
Can Pre-Trained Text-to-Image Models Generate Visual Goals for Reinforcement Learning?

Anonymous Author(s)

Affiliation

Address

email

Abstract

Pre-trained text-to-image generative models can produce diverse, semantically rich, and realistic images from natural language descriptions. Compared with language, images usually convey information with more details and less ambiguity. In this study, we propose Learning from the Void (LfVoid), a method that leverages the power of pre-trained text-to-image models and advanced image editing techniques to guide robot learning. Given natural language instructions, LfVoid can edit the original observations to obtain goal images, such as “wiping” a stain off a table. Subsequently, LfVoid trains an ensembled goal discriminator on the generated image to provide reward signals for a reinforcement learning agent, guiding it to achieve the goal. The ability of LfVoid to learn with zero in-domain training on expert demonstrations or true goal observations (the void) is attributed to the utilization of knowledge from web-scale generative models. We evaluate LfVoid across three simulated tasks and validate its feasibility in the corresponding real-world scenarios. In addition, we offer insights into the key considerations for the effective integration of visual generative models into robot learning workflows. We posit that our work represents an initial step towards the broader application of pre-trained visual generative models in the robotics field. Our project page: LfVoid.github.io.

1 Introduction

What is the simplest way to provide guidance or goals to a robot? This question is answered multiple times: a set of expert demonstrations [1, 2, 3, 4], goal images [5, 6, 7], or natural language instructions [8, 9]. However, these answers either require a laborious and sometimes prohibitive effort to collect data, or contain ambiguity across modalities. The desire to alleviate the challenges begs the question: can we generate goal images from natural language instructions directly, without physically achieving the goals for robots?

Large text-to-image generative models [10, 11, 12, 13] have achieved exciting breakthroughs, demonstrating an unprecedented ability to generate plausible images that are semantically aligned with given text prompts. Trained on extensive datasets, these models are thought to possess a basic understanding of the world [14]. Therefore, we aim to harness the power of these large generative models to provide unambiguous visual goals for robotic tasks without any in-domain training.

In the past, a common approach to leverage off-the-shelf large generative models in robot learning involves using these models for data augmentation on expert datasets to improve policy generalization [15, 16, 17]. Despite the success, these methods still rely on human demonstrations, not generative models, as the main source of guidance.

35 DALL-E-Bot [18], an early attempt to utilize web-scale generative models in a zero-shot manner for
36 guiding manipulation tasks, utilizes DALL-E 2[10] to generate goal images for object rearrangement
37 tasks. However, the generated images are often too diverse and largely disagree with real-world
38 scenarios, requiring segmentation masks for object matching and a rule-based transformation planner
39 to close the visual discrepancy.

40 The recent uprise of language-based image editing methods [19, 20, 21] provides a new paradigm:
41 editing the generated images to align with language descriptions, such as “wiping” a stain off a table,
42 while keeping the visual appearance of the other objects and the background roughly unchanged.
43 However, the edited images are usually examined by visual appearance rather than embodied tasks.

44 In this work, we introduce Learning from the Void (LfVoid), a method that uses image editing
45 techniques on pre-trained diffusion models to generate visual goals for reinforcement learning. LfVoid
46 contains a unique image editing module capable of performing appearance-based and structure-based
47 editing to generate goal images according to language instructions. With these goal images, LfVoid
48 improves upon example-based RL methods [5, 6] to solve several robot control tasks without the
49 need for any reward function or demonstrations.

50 We evaluate LfVoid on three simulated environment tasks and the corresponding real-world scenarios.
51 Empirical results show that LfVoid can generate goal images with higher fidelity to both the editing
52 instructions and the source images when compared to existing image editing techniques, leading to
53 better downstream control performance in comparison to language-guided and other image-editing-
54 based methods. We also validate the feasibility of LfVoid in real-world settings: LfVoid can retrieve
55 meaningful reward signals from generated images comparable to those from true goal images. Based
56 on these observations, we discuss key considerations for the current generative model research when
57 aiming for real-world robotic applications.

58 Our contributions are threefold: First, we propose an effective approach to leverage the knowledge
59 encapsulated in large pre-trained generative models and apply it in a zero-shot manner to guide robot
60 learning tasks. Second, we provide empirical evidence that suggests LfVoid-edited images provide
61 guidance more effectively than other methods including the direct usage of text prompts. Lastly,
62 we identify the existing gap between the capabilities of current text-conditioned image generation
63 models and the demands in robotic applications, thus shedding light on a potential direction for future
64 research in this domain.

65 2 Related work

66 **Image editing for diffusion models.** Recent work in text-to-image diffusion models has
67 demonstrated a strong ability to generate images that are semantically aligned with given text
68 prompts[22, 23, 10, 11, 12, 24, 25]. Based on these generative models, Prompt-to-Prompt [19]
69 proposes to perform text-conditioned editing through injection of the cross-attention maps. Imagic
70 [20] proposes to perform editing by using interpolation between source and target embeddings to
71 generate images with edited effects. InstructPix2Pix [26] demonstrates another approach for training
72 a diffusion model on the paired source and target images and the corresponding editing instructions.
73 Apart from editing techniques, several methods focusing on controlled editing have been proposed.
74 Textual-Inversion [27] and DreamBooth [28] both aim to learn a special token corresponding to a
75 user-specified object and at inference time can generate images containing that object with a diverse
76 background. Directed Diffusion [21] focuses on object placement generation, which can control the
77 location of a specified object to reside in a given area.

78 **Example-based visual RL.** Standard reinforcement learning methods require a predefined reward
79 function to guide the agent towards desired goal states, yet these reward functions may not always
80 capture the essence of the task at hand or be readily available. Example-based RL methods aim
81 to utilize the observations of the goal states to guide the learning process. VICE, proposed by
82 Fu et al. [5], establishes a general framework for learning from goal observations only: it trains
83 a discriminator with the observations in the replay buffer as negative samples and goal images as
84 positive samples. The positive logits for new observations are used as a reward to guide the agent
85 toward the goal observation. Building upon this, Singh et al. [6] introduce VICE-RAQ, integrating
86 the VICE algorithm with active querying methods through periodically asking a human to provide
87 labels for ambiguous observations. A label smoothing technique is also proposed to enhance reward

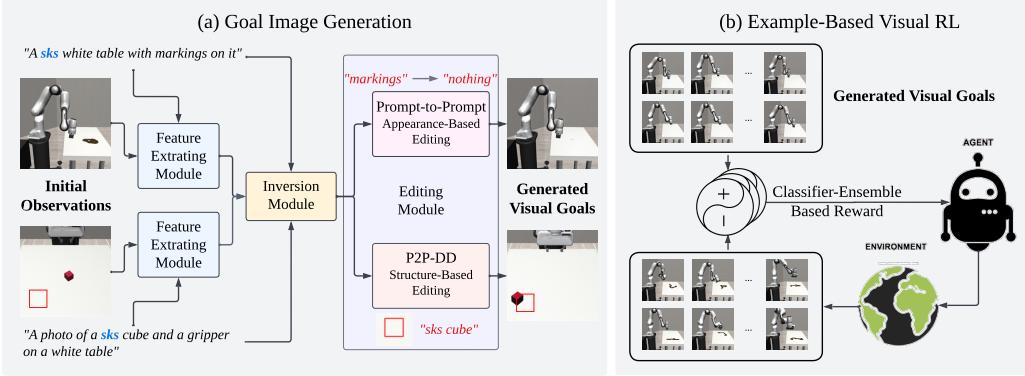


Figure 1: **An overview of LfVoid.** LfVoid consists of two parts: (a) Goal image generation, where we apply image editing on the initial observations according to different editing instructions to obtain a visual goal dataset; (b) Example-Based Visual RL, where we perform reinforcement learning on the generated dataset to achieve the desired goal image in various environments.

shaping. Other works [29, 30] are proposed to improve the reward shaping of VICE as well. In a parallel direction, Eysenbach et al. [7] develop Recursive Classification of Examples (RCE), a method that learns the Q function of actor-critic methods directly from the goal observations. RCE derives the Bellman equation for discriminator-based rewards and makes temporal difference updates on them to get an estimation of the discriminator-deduced Q function directly.

Large generative models for robot control. Current application of large generative models in robotic tasks mainly focuses on using Large Language Models for planning [8, 9, 31, 32, 33] or learning a language conditioned policy [34, 35], while the application of pre-trained image generative models for robotics tasks are limited. A common way of leveraging large-scale diffusion models for robot learning is by using diffusion models to perform data augmentation on training data, as suggested in CACTI [15], GenAug [16] and ROSIE [17]. Another line of work suggests using diffusion models to generate plans to solve robotics tasks [36, 37, 38]. More recently, UniPi[39] trains a text-to-video generation model on web-scale datasets and expert demonstrations to generate image sequences for planning and inverse modeling. The most related work is DALL-E-Bot [18], which uses the DALL-E 2[10] model to generate goal plans and performs object rearrangement according to these goals. However, unlike our work, DALL-E-Bot does not apply image editing techniques for goal generation and requires rule-based matching and predefined pick-and-place actions to bridge the visual gap between generated images and true observations.

3 Method

In this study, we aim to provide zero-shot visual goals to reinforcement learning agents for manipulation tasks using only text prompts. Our approach utilizes the information grounded in large-scale pre-trained visual generative models to create semantically meaningful visual goals. We first generate a synthetic goal image dataset from raw observations using image editing techniques, then employ example-based visual reinforcement learning that is optimized for our task with the generated dataset.

The structure of this section unfolds as follows: Section 3.1 elucidates the image editing techniques and adjustments applied to create goal images from text prompts and initial image observations, as depicted in Figure 1(a). Section 3.2 discusses the execution of example-based visual reinforcement learning using the generated goal images, as outlined in Figure 1(b).

3.1 Visual goal generation

In the first phase of our methodology, we edit and generate goal images from raw observations based on natural language guidance. Given a source prompt \mathcal{P} , a source image x_{src} , and an editing instruction, LfVoid synthesizes a target image x_{tgt} using a pre-trained Latent Diffusion Model (LDM) [12]. We consider two types of editing instructions: either a target prompt \mathcal{P}^* describing

121 the appearance changes in the target image, or a bounding-box region \mathcal{B} and a set of tokens \mathcal{I}
122 corresponding to the object to be relocated when structural changes are needed.

123 To provide sufficient visual guidance for downstream reinforcement learning tasks, the generated
124 goal images should highlight the visual changes while preserving the irrelevant scene as much as
125 possible, which is a challenging requirement even for state-of-the-art image editing techniques. To
126 this end, we integrate a number of different techniques in the visual goal generation pipeline of
127 LfVoid, consisting of a feature extracting module, an inversion module, and an editing module, as
128 outlined in Figure 1(a). Following conventional notations, we denote x_t as the synthesized LDM
129 image at time step $t \in \{T, T-1, \dots, 0\}$, where T is the total diffusion time steps. Specifically, x_T
130 denotes the Gaussian noise, and x_0 denotes the generated image. The techniques used in each module
131 and their purpose are described as follows.

132 3.1.1 Feature extracting module

133 To ensure high fidelity to the source image x_{src} , we learn a unique token sks that encapsulates the
134 visual features of objects within x_{src} as in DreamBooth [28]. This process involves optimizing the
135 diffusion model parameters along with the special token sks , using a set of images that contain the
136 target object. As a result, we are able to derive a specialized model that can accurately retain key
137 details of x_{src} , such as the color and texture of a cube to be positioned, or the shape of a Franka robot
138 arm for manipulation. As later demonstrated in our experiments, this module substantially boosts the
139 resemblance of our edited images to the corresponding source images.

140 3.1.2 Inversion module

141 After employing the special token sks to capture the essential features of the source scene, we
142 invert the provided source image x_{src} to a diffusion process. A simple inversion technique using
143 the Denoising Diffusion Implicit Model (DDIM) [23] sampling scheme calculates the samples
144 x_T, x_{T-1}, \dots, x_0 by reversing the ODE process. However, the inverted image x_0 obtained this way
145 often deviates from x_{src} due to cumulative errors. Null-text inversion [40] aims to mitigate this
146 discrepancy by fine-tuning the "null-text" embedding used in the classifier-free guidance [41] Φ_t
147 for each $t \in \{T, T-1, \dots, 1\}$ to control the generation process. Thus, with the initial noise sample
148 x_T obtained through DDIM inversion and the optimized "null-text" embedding $\Phi_t, \Phi_{t-1}, \dots, \Phi_1$, the
149 diffusion process is able to generate the image x'_{src} , an approximation of the source image x_{src} .
150 Based on the diffusion process obtained with the inversion module, LfVoid can perform precise and
151 detailed editing control on the source image, which will be discussed in Section 3.1.3.

152 3.1.3 Editing module

153 Depending on the specific requirements, we divide the editing tasks into two distinct categories:
154 appearance-based and structure-based image editing. Appearance-based editing involves maintaining
155 the structural layout of the image while altering certain visual aspects, such as cleaning the surface of
156 a table or lighting up an LED bulb. In contrast, structure-based editing involves changes in the layout
157 of the image, such as relocating objects within the image from one area to another.

158 **Appearance-based editing.** In tasks involving appearance-based editing, we employ the Prompt-to-
159 Prompt [19] editing control technique. When the Latent Diffusion Model (LDM) generates an image
160 x_0 conditioned on a text prompt \mathcal{P} , the information encapsulated in \mathcal{P} influences the diffusion process
161 via the cross-attention layers [12, 11, 19]. Prompt-to-Prompt reveals that the spatial configuration
162 of the image x_0 largely depends on the cross-attention maps M_t within these cross-attention layers,
163 especially during the initial diffusion steps. Consequently, it suggests replacing the attention maps
164 M_t^* of the target image diffusion process with the attention maps M_t from the source image diffusion
165 process for the first N time steps, beginning at time step $t = T$:

$$P2PEdit(M, M^*, t) = \begin{cases} M_t & \text{if } t > T - N \\ M_t^* & \text{otherwise} \end{cases} \quad (1)$$

166 This method ensures that the structure of the source image encapsulated in the attention maps M_t is
167 conserved in the target image in a more controlled manner, which is crucial in downstream RL tasks.

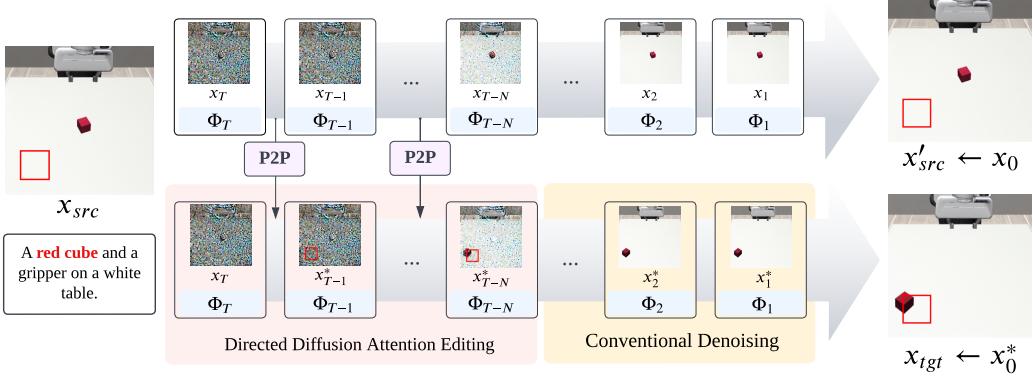


Figure 2: **Image editing with structural changes.** The details of P2P-DD technique used in the editing module of LfVoid to perform structural changes. LfVoid performs a combination of P2P attention map injection with Directed Diffusion attention editing on the diffusion process generating the target image x_{tgt} , based on the Gaussian noise x_T and the optimized null-text embeddings $\Phi_T, \Phi_{T-1}, \dots, \Phi_1$ obtained through the Inversion Module of LfVoid.

168 **Structure-based editing.** A notable limitation of Prompt-to-Prompt editing is the incapacity to
 169 spatially move existing objects across the image, i.e., implement structural changes. Therefore, we
 170 introduce a novel approach termed P2P-DD for structure-based editing tasks. P2P-DD combines
 171 Prompt-to-Prompt control with the idea of Directed Diffusion [21], extending its capability to perform
 172 object replacement.

173 Directed Diffusion [21] illustrates the possibility of achieving object replacement through direct
 174 editing of attention maps during the first N time steps of the generation process. The method performs
 175 attention strengthening on the cross-attention maps corresponding to the tokens in \mathcal{I} (recall that \mathcal{I} is
 176 the token sets representing the object of interest), through calculating a Gaussian strengthening mask
 177 (SM) of the bounding-box region \mathcal{B} . It also performs attention annealing to the remaining area $\bar{\mathcal{B}}$ by
 178 applying a constant weakening mask (WM), and the two masks are weighted by a scalar c :

$$\text{DDEdit}(M_t, \mathcal{B}, \mathcal{I}) = \begin{cases} M_t \odot \text{WM}(\bar{\mathcal{B}}, \mathcal{I}) + c \cdot \text{SM}(\mathcal{B}, \mathcal{I}) & \text{if } t > T - N \\ M_t & \text{otherwise} \end{cases} \quad (2)$$

179 Our proposed P2P-DD method aims to encourage the background and other details of the image
 180 generated by Directed Diffusion with a higher resemblance to the source image. P2P-DD first injects
 181 cross-attention maps M_t into M_t^* as suggested by Prompt-to-Prompt. This will preserve the structure
 182 and background information from the source image, providing a decent starting point for attention
 183 map editing. Next, P2P-DD performs pixel value strengthening and weakening on the attention maps
 184 M_t according to the bounding-box \mathcal{B} and the tokens in \mathcal{I} , as suggested in the Directed Diffusion, to
 185 achieve the desired object placement:

$$\text{P2P-DDEdit}(M_t, M_t^*, \mathcal{B}, \mathcal{I}) = \text{DDEdit}(\text{P2PEdit}(M, M^*, t), \mathcal{B}, \mathcal{I}) \quad (3)$$

186 This combined editing control is only applied in the first N diffusion time steps, i.e. $\{T, \dots, T - N + 1\}$.
 187 After the first N time steps, we cease any control on attention maps and allow the target diffusion
 188 process to complete in a conventional denoising manner. The details of P2P-DD are shown in Figure
 189 2. Attention visualization of the generation process can be found in Appendix A.1, and we provide a
 190 detailed description of the full algorithm in Appendix A.2.

191 3.2 Example-based visual reinforcement learning

192 We modify and extend the approach of VICE [5] to devise our method for example-based visual
 193 reinforcement learning. At the core of our process, we employ the image editing methods described
 194 in Section 3.1 on observations gathered from randomly initialized environments, producing a dataset
 195 of 1024 target images for each task. During training, these target images serve as positive samples,
 196 while the images sampled from the agent’s replay buffer after a number of random exploration
 197 steps are treated as negative samples. We train a discriminator using these instances with the binary

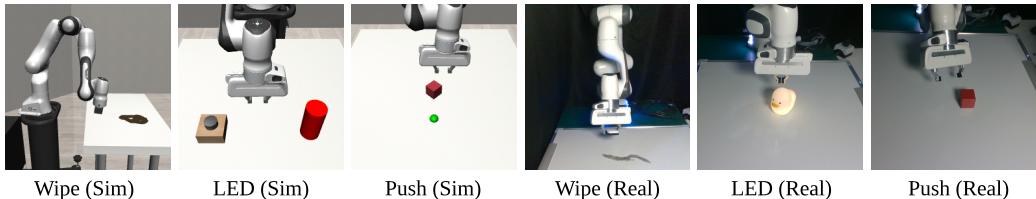


Figure 3: **Visualization of robot manipulation tasks.** We evaluate LfVoid on three simulation tasks: Wipe, Push, and LED, as well as three real-robot environments correspondingly.

198 cross-entropy loss and use the output of the discriminator for new observations for the positive class
 199 as the agent reward.

200 The discriminator shares the CNN encoder with the reinforcement learning agent and performs
 201 classification over the output latent representations. To improve the reward shaping, we utilize
 202 the label mixup method [6], which performs a random linear interpolation of the 0-1 labels and
 203 their corresponding hidden vectors to obtain continuous labels between 0 and 1. We discover that
 204 restricting the negative instances from the recent portion of the replay buffer (the last 5%) promotes
 205 the discriminator’s ability to discern subtle differences between the target and current observations.
 206 Additionally, we find that our method enjoys an ensemble of discriminators for classification results
 207 (i.e., RL rewards), which gives the agent more representative rewards.

208 For the reinforcement learning backbone algorithm, we use DrQ-v2 [42] for visual RL training. Based
 209 on Twin Delayed DDPG (TD3) [43], DrQ-v2 enhances image representation by applying random
 210 crop augmentation to the input images. Collectively, these refinements constitute our approach to
 211 example-based visual reinforcement learning, making the “learning” from an initial observation and
 212 language prompts possible.

213 4 Experiments

214 In this section, we present a comprehensive evaluation of LfVoid across both simulated environments
 215 and real-world robotic tasks. An illustration of each task can be found in Figure 3.

216 The environments we use are: 1) LED-light, where the robot reaches for a switch or touches the
 217 light directly to turn the light from red to green; 2) Wipe, where the robot needs to wipe out stains
 218 from the table; 3) Push, where the robot needs to push a red cube to a goal position indicated by
 219 a green dot. The simulated tasks are developed based on the Robosuite benchmark [44], while we
 220 provide corresponding real-world tasks for each environment. A full description of the environments
 221 is provided in Appendix B.1.

222 4.1 Goal generation

223 In this section, we evaluate the ability of LfVoid and other baseline methods to generate goal images
 224 according to two types of editing instructions: appearance-based editing and structure-based editing.

225 **Appearance-based goal image editing results.** In Figure 4, we present the visual goal generation
 226 results of the “Wipe” and “LED” tasks, focusing on editing object appearances in both simulated
 227 and real-world environments. We compare our method with two existing image editing methods:
 228 Imagic [20] and Instruct-pix2pix [26]. Imagic struggles to preserve the texture and details of the
 229 image background, often resulting in significant alterations to the source image. Instruct-pix2pix fails
 230 to effectively locate the specified regions indicated by the editing instructions, frequently causing
 231 global color changes and failing to perform the proper edits. Our method demonstrates superior
 232 performance in terms of maintaining the appearance of unrelated parts while effectively editing
 233 specified regions. As mentioned in Section 3.1.3, through directly applying control and editing on the
 234 attention maps, LfVoid is capable of locating the target area and performing the desired edit based on
 235 language instructions.

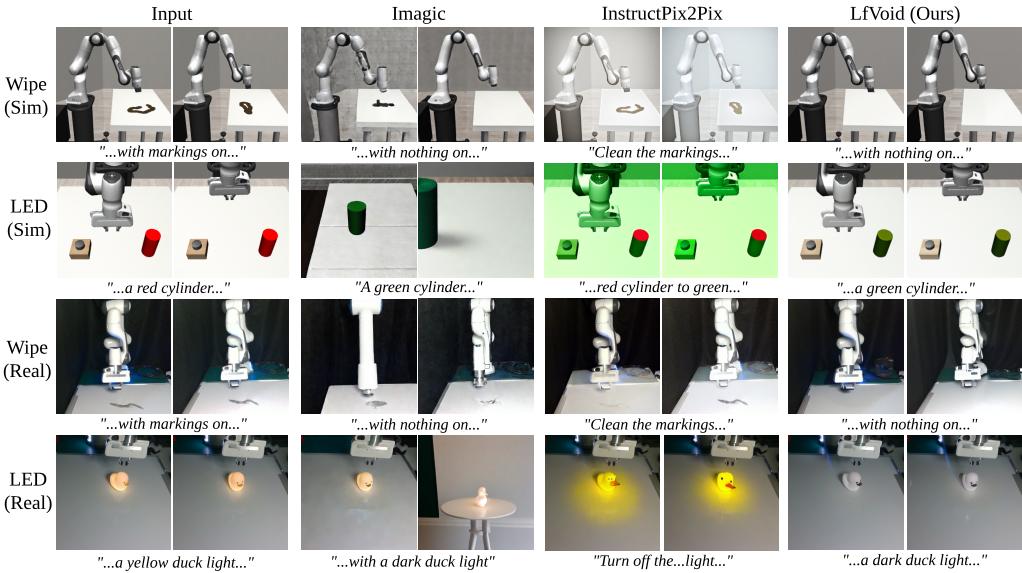


Figure 4: **Appearance-based goal image editing results.** We compare LfVoid against two recent image editing methods, Imagic and InstructPix2Pix, on four editing tasks. In both simulation and real settings, LfVoid can generate images that are better aligned with given text prompts while preserving the remaining source scene. See Appendix C for more examples.

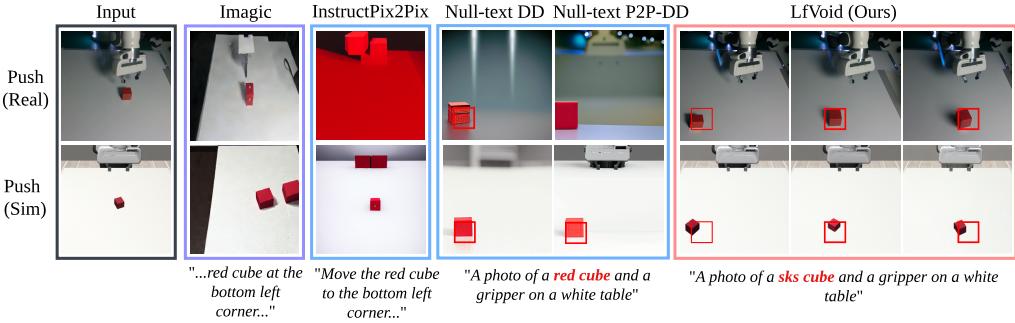


Figure 5: **Structure-based goal image editing results.** For structure-based editing tasks, we compare LfVoid against Imagic and InstructPix2Pix, as well as two ablations, Null-text Directed Diffusion (Null-text DD) and Null-text P2P-DD. LfVoid can successfully perform object displacement and preserves the background and details of the source image, while other methods fail to do so. See Appendix C for more examples.

236 **Structure-based goal image editing results.** Next, we report the performance of LfVoid in editing
 237 tasks with structural changes ("Push") in both simulated and real-world settings. In Figure 5, we
 238 first present the editing results of Imagic and Instruct-pix2pix: Imagic distorts the robot arm and
 239 table significantly and often introduces multiple objects into the image, rather than moving the
 240 existing one. Instruct-pix2pix only changes the overall color of the image or replaces the gripper
 241 of the robot arm with a geometric body, failing to achieve the required object displacement. In
 242 addition, we also present ablative results of LfVoid when certain modules are removed: Null-text DD
 243 removes the feature extracting module and P2P in the editing module. Although it can successfully
 244 move objects to the lower-left corner, the background and the robot arm nearly disappear. Null-text
 245 P2P-DD removes only the feature extracting module of LfVoid, and results show that it improves
 246 the preservation of the source image's background compared to Null-text DD, but the shape of the
 247 moved object is inconsistent with its original form. Our method, when compared to all baselines, can
 248 perform the object displacement task successfully while better preserving the visual features of the
 249 object and the remaining scene. Moreover, we also demonstrate LfVoid's ability to perform object

displacement at different locations, as shown in the last three columns of Figure 5. These experiments reveal the importance of each component in our method. The creative integration of different image editing techniques provides a powerful tool for image editing with structural changes, demonstrating a significant improvement over existing methods.

Human user evaluation. In Table 1, we present the user study results in terms of normalized Elo score. The evaluators show a strong preference for images generated by LfVoid compared to other methods. Details of the user study and the Elo metrics can be found in Appendix B.4.

	Imagic	Instruct-pix2pix	NT-DD	NT-P2P-DD	Ours
Appearance-Based (\uparrow)	47.7	42.4	N/A	N/A	87
Structure-Based (\uparrow)	43.6	18.2	58	63.8	91.6

Table 1: **User study on visual goal generation.** The Elo score shows user preference among multiple choices, higher is better.

4.2 Example-based visual reinforcement learning

We now dive into the results obtained from applying example-based visual reinforcement learning using the goal images generated by LfVoid. We select three baselines for our comparison. The first baseline uses language guidance by leveraging the Contrastive Language-Image Pre-training (CLIP) score [45] of the task prompt and image observations as a reward to train an RL agent. This baseline was chosen to verify the effectiveness of translating prompts into images against raw language hints. The second baseline uses InstructPix2Pix (IP2P) for image editing and then applies the same example-based RL as in LfVoid on the resulting images. It allows us to compare the performance of LfVoid with a straightforward image editing technique. We also present an upper bound performance of our pipeline (marked as Real Goal), where we manually create ideal goal images, such as erasing the stains directly in the simulator for the "Wipe" task. This strategy illustrates the potential of our algorithm and indicates the gap between the generated goals and the real ones.

Simulation results. We report both episode reward (Figure 6) and numerical metrics (Table 2 on three simulation tasks. CLIP baseline achieves low rewards in all the tasks, showing that directly using language guidance through CLIP embeddings cannot provide sufficient guidance. Moreover, LfVoid consistently outperforms the IP2P baseline and gains comparable performance with the real goal upper bound, providing evidence that generated goal images with high fidelity to both the original scene and the editing instructions is crucial for the deployment of example-based RL methods.

Real-world results. In this experiment, we aim to validate the feasibility of LfVoid in more complex real-world settings through visualization of reward functions. As in our previous experiment, we initially generate goal images using InstructPix2Pix (IP2P) and LfVoid, and concurrently collect a real goal dataset as an upper-bound (Real Goal). We then train discriminators on these datasets separately. We evaluate the discriminators on five successful trajectories for each task on a real robot

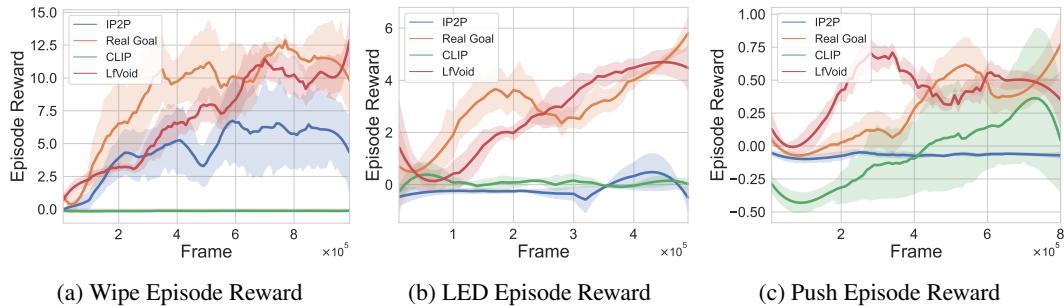


Figure 6: **Episode reward curve of simulation tasks.** The results of the CLIP baseline (CLIP), InstructPix2Pix baseline (IP2P), using real goal image as an upper bound (Real Goal), and LfVoid (Ours)

	CLIP	Instruct-pix2pix	Real Goal	Ours
Wipe (Cleaned Stains / Patch)	0.0±0.0	1.7±1.9	22.0±5.8	21.3±5.8
LED (Success Rate / %)	10.0±20.0	0.0±0.0	93.3±11.5	75.0±50.0
Push (Success Rate / %)	12.5±19.1	3.3±3.5	30.0±13.7	27.9±12.7

Table 2: **Numerical metrics of simulation tasks.** We report the success rate for LED and Push, and the number of stain patches cleaned for Wipe, please refer to Appendix B.1 for details.

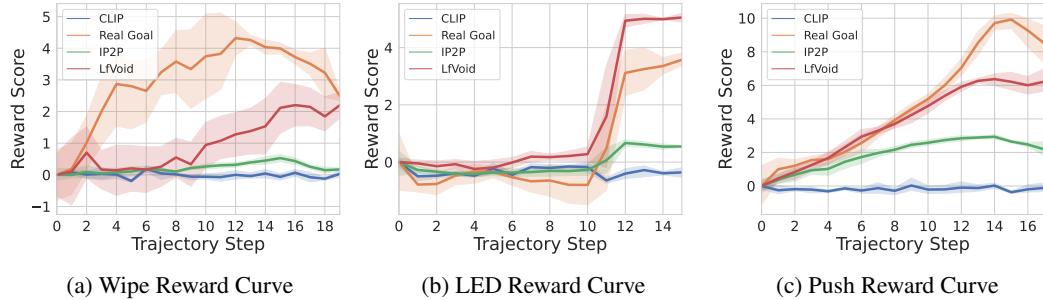


Figure 7: **Reward curve on successful trajectories for real-robot tasks.** Visualization of the reward function of LfVoid (Ours) compared with baselines.

and obtain the reward curve. In addition, we compare our approach with language-based guidance, which uses the distance of the CLIP embedding between the task prompt and the current observation image as rewards.

As illustrated in Figure 7, the reward curve derived from LfVoid can provide near monotonic dense signals comparable to those from the real goals. On the contrary, the curves obtained by IP2P and the CLIP score fail to accurately track the progress of each task. This demonstrates that our method offers a more instructive learning signal for real-world robotic tasks than either IP2P or CLIP.

5 Discussion

In this work, we present LfVoid, an effective approach for leveraging the knowledge of large-scale text-to-image models and enabling zero-shot reinforcement learning from text prompts and raw observation images. We have identified and tackled numerous challenges that arise when applying state-of-the-art image editing technologies to example-based RL methods. Our work highlights the potential for the adaptation and application of image-generation techniques in the realm of robotics. Our findings not only enhance the understanding of image editing for robotic applications but also provide a clear direction for the image generation community to address real-world challenges.

Although our proposed method, LfVoid, has shown promising results, we acknowledge that it is not without limitations. While LfVoid has succeeded in generating goal images to train example-based RL policies, a certain level of prompt tuning is needed in order to achieve optimal editing performance. We refer the readers to Appendix D for an analysis of failure cases.

We would also like to highlight that there exists a gap between the ability of text-to-image generation models and the need for robot learning. For example, large diffusion models exhibit poor understandings of the spatial relationships between objects [46]; therefore, both the generative models and the editing methods struggle to handle object displacement solely through language prompts. The Directed Diffusion technique, while being able to perform movements of objects, requires a user-defined bounding box to achieve precise control. Furthermore, large generative models sometimes struggle to generate images that are considered valid under physics laws. When asked to pick up a bottle with a robot arm, the model simply stretches the shape of the arm to reach the bottle rather than changing the joint position. Lastly, AI-generated images inevitably introduce alterations to the details of objects. The robustness of current visual Reinforcement Learning (RL) algorithms to such changes remains an open question.

310 **References**

- 311 [1] Georgios Papagiannis and Yunpeng Li. Imitation learning with sinkhorn distances. In *Machine*
312 *Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2022,*
313 *Grenoble, France, September 19–23, 2022, Proceedings, Part IV*, pages 116–131. Springer,
314 2023.
- 315 [2] Robert Dadashi, Léonard Hussenot, Matthieu Geist, and Olivier Pietquin. Primal wasserstein
316 imitation learning. *arXiv preprint arXiv:2006.04678*, 2020.
- 317 [3] Ilya Kostrikov, Kumar Krishna Agrawal, Debidatta Dwibedi, Sergey Levine, and Jonathan
318 Tompson. Discriminator-actor-critic: Addressing sample inefficiency and reward bias in
319 adversarial imitation learning. *arXiv preprint arXiv:1809.02925*, 2018.
- 320 [4] Siddhant Haldar, Vaibhav Mathur, Denis Yarats, and Lerrel Pinto. Watch and match: Super-
321 charging imitation with regularized optimal transport. In *Conference on Robot Learning*, pages
322 32–43. PMLR, 2023.
- 323 [5] Justin Fu, Avi Singh, Dibya Ghosh, Larry Yang, and Sergey Levine. Variational inverse control
324 with events: A general framework for data-driven reward definition. *Advances in neural*
325 *information processing systems*, 31, 2018.
- 326 [6] Avi Singh, Larry Yang, Kristian Hartikainen, Chelsea Finn, and Sergey Levine. End-to-end
327 robotic reinforcement learning without reward engineering. *arXiv preprint arXiv:1904.07854*,
328 2019.
- 329 [7] Ben Eysenbach, Sergey Levine, and Russ R Salakhutdinov. Replacing rewards with examples:
330 Example-based policy search via recursive classification. *Advances in Neural Information*
331 *Processing Systems*, 34:11541–11552, 2021.
- 332 [8] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David,
333 Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, et al. Do as i can, not
334 as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.
- 335 [9] Anthony Brohan, Noah Brown, Justice Carbalal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn,
336 Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics
337 transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.
- 338 [10] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical
339 text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- 340 [11] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton,
341 Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al.
342 Photorealistic text-to-image diffusion models with deep language understanding. *Advances in*
343 *Neural Information Processing Systems*, 35:36479–36494, 2022.
- 344 [12] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer.
345 High-resolution image synthesis with latent diffusion models, 2021.
- 346 [13] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay
347 Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, Benton C. Hutchinson, Wei Han,
348 Zarana Parekh, Xin Li, Han Zhang, Jason Baldridge, and Yonghui Wu. Scaling autoregressive
349 models for content-rich text-to-image generation. *ArXiv*, abs/2206.10789, 2022.
- 350 [14] Stephen Adams, Tyler Cody, and Peter A Beling. A survey of inverse reinforcement learning.
351 *Artificial Intelligence Review*, 55(6):4307–4346, 2022.
- 352 [15] Zhao Mandi, Homanga Bharadhwaj, Vincent Moens, Shuran Song, Aravind Rajeswaran, and
353 Vikash Kumar. Cacti: A framework for scalable multi-task multi-scene visual imitation learning.
354 *arXiv preprint arXiv:2212.05711*, 2022.
- 355 [16] Zoey Chen, Sho Kiami, Abhishek Gupta, and Vikash Kumar. Genaug: Retargeting behaviors to
356 unseen situations via generative augmentation. *arXiv preprint arXiv:2302.06671*, 2023.

- 357 [17] Tianhe Yu, Ted Xiao, Austin Stone, Jonathan Tompson, Anthony Brohan, Su Wang, Jaspiar
358 Singh, Clayton Tan, Jodilyn Peralta, Brian Ichter, et al. Scaling robot learning with semantically
359 imagined experience. *arXiv preprint arXiv:2302.11550*, 2023.
- 360 [18] Ivan Kapelyukh, Vitalis Vosylius, and Edward Johns. Dall-e-bot: Introducing web-scale
361 diffusion models to robotics. *arXiv preprint arXiv:2210.02438*, 2022.
- 362 [19] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or.
363 Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*,
364 2022.
- 365 [20] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri,
366 and Michal Irani. Imagic: Text-based real image editing with diffusion models. *arXiv preprint
367 arXiv:2210.09276*, 2022.
- 368 [21] Wan-Duo Kurt Ma, JP Lewis, W Bastiaan Kleijn, and Thomas Leung. Directed diffusion: Direct
369 control of object placement through attention guidance. *arXiv preprint arXiv:2302.13153*, 2023.
- 370 [22] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances
371 in Neural Information Processing Systems*, 33:6840–6851, 2020.
- 372 [23] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv
373 preprint arXiv:2010.02502*, 2020.
- 374 [24] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis.
375 *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021.
- 376 [25] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew,
377 Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing
378 with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- 379 [26] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow
380 image editing instructions. *arXiv preprint arXiv:2211.09800*, 2022.
- 381 [27] Parsa Mahmoudieh, Deepak Pathak, and Trevor Darrell. Zero-shot reward specification via
382 grounded natural language. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepes-
383 vari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference
384 on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages
385 14743–14752. PMLR, 17–23 Jul 2022.
- 386 [28] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman.
387 Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. *arXiv
388 preprint arXiv:2208.12242*, 2022.
- 389 [29] Kevin Li, Abhishek Gupta, Ashwin Reddy, Vitchyr H Pong, Aurick Zhou, Justin Yu, and Sergey
390 Levine. Mural: Meta-learning uncertainty-aware rewards for outcome-driven reinforcement
391 learning. In *International conference on machine learning*, pages 6346–6356. PMLR, 2021.
- 392 [30] Daesol Cho, Seungjae Lee, and H Jin Kim. Outcome-directed reinforcement learning by uncer-
393 tainty & temporal distance-aware curriculum goal generation. *arXiv preprint arXiv:2301.11741*,
394 2023.
- 395 [31] Maria Attarian, Advaya Gupta, Ziyi Zhou, Wei Yu, Igor Gilitschenski, and Animesh Garg.
396 See, plan, predict: Language-guided cognitive planning with video prediction. *arXiv preprint
397 arXiv:2210.03825*, 2022.
- 398 [32] Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng,
399 Jonathan Tompson, Igor Mordatch, Yevgen Chebotar, et al. Inner monologue: Embodied
400 reasoning through planning with language models. *arXiv preprint arXiv:2207.05608*, 2022.
- 401 [33] Yan Ding, Xiaohan Zhang, Chris Paxton, and Shiqi Zhang. Task and motion planning with
402 large language models for object rearrangement. *arXiv preprint arXiv:2303.06247*, 2023.
- 403 [34] Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Cliport: What and where pathways for robotic
404 manipulation. In *Conference on Robot Learning*, pages 894–906. PMLR, 2022.

- 405 [35] Yunfan Jiang, Agrim Gupta, Zichen Zhang, Guanzhi Wang, Yongqiang Dou, Yanjun Chen,
406 Li Fei-Fei, Anima Anandkumar, Yuke Zhu, and Linxi Fan. Vima: General robot manipulation
407 with multimodal prompts. *arXiv preprint arXiv:2210.03094*, 2022.
- 408 [36] Weiyu Liu, Yilun Du, Tucker Hermans, Sonia Chernova, and Chris Paxton. Structdiffusion:
409 Language-guided creation of physically-valid structures using unseen objects, 2023.
- 410 [37] Nikolaos Gkanatsios, Ayush Jain, Zhou Xian, Yunchu Zhang, Christopher Atkeson, and Katerina
411 Fragkiadaki. Energy-based models as zero-shot planners for compositional scene rearrangement.
412 *arXiv preprint arXiv:2304.14391*, 2023.
- 413 [38] Utkarsh A Mishra and Yongxin Chen. Reorientdiff: Diffusion model based reorientation for
414 object manipulation. *arXiv preprint arXiv:2303.12700*, 2023.
- 415 [39] Yilun Dai, Mengjiao Yang, Bo Dai, Hanjun Dai, Ofir Nachum, Josh Tenenbaum, Dale Schuur-
416 mans, and Pieter Abbeel. Learning universal policies via text-guided video generation. *arXiv
417 preprint arXiv:2302.00111*, 2023.
- 418 [40] Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion
419 for editing real images using guided diffusion models. *arXiv preprint arXiv:2211.09794*, 2022.
- 420 [41] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint
421 arXiv:2207.12598*, 2022.
- 422 [42] Denis Yarats, Ilya Kostrikov, and Rob Fergus. Image augmentation is all you need: Regu-
423 larizing deep reinforcement learning from pixels. In *International Conference on Learning
424 Representations*, 2021.
- 425 [43] Scott Fujimoto, Herke van Hoof, and David Meger. Addressing function approximation error in
426 actor-critic methods. *ArXiv*, abs/1802.09477, 2018.
- 427 [44] Yuke Zhu, Josiah Wong, Ajay Mandlekar, Roberto Martín-Martín, Abhishek Joshi, Soroush
428 Nasiriany, and Yifeng Zhu. robosuite: A modular simulation framework and benchmark for
429 robot learning. In *arXiv preprint arXiv:2009.12293*, 2020.
- 430 [45] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
431 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
432 models from natural language supervision. In *International conference on machine learning*,
433 pages 8748–8763. PMLR, 2021.
- 434 [46] Colin Conwell and Tomer Ullman. Testing relational understanding in text-guided image
435 generation. *arXiv preprint arXiv:2208.00005*, 2022.