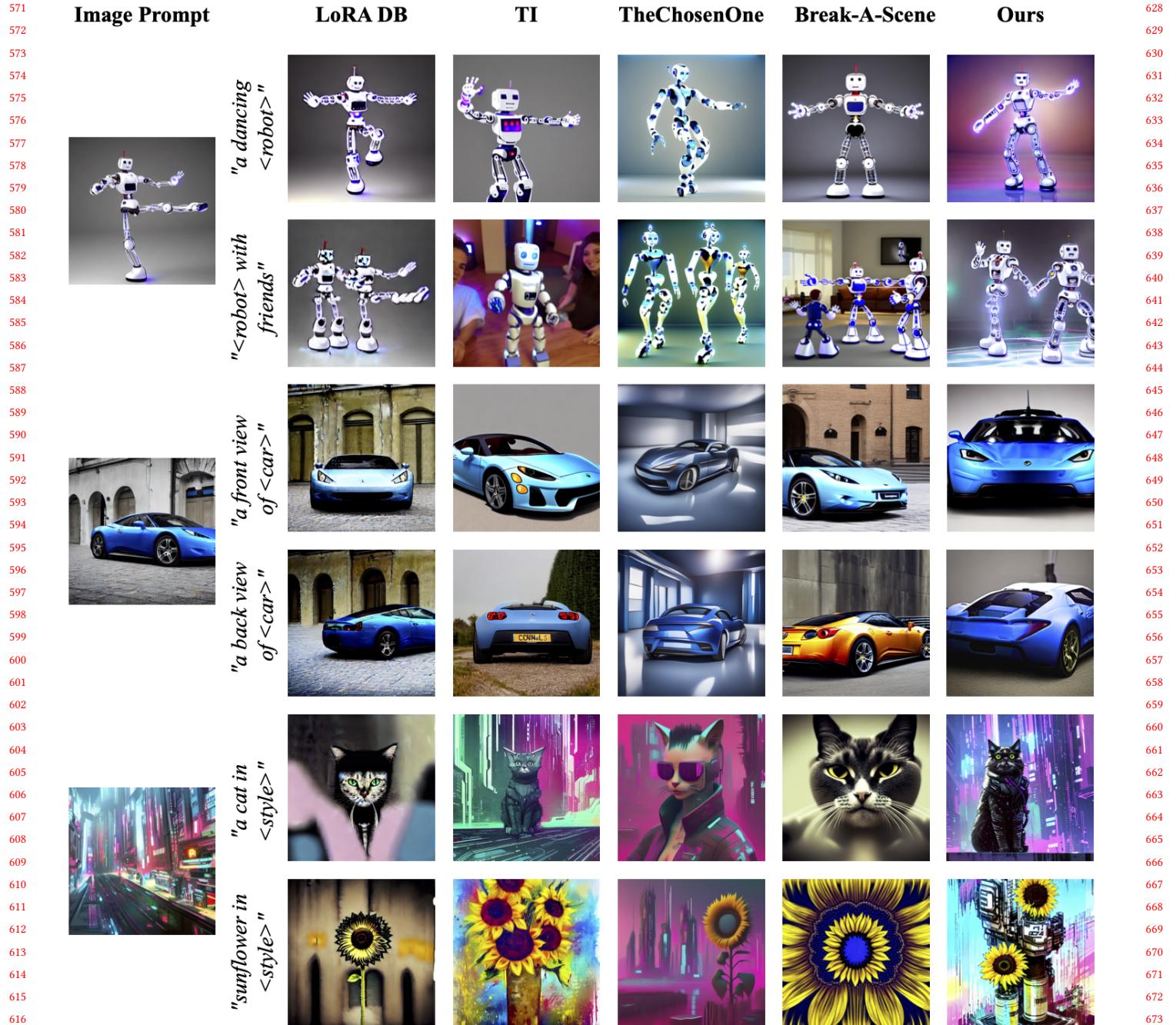
A visual comparison with the baselines is shown in Figure 4. Our model is able to generate results that are consistent to the text prompt while maintaining the original object identity (<robot> and <car>). Furthermore, our model can accurately extract abstract concepts like <style> in the last two rows of Figure 4.

### 4.2 User Study

To thoroughly evaluate our model based on a subjective assessment by human participants, we conducted a user study to gauge the consistency rate of the generated samples. Users were asked to rank each candidate from 1 to 5 for all five methods by independently considering the prompt consistency and identity consistency. We used the same platform, Amazon Mechanical Turk (AMT), to distribute the studies to participants. After averaging the scores across all categories, we obtained the final score, which is presented as the user study score, shown in the right plot in Figure 5.

### 4.3 Ablation Study

We performed an ablation study to determine the optimal number of variation inputs for training identity concept inversion. As shown in Figure 6, the results, indicate that training with a single image often leads to overfitting, resulting in the generation of images with limited diversity and compromised image quality during inference. Conversely, training with an excessive number of pseudo words (up to 16) can cause the model to overfit, which negatively impacts the diversity and fidelity of the background content. To ensure a



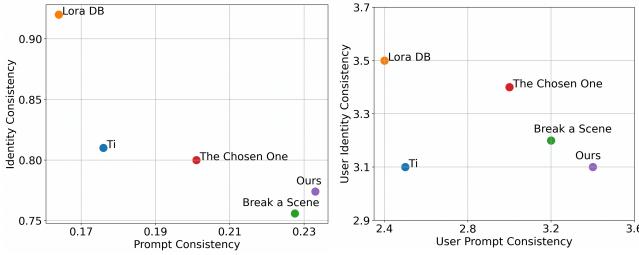


Fig. 5. Quantitative result (left) and User study (right) comparing with various baselines. We show that our method can get a consistent favorable choice under both feature similarity and user choice under random choice prompts and characters.

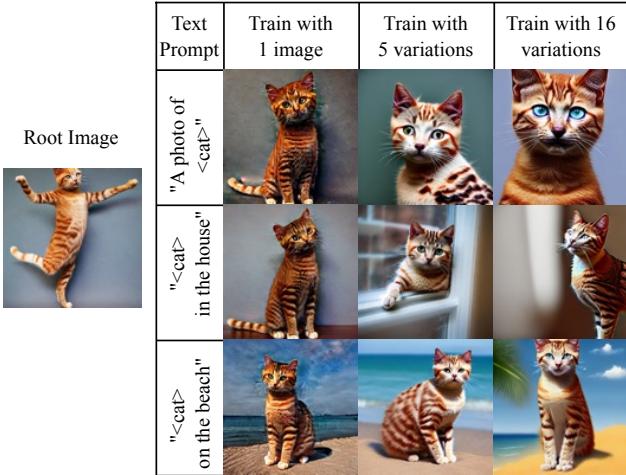


Fig. 6. Ablations on number of training variations

balanced comparison, we varied the number of textual inversion training steps in conjunction with LoRA training steps (500+500, 1000+1000, 1500+1500) and selected the best visual outcomes for the figure. Based on our observations, employing 5-8 variations and 1000+1000 steps within our framework yielded the most effective results in terms of concept consistency, prompt alignment and overall image quality.

#### 4.4 Applications

**4.4.1 General Concept Extraction.** Figure 8 presents various disentangled concepts derived from the examples in Figure 3. A key observation is that abstract concepts, although not visually explicit like ‘Ballet’, ‘Color’, and ‘Happy’, can be successfully extracted and applied to a variety of identity concepts. For additional examples demonstrating how our extracted concepts can be generalized and combined with various textual guidance during inference please see the Appendix.

**4.4.2 Multiple Pseudowords Combination.** Figure 1 illustrates the integration of multiple extracted assets from different experiments, showcasing the versatility of our method in applications such as

storytelling, comic book creation, poster design, and more. However, we find that combining an excessive number of concepts (more than four) in a single image is challenging. Additionally, we note that concepts related to style and color are more readily combined compared to visual identities. These limitations are further discussed in Section 5.

## 5 LIMITATIONS AND CONCLUSIONS



Fig. 7. Limitations examples.

Because our method does not modify the underlying diffusion model in any way, it exhibits the same issues those models do. This includes problems like catastrophic neglect and incorrect attribute binding, described in prior literature [5]. This manifests as e.g. inability to generate two concepts from the same class during inference, particularly when these concepts are textually similar, as illustrated in Figure 7. However, the same prior work describes solutions to this problem and may be combined with our method, again because we do not modify the underlying diffusion model and so retain compatibility with similar enhancements.

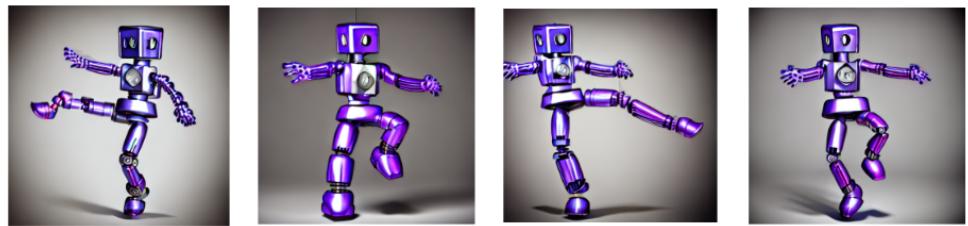
In conclusion, our method presents an innovative interactive method for text-driven diffusion models, significantly enhancing the textual inversion process to capture and retain unique ‘custom assets.’ This approach is user-friendly, facilitating effortless editing and application of these concepts in diverse areas such as character design and personalized image generation, all within a streamlined one-shot framework. Our method surpasses the limitations of merely achieving visual consistency, introducing a lightweight and adaptable solution that effectively incorporates user input without altering the underlying model architecture. Extensive evaluations, including comparisons with existing baselines and a user study, have validated the efficacy and adaptability of our method in managing a broad spectrum of concepts.

## 799 REFERENCES

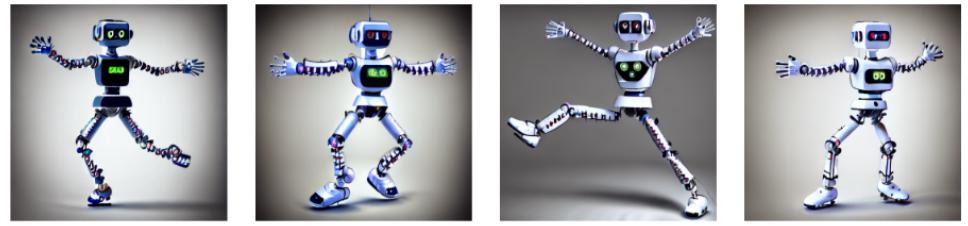
- 800 [1] Omri Avrahami, Kfir Aberman, Ohad Fried, Daniel Cohen-Or, and Dani Lischinski.  
801 2023. Break-A-Scene: Extracting Multiple Concepts from a Single Image. *arXiv*  
802 preprint arXiv:2305.16311 (2023).
- 803 [2] Omri Avrahami, Amir Hertz, Yael Vinker, Moab Arar, Shlomi Fruchter, Ohad  
804 Fried, Daniel Cohen-Or, and Dani Lischinski. 2023. The Chosen One: Consistent  
805 Characters in Text-to-Image Diffusion Models. *arXiv preprint arXiv:2311.10093*  
806 (2023).
- 807 [3] Omri Avrahami, Amir Hertz, Yael Vinker, Moab Arar, Shlomi Fruchter, Ohad  
808 Fried, Daniel Cohen-Or, and Dani Lischinski. 2023. The Chosen One: Consistent  
809 Characters in Text-to-Image Diffusion Models. *arXiv preprint arXiv:2311.10093*  
810 (2023).
- 811 [4] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora,  
812 Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma  
813 Brunsell, et al. 2021. On the opportunities and risks of foundation models. *arXiv*  
814 preprint arXiv:2108.07258 (2021).
- 815 [5] Hila Chefer, Yuval Alaluf, Yael Vinker, Lior Wolf, and Daniel Cohen-Or. 2023.  
816 Attend-and-Excite: Attention-Based Semantic Guidance for Text-to-Image Dif-  
817 fusion Models. *ACM Trans. Graph.* 42, 4, Article 148 (jul 2023), 10 pages.  
818 https://doi.org/10.1145/3592116
- 819 [6] Katherine Crowson, Stella Biderman, Daniel Kornis, Dashiell Stander, Eric Hall-  
820 han, Louis Castricato, and Edward Raff. 2022. Vqgan-clip: Open domain im-  
821 age generation and editing with natural language guidance. *arXiv preprint arXiv:2204.08583*  
822 (2022).
- 823 [7] Xiaoliang Dai, Ji Hou, Chih-Yao Ma, Sam Tsai, Jialiang Wang, Rui Wang, Peizhao  
824 Zhang, Simon Vandenhende, Xiaofang Wang, Abhimanyu Dubey, et al. 2023. Emu:  
825 Enhancing image generation models using photogenic needles in a haystack.  
826 *arXiv preprint arXiv:2309.15807* (2023).
- 827 [8] Prafulla Dhariwal and Alexander Nichol. 2021. Diffusion models beat gans on  
828 image synthesis. *Advances in Neural Information Processing Systems* 34 (2021),  
829 8780–8794.
- 830 [9] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal  
831 Chechik, and Daniel Cohen-Or. 2022. An image is worth one word: Personalizing  
832 text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618*  
833 (2022).
- 834 [10] Rinon Gal, Moab Arar, Yuval Atzmon, Amit H Bermano, Gal Chechik, and Daniel  
835 Cohen-Or. 2023. Encoder-based domain tuning for fast personalization of text-to-  
836 image models. *ACM Transactions on Graphics (TOG)* 42, 4 (2023), 1–13.
- 837 [11] Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic  
838 models. *Advances in Neural Information Processing Systems* 33 (2020), 6840–6851.
- 839 [12] Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic  
840 models. *Advances in neural information processing systems* 33 (2020), 6840–6851.
- 841 [13] Emiel Hoogeboom, Jonathan Heek, and Tim Salimans. 2023. simple diffusion:  
842 End-to-end diffusion for high resolution images. *arXiv preprint arXiv:2301.11093*  
843 (2023).
- 844 [14] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean  
845 Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large  
846 language models. *arXiv preprint arXiv:2106.09685* (2021).
- 847 [15] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. [n. d.]. Elucidating  
848 the Design Space of Diffusion-Based Generative Models. In *Advances in Neural*  
849 *Information Processing Systems*.
- 850 [16] Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu.  
851 2023. Multi-concept customization of text-to-image diffusion. In *Proceedings of*  
852 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1931–1941.
- 853 [17] Xingchao Liu, Chengyue Gong, Lemeng Wu, Shujian Zhang, Hao Su, and Qiang  
854 Liu. 2021. Fusedream: Training-free text-to-image generation with improved  
855 clip+ gan space optimization. *arXiv preprint arXiv:2112.01573* (2021).
- 856 [18] Rodrigo Mello, Filipe Calegario, and Geber Ramalho. 2023. ELODIN: Naming  
857 Concepts in Embedding Spaces. *arXiv preprint arXiv:2303.04001* (2023).
- 858 [19] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sand-  
859 hini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al.  
860 2021. Learning transferable visual models from natural language supervision.  
861 *arXiv preprint arXiv:2103.00020* (2021).
- 862 [20] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec  
863 Radford, Mark Chen, and Ilya Sutskever. 2021. Zero-shot text-to-image generation.  
864 In *International Conference on Machine Learning*. PMLR, 8821–8831.
- 865 [21] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn  
866 Ommer. 2021. High-Resolution Image Synthesis with Latent Diffusion Models.  
867 *arXiv:2112.10752 [cs.CV]*
- 868 [22] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn  
869 Ommer. 2022. High-resolution image synthesis with latent diffusion models. In  
870 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*.  
871 10684–10695.
- 872 [23] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and  
873 Kfir Aberman. 2023. Dreambooth: Fine tuning text-to-image diffusion models for  
874 subject-driven generation. In *Proceedings of the IEEE/CVF Conference on Computer*  
875 *Vision and Pattern Recognition*. 22500–22510.
- 876 [24] Simo Ryu. 2023. Low-rank adaptation for fast text-to-image diffusion fine-tuning.
- 877 [25] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily  
878 Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi,  
879 Rapha Gontijo Lopes, Tim Salimans, Tim Salimans, Jonathan Ho, David J Fleet,  
880 and Mohammad Norouzi. 2022. Photorealistic Text-to-Image Diffusion Models  
881 with Deep Language Understanding. *arXiv preprint arXiv:2205.11487* (2022).
- 882 [26] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli.  
883 2015. Deep unsupervised learning using nonequilibrium thermodynamics. In  
884 *International Conference on Machine Learning*. PMLR, 2256–2265.
- 885 [27] Jiaming Song, Chenlin Meng, and Stefano Ermon. 2020. Denoising diffusion  
886 implicit models. *arXiv preprint arXiv:2010.02502* (2020).
- 887 [28] Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. 2021. Maximum like-  
888 lihood training of score-based diffusion models. *Advances in Neural Information  
889 Processing Systems* 34 (2021), 1415–1428.
- 890 [29] Yang Song and Stefano Ermon. 2019. Generative modeling by estimating gradients  
891 of the data distribution. *Advances in Neural Information Processing Systems* 32  
892 (2019).
- 893 [30] Yang Song and Stefano Ermon. 2020. Improved techniques for training score-  
894 based generative models. *Advances in neural information processing systems* 33  
895 (2020), 12438–12448.
- 896 [31] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano  
897 Ermon, and Ben Poole. [n. d.]. Score-Based Generative Modeling through Stochastic  
898 Differential Equations. In *International Conference on Learning Representations*.
- 899 [32] Yufan Zhou, Ruiyi Zhang, Changyou Chen, Chunyuan Li, Chris Tensmeyer, Tong  
900 Yu, Jiuxiang Gu, Jinhui Xu, and Tong Sun. 2021. Lafite: Towards language-free  
901 training for text-to-image generation. *arXiv preprint arXiv:2111.13792* (2021).
- 902 [33] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano  
903 Ermon, and Ben Poole. [n. d.]. Score-Based Generative Modeling through Stochastic  
904 Differential Equations. In *International Conference on Learning Representations*.
- 905 [34] Yufan Zhou, Ruiyi Zhang, Changyou Chen, Chunyuan Li, Chris Tensmeyer, Tong  
906 Yu, Jiuxiang Gu, Jinhui Xu, and Tong Sun. 2021. Lafite: Towards language-free  
907 training for text-to-image generation. *arXiv preprint arXiv:2111.13792* (2021).
- 908 [35] Yufan Zhou, Ruiyi Zhang, Changyou Chen, Chunyuan Li, Chris Tensmeyer, Tong  
909 Yu, Jiuxiang Gu, Jinhui Xu, and Tong Sun. 2021. Lafite: Towards language-free  
910 training for text-to-image generation. *arXiv preprint arXiv:2111.13792* (2021).
- 911 [36] Yufan Zhou, Ruiyi Zhang, Changyou Chen, Chunyuan Li, Chris Tensmeyer, Tong  
912 Yu, Jiuxiang Gu, Jinhui Xu, and Tong Sun. 2021. Lafite: Towards language-free  
913 training for text-to-image generation. *arXiv preprint arXiv:2111.13792* (2021).

913 <Ballet> 971  
 914 972  
 915 973  
 916 974  
 917 975  
 918 976  
 919 977  
 920 978  
 921 979  
 922  
 923 <color> 980  
 924 981  
 925 982  
 926 983  
 927 984  
 928 985  
 929 986  
 930 987  
 931 988  
 932  
 933 <happy> 990  
 934 991  
 935 992  
 936 993  
 937 994  
 938 995  
 939 996  
 940 997  
 941 998  
 942  
 943  
 944  
 945  
 946  
 947  
 948  
 949  
 950  
 951  
 952  
 953  
 954  
 955  
 956  
 957  
 958  
 959  
 960  
 961  
 962  
 963  
 964  
 965  
 966  
 967  
 968  
 969

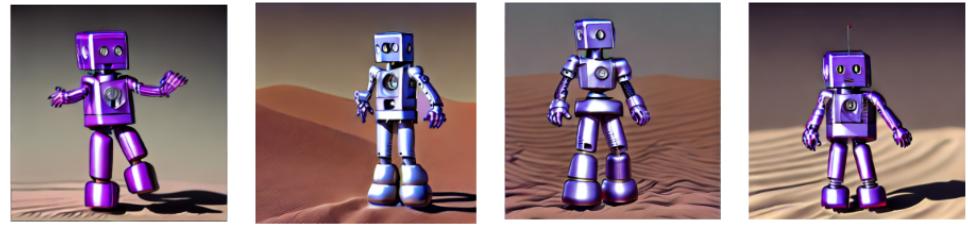
A DANCE (ROBOT 1)



A DANCE (ROBOT 2)



&lt;ROBOT 1&gt; ON MARS



&lt;ROBOT 2&gt; ON MARS

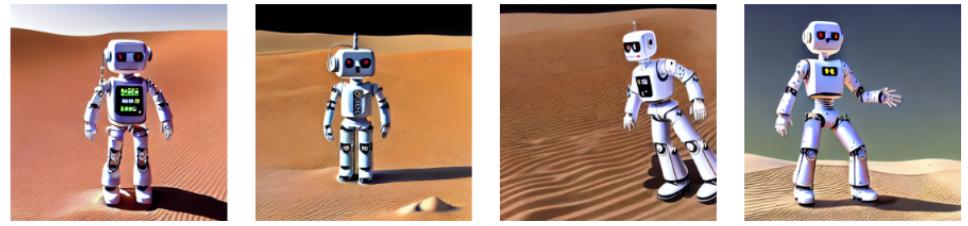


Fig. 9. More visualization examples, we show our method can generalize in different cases.

1141

1142

1143

1144

1145

1146

1147

**HAPPY <BOY 1>**

1148

1149

1150

1151

1152

1153

1154

1155

1156

**HAPPY <BOY 2>**

1157

1158

1159

1160

1161

1162

1163

1164

**ASTRONAUT <BOY 1>**

1165

1166

1167

1168

1169

1170

1171

1172

**ASTRONAUT <BOY 2>**

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1242

1243

1244

1245

1246

1247

1248

1249

1250

1251

1252

1253

1254

Fig. 10. More visualization examples, we show our method can generalize in different cases.

1255

1256

1257

1258

1259

1260

1261

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295

1296

1297

1298

1299

1300

1301

1302

1303

1304

1305

1306

1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

1350

1351

1352

1353

1354

1355

1356

1357

1358

1359

1360

1361

1362

1363

1364

1365

1366

1367

1368

LOTUS IN &lt;STYLE 1&gt;



LOTUS IN &lt;STYLE 2&gt;



WOMEN IN &lt;STYLE 1&gt;



WOMEN IN &lt;STYLE 2&gt;

