

EraseDraw: Learning to Draw Step-by-Step via Erasing Objects from Images

Alper Canberk, Maksym Bondarenko, Ege Ozguroglu, Ruoshi Liu, and Carl Vondrick

Columbia University
erasedraw.cs.columbia.edu

Abstract. Creative processes such as painting often involve creating different components of an image one by one. Can we build a computational model to perform this task? Prior works often fail by making global changes to the image, inserting objects in unrealistic spatial locations, and generating inaccurate lighting details. We observe that while state-of-the-art models perform poorly on object insertion, they can remove objects and erase the background in natural images very well. Inverting the direction of object removal, we obtain high-quality data for learning to insert objects that are spatially, physically, and optically consistent with the surroundings. With this scalable automatic data generation pipeline, we can create a dataset for learning object insertion, which is used to train our proposed text-conditioned diffusion model. Qualitative and quantitative experiments have shown that our model achieves state-of-the-art results in object insertion, particularly for in-the-wild images. We show compelling results on diverse insertion prompts and images across various domains. In addition, we automate iterative insertion by combining our insertion model with beam search guided by CLIP.

Keywords: Object Insertion · Diffusion Models · Image Editing

1 Introduction

Inserting objects into an image based on a language prompt is a challenging task but has many applications in image editing and content creation in general. There are many fundamental challenges in solving this task. Inserted objects need to appear at physically plausible spatial locations that respect the natural distribution of images. Existing objects in the scene must be preserved. The appearance of the inserted objects needs to be consistent with that of the context. The lighting details should match the environmental lighting. Any one of these is traditionally challenging in computer vision, let alone combined.

Prior works represented by InstructPix2Pix [61] approached this task by generating pairs of images from pairs of prompts using Prompt2Prompt [17], which is subsequently used to train a model for language-conditioned image editing tasks. While achieving impressive results on image editing tasks such as changing styles or settings, these methods often fail on the tasks of object insertion by making global edits to the image, replacing an existing object when inserting new ones, and struggling to spatially reason

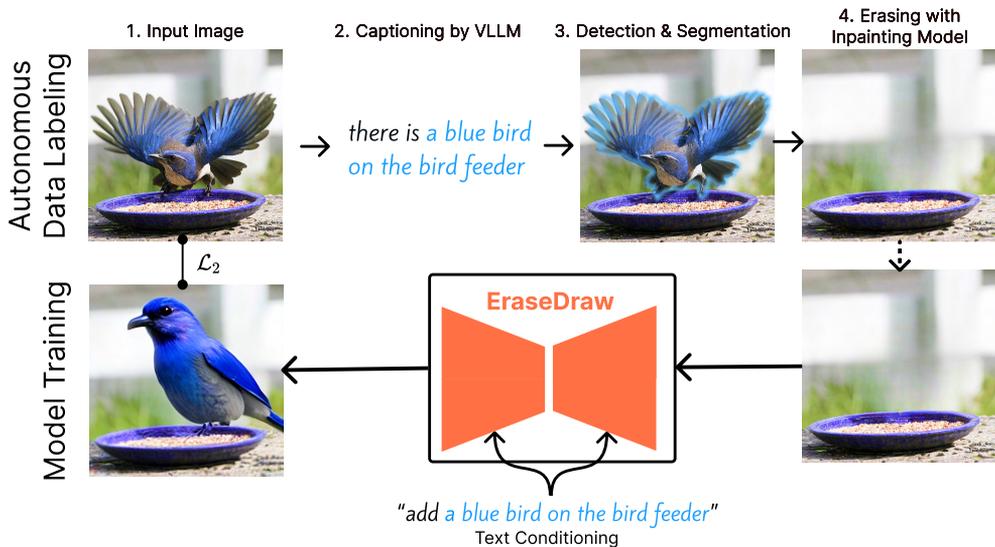


Fig. 1: EraseDraw We leverage advancements in image understanding and inpainting to train a model that can insert an object given a language instruction.

(see Figure 2). Another line of work attempts to create the training data by rendering from a simulation environment and synthetic objects [69]. However, a significant sim2real gap exists, preventing it from generalizing well to in-the-wild images.

In this paper, we propose EraseDraw, a scalable system for learning the task of language-conditioned object insertion into images. We observe that with modern segmentation, captioning and inpainting models, we can perform the task of object removal with much higher photo and physical realism than object insertion. With this observation, we propose an autonomous data generation pipeline for generating input-output pairs for the task of instruction-guided object insertion. We modify a wide distribution of images from the internet by erasing objects from them using inpainting models and describing the attributes and locations of the erased objects using vision language models. As a result, we created a large-scale dataset of paired images as well as language prompts describing the object insertion.

With this dataset, we train a language-conditioned diffusion model to perform the task of object insertion. Experiments show that our model achieves state-of-the-art results, outperforming baselines trained with orders of magnitude more computing resources. Due to our autonomous data generation pipeline, we are able to create training data directly from in-the-wild images, resulting in high-quality object insertion data.

The primary contribution of this paper is a system for learning to insert objects in images and a scalable pipeline for generating the training data for this task from natural images. Insertion furthermore equips us with the ability to plan the iterative generation of novel images in a “step-by-step” manner. We believe the ability to insert

objects into natural images will have a significant impact on content creation, as well as other related areas, including computer graphics, AR/VR, and robotics.

2 Related Work

2.1 Image Editing

Generative models such as GANs [16, 23, 24] and diffusion models [19, 37, 73] have enabled various image editing tasks such as style transfer [20], image-to-image translation [22, 57], and latent space manipulation [46, 50]. More recently, text-guided diffusion models [70, 72, 73, 74] have allowed intuitive editing of images based on textual prompts.

Several methods have been proposed to enhance the controllability and precision of text-based image editing, however none of them have data necessary to perform object insertion. SDEdit [68] employs a stochastic differential equation for iterative denoising to increase the realism of user-provided pixel edits. Prompt2Prompt [18] and Null Text Inversion [35] modify cross-attention maps to enable both local and global editing. Imagic [64], EDICT [77] and Plug-and-Play [51] optimize text embeddings for better alignment between the input image and target description. Text2LIVE [5] and Blended Diffusion [59] train models to add edit layers or blend edited regions along the diffusion process. Imagen Editor [25] finetunes the diffusion model by inpainting masked objects.

Image Sculpting [75] presents a novel framework for editing 2D images by converting objects into 3D, allowing direct manipulation of their geometry, and re-rendering them back into the 2D image. Emu Edit [3] introduces a multi-task image editing model that achieves state-of-the-art results in instruction-based editing by training on a wide range of tasks and utilizing learned task embeddings. MagicBrush [78] improves instruction-based editing by finetuning InstructPix2Pix on a manually-annotated dataset collected using an online editing tool.

To provide more intuitive editing interfaces, InstructPix2Pix [61] introduced an instruction-based editing model trained on a synthetic dataset. MagicBrush [78] further improved it by finetuning on a manually-annotated dataset. However, these methods still struggle with accurately interpreting and precisely executing editing instructions, especially for object insertions. In contrast, our approach leverages the strength of inpainting models to automatically generate high-quality training data for learning object insertion.

2.2 Object Insertion

Determining where to place objects in images is a key problem for many editing tasks. Traditional approaches in computer graphics rely on manual specification [55] or synthetic data-driven methods [14]. In computer vision, early work used contextual information to predict likely object locations [11, 30, 56].

More recently, deep generative models have been used to learn object placements from data. Compositing GAN [60] generates realistic object composites by predicting



Fig. 2: State-of-the-art image editing methods fail to correctly insert objects into visual scenes. They perform global edits that don’t preserve scene context (**left**) [61], replace existing objects (**middle**) [78], and struggle to spatially reason (**right**) [71]. You may see how we did on these examples in Figure 17 of the Appendix.

geometric and appearance adjustments. RelaxedPlacement [66] optimizes for object positions and sizes to satisfy inter-object relationships described in scene graphs. Object-3DIT [69] studies 3D-aware object insertion using language instructions on synthetic data.

However, existing methods are often limited in their ability to handle complex, real-world object placements. Our key insight is that inpainting models can be used to erase objects from real images, providing valuable training data to learn meaningful object placements. By inverting this process, we show that we can train models to realistically insert diverse objects into images based on language instructions.

2.3 Diffusion Models

Denoising Diffusion Probabilistic Models (DDPM) [19] have marked their place in computer vision as a preferred generative architecture, thanks to their adeptness in handling multi-modal distributions, ensuring training stability, and offering scalability. First, diffusion models were shown to outperform GANs [16] in image generation [13], followed by Stable Diffusion (SD) demonstrating the efficient scalability of the approach [73]. SD achieved this through training on the internet-scale LAION-5B dataset [45], while employing diffusion in the latent space of a Variational Autoencoder (VAE) [27]. More recently, DDPMs were shown to be effective in video generation too, notably by Stable Video Diffusion [6] and Sora [9].

With the rise of internet-scale diffusion models, multiple works have studied their excellent ability to represent the natural image manifold, enabling zero-shot generalization in tasks such as 3D-reconstruction [12, 31, 32, 44, 52], segmentation [1, 53], amodal perception [38, 54], recognition [10, 29], as well as image editing [8, 15, 41]. However, these image-editing methods (2.1) significantly struggle in performing object-insertion. In this work, we leverage the ability of pre-trained diffusion models to represent the natural and multi-modal distribution of people and objects in visual scenes, thereby allowing us to perform object removal first (3.1), which by inversion enables insertion (3.2).

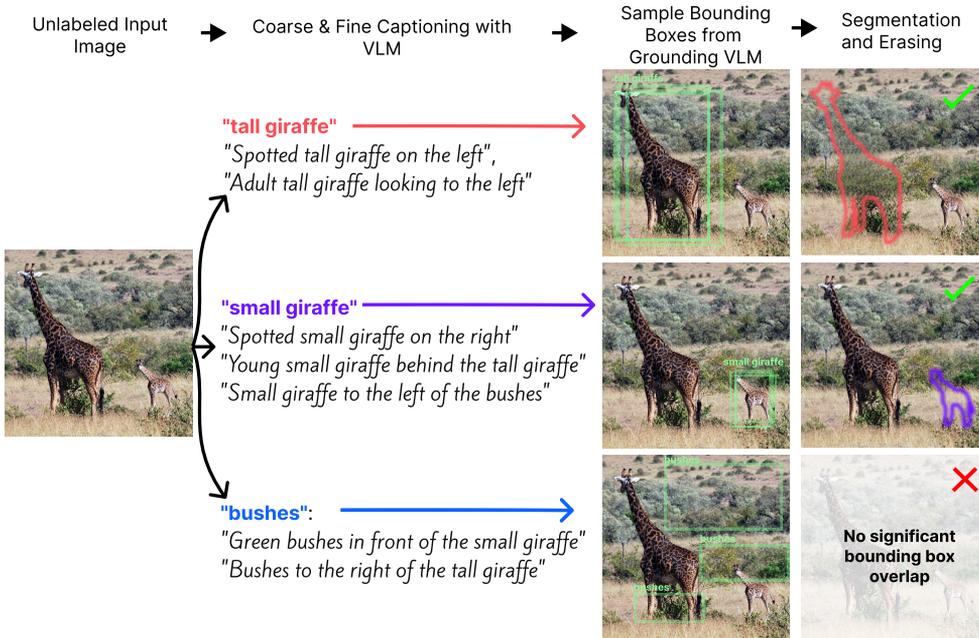


Fig. 3: EraseDraw Data Generation Pipeline (i) An unlabeled image is sampled taken from a dataset (ii) The images are given to a captioning model, which describes the objects in the image (iii) Objects are detected using the coarse caption from the captioning model, and the objects that are confidently detected are (iv) segmented, (v) and erased. The final images are added to the dataset along with the captions corresponding to them.

3 Method

Our main contribution is a framework for learning object insertion from natural images. This framework brings together the best aspects of existing pre-trained models: combining the powerful descriptive capabilities of multimodal LLMs, complementing it with attribute bound detection abilities of Grounding VLMs, and modifying images with precise segmentation and erasing models. After generating a synthetically annotated dataset of 65,000 images, we fine-tune a large pre-trained diffusion model on our dataset, the procedure for which we outline in Sec. 3.2.

3.1 Automatic Dataset Generation

For training a language-conditioned object insertion model in a supervised manner, one needs triplets in the form of (c_I, c_T, \mathbf{x}) , where c_I represents the source image, c_T describes the identity and location of the object to be inserted, and \mathbf{x} is modified version of c_I , which includes the object described in c_T . Our key insight is that we can let natural images from the internet be the target images \mathbf{x} , and we can autonomously derive the context c_T and c_I by detecting and erasing objects from these images. We

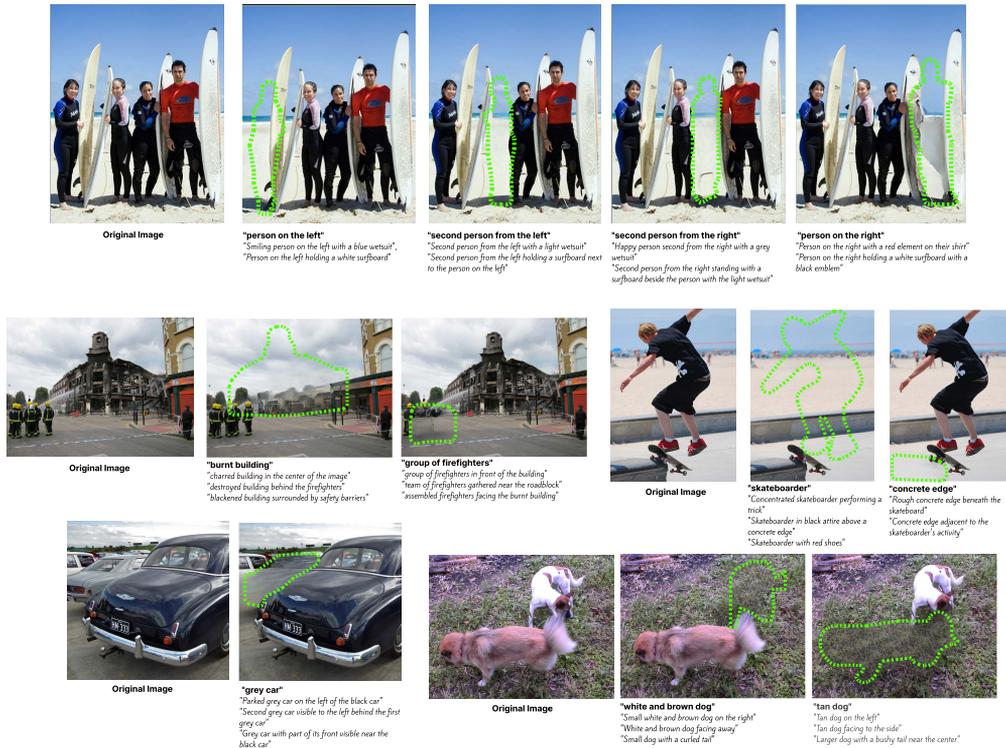


Fig. 4: We show examples from our EraseDraw Dataset.

illustrate the data generation framework in the following subsections, and we precisely describe our **algorithm** in the appendix.

Coarse & Fine Captioning To enhance object identification in images, a powerful Vision-Language Model (VLM), which is GPT-4 [58] in our implementation, is used to generate both coarse and fine captions for each object. Coarse captions provide a simple, yet unique description focusing on basic attributes like shape and color, aiding in the detection phase. Fine captions offer complex, detailed descriptions, enriching the object’s representation. This approach leverages the strength of grounding VLMs, such as CogVLM [80], in detecting objects using straightforward descriptions, while at training time, objects identified through coarse captions can be associated with any of their detailed fine captions. For example, a "red bowl" might be simply identified, but further described in detail as "a dark red bowl on top of the shelf." This methodology ensures efficient object detection while enabling rich, descriptive training data for VLMs.

Detection with Grounding VLM & Segmentation Standard open-vocabulary object detectors struggle with accurately identifying colors, sizes, or spatial relation-

ships, leading to inaccuracies in datasets with multiple instances of the same object type. To address this, we utilize CogVLM’s attribute-binding capabilities, which excel in recognizing these features, as demonstrated in Fig. 3 and 4. By inputting coarse captions into CogVLM, we generate bounding boxes for objects, though errors can still arise from VLM hallucinating objects or difficulty in identification (e.g., bushes scenario in Fig 3).

To mitigate these issues, we employ rejection sampling based on CogVLM’s uncertainty, using it as a transformer-based probabilistic model. By sampling three bounding boxes at a temperature of 0.2 and accepting those with an Intersection over Union (IoU) exceeding 80%, we ensure the object’s presence and identifiability by the consistency of bounding box locations. Objects with scattered bounding boxes are excluded from the dataset. Identified objects are then segmented using SAM [65] and prepared for the erasing stage, refining the dataset’s accuracy for training purposes.

Erasing We employ the LaMa inpainting model [79] for erasing objects using segmentation masks. Although LaMa is less versatile compared to advanced models like Stable Diffusion for inpainting, it is effective in consistently erasing objects. A minor issue with LaMa is artifact creation when erasing large objects, but this does not greatly impact our process due to two main reasons: firstly, images with artifacts serve as conditional inputs rather than direct supervision; secondly, artifacts often diminish after processing through the VAE for latent diffusion, mitigating their presence.

3.2 Diffusion Model

Our goal is to sample from $p(\mathbf{x}|c_I, c_T)$, the distribution of target images conditioned on a source image and a text instruction. To this end, we train a latent conditional diffusion model $\epsilon_\theta(\cdot)$, which estimates the score function of $p(\mathbf{x}|c_I, c_T)$ by optimizing for the simplified variational lower bound objective

$$\min_{\theta} \mathbb{E}_{\mathcal{E}(\mathbf{x}), \mathcal{E}(c_I), c_T, \epsilon \sim \mathcal{N}(0,1), t} [\|\epsilon_\theta(\mathbf{z}_t, t, \mathcal{E}(c_I), c_T) - \epsilon\|_2^2]$$

where \mathbf{x} is an image sampled from the dataset, $\mathcal{E}(\cdot)$ is a VAE, and $\mathbf{z}_t = \mathcal{E}(\mathbf{x})$ is a noisy latent embedding of \mathbf{x} where the noise level increases with $t \in T$.

In essence, the diffusion model learns to predict the noise ϵ added to a ground-truth latent image \mathbf{z}_0 . To sample from the model, we start with pure Gaussian noise tensor and repeatedly invoke the diffusion model to predict and subtract noise, which iteratively denoises the pure noise tensor into a latent image. Finally, we invoke the decoder \mathcal{D} to obtain a full-resolution image from the latent image.

Following [61], we initialize the weights of our network ϵ_θ from Stable Diffusion 1.5, a text-to-image model pretrained on a large-scale dataset of images. Similarly, we allow for conditioning on images by expanding the number of input channels of the convolutional input layer. Initialization from a pretrained model allows our fine-tuned model to generate objects from an open-vocabulary, well beyond the set of objects that are in our fine-tuning dataset of 65,000 images.

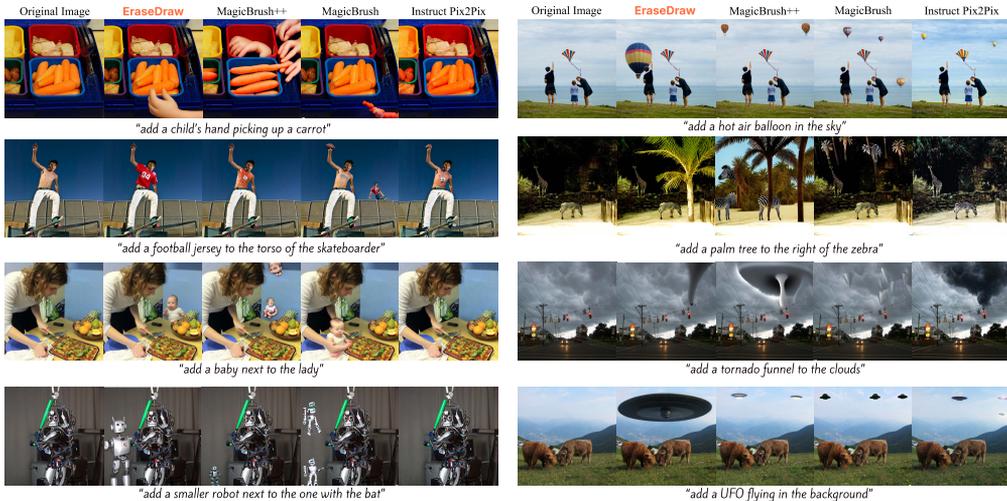


Fig. 5: Qualitative Results on EmuEdit Benchmark on inserting people and outdoor objects.

3.3 Iterative Generation and Planning

Our model’s ability to precisely insert an object into an existing image allows the composition of novel images in a "step-by-step" manner.

In particular, given a background image and a sequence of desired edits, one can repeatedly invoke EraseDraw, sequentially prompting the model with each edit, and feeding in the output image of the last step as the input to the next.

Automating Step-by-Step Composition To compose a complex image by its constituents, a user needs to decide the order in which to insert each object. For each inserted object, the user may also choose to review a multitude of samples from the diffusion model to pick which edited image should be the input to the next step.

Alternatively, the decision-making and verification process can be offloaded to off-the-shelf vision and language models. We instantiate such an application by implementing *beam search* with a CLIP Score heuristic. Starting with an unmodified image and a pre-specified sequence of edits, for each step, we sample N images for each of the k 'beams', resulting in Nk images. Then, we compute text-image CLIP alignment between the Nk generated images and the editing instruction to rank images according to their quality and continue onto the next step with the top k .

The goal of our experiments is to evaluate our model’s ability to make precise, high-fidelity insertions that are faithful to the given prompt. To this end, we conduct quantitative and qualitative comparisons with publicly available language-guided image editing models on single-step object insertion (3.4). Then, we demonstrate our model’s ability to iteratively compose scenes step-by-step (3.5). Finally, we present quantitative ablations for our dataset recipe.



Fig. 6: Qualitative Results on EmuEdit Benchmark on inserting animals and household objects.

3.4 Single-Step Object Insertion

Setup Given an RGB image and a text prompt, the task of single-step object insertion aims to generate a new image with the specified object inserted while faithfully preserving the input scene. We compare the performance of EraseDraw with two similar-architecture models: InstructPix2Pix [61], trained on 450,000 auto-labeled images, and MagicBrush [78], a fine-tuned version of InstructPix2Pix with 5,000 human-generated examples. We are not able to compare against EmuEdit [76] since the model is not publicly available.

Metrics For evaluating object insertion, we use a part of EmuEdit’s [76] validation set designed for "add" instructions, comprising tuples of captions and images before and after edits. Models are judged by $CLIP_{out}$, comparing output captions with generated images, and $CLIP_{dir}$, evaluating the alignment between changes in captions and images. Lastly, to confirm our autonomously evaluated results, ten participants performed pairwise ranking tasks across 15 images and 4 models (90 pairwise comparisons per user) on a random subset of the EmuEdit insertion tasks. We report the win rate of each model.

Results Table 1 shows the effectiveness of our model on the EmuEdit object insertion benchmark, as indicated by higher $CLIP_{out}$ and $CLIP_{dir}$ scores, despite our model being trained on one order of magnitude less data. In Fig. 5 and 6, we qualitatively contrast our system with baselines, showcasing humans, animals, and indoor and outdoor scenes. These illustrate the nuanced, realistic, and vibrant edits of EraseDraw compared to the baseline models, which face numerous challenges as depicted in Fig. 2. In Fig. 7, we

Method	Dataset Size	CLIP _{out} ↑	CLIP _{dir} ↑	Win Rate (%)
Erasedraw	65,000	0.0930	0.1383	82.22
Magicbrush++	-	0.0896	0.1304	55.06
Magicbrush	455,000	0.0854	0.0930	48.64
InstructPix2Pix	450,000	0.0746	0.0976	12.56

Table 1: Comparison with language-guided editing baselines on the EmuEdit dataset

further showcase intricate edits like water flow simulation in a sink, object placement behind existing items, human-object interaction predictions, and lighting adjustments.

Additionally, in Fig. 8, we display our model’s capability to accurately reflect the natural spatial distribution of objects and people in a given scene. This capability highlights the advantage of learning from natural images, as opposed to simulated objects in prior work [69].

3.5 Iterative Generation

Setup In this section, we display qualitative results of the iterative object insertion process in both human-guided and automated settings. We start with empty background images and iteratively invoke our model with a predetermined sequence of instructions. To compare this process with the conventional one-shot text-to-image method, we combine the set of edit instructions into a single prompt and generate an image from Stable Diffusion 1.5.

Results According to Fig. 9, we notice that Stable Diffusion 1.5 fails to follow the entire prompt precisely, confusing attribute bindings or omitting insertions that would appear unnatural in the described context such as the ‘giraffe in the office. In comparison, our method results in complex final images that accurately match the edit instructions. This is a surprising result given that our base model is Stable Diffusion 1.5. Iteratively composing the image appears to have addressed important problems of attribute binding and out-of-distribution robustness. We attribute this to the fact that performing a single step of insertion is a much easier task than trying to satisfy all constraints of one-shot text-to-image image generation.

3.6 Beam Search

Setup We instantiate the beam search procedure outlined in 3.3 with branching factor $N = 4$, beam width of $k = 3$, for 3 iterations. For clarity, we only trace and display the top beam and the top-3 generations from that beam at each step. As a baseline, we sample TNk images from SD1.5 (our base model) using the combined prompt, and present the image with the highest CLIP score.

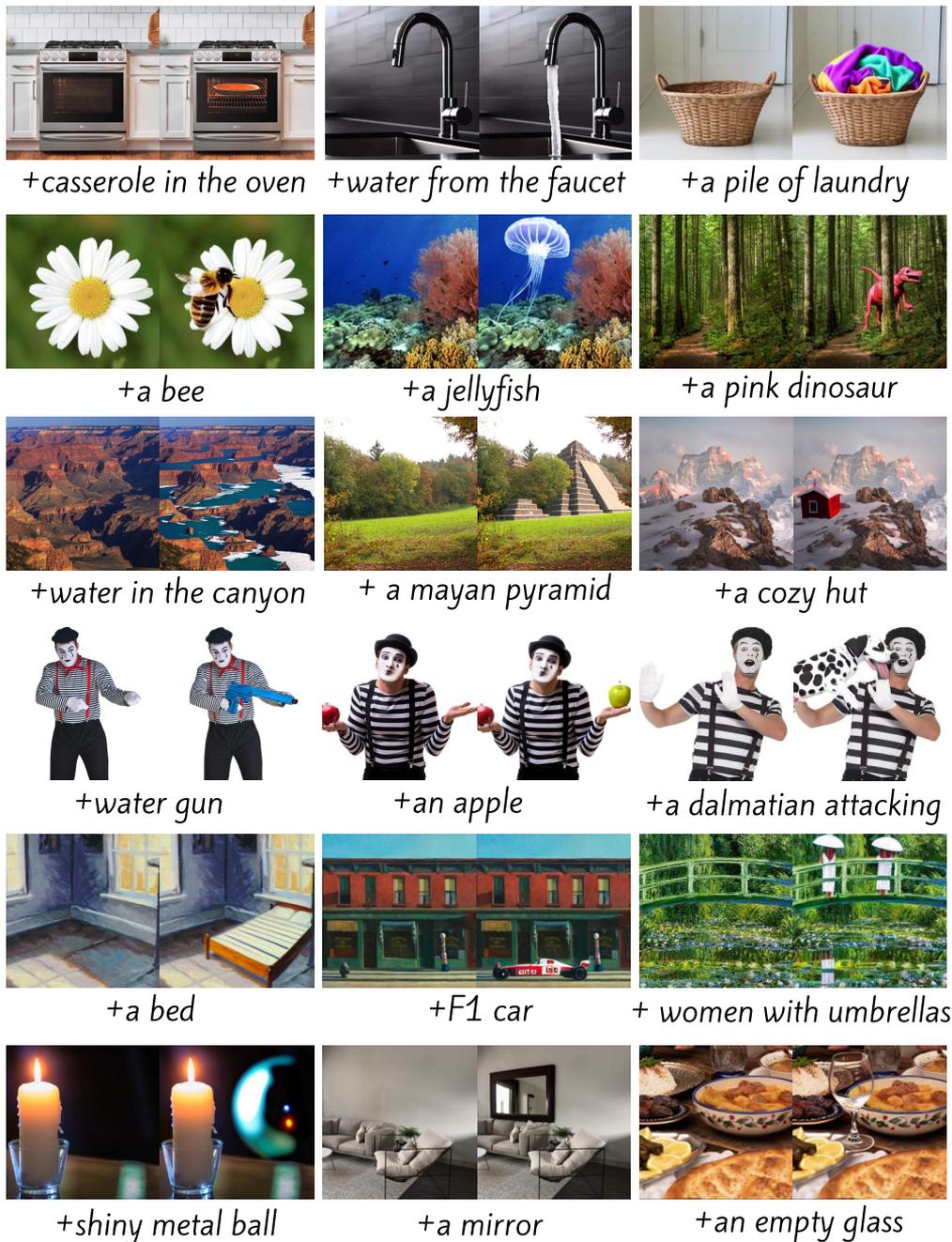


Fig. 7: Single-Step Generation. EraseDraw can perform complex object insertion tasks such as inserting water flowing down a sink, placing objects behind occluders, predicting hand-poses with correct affordance, and accounting for lighting effects.



Fig. 8: Modeling the Multimodal Distribution of People and Objects. Sampling from our model reveals where objects naturally appear in the world. This opens up applications where commonsense knowledge about object placements are required, such as embodied agents.

Results We keep the inference setting consistent with Sec. 3.5. Fig. 10 shows that CLIP score successfully finds the best candidates for the next step of insertion, resulting in a final image that composes an accurate image without human intervention. Meanwhile, SD1.5 mixes up attributes to their respective objects. While we have shown only 3 steps of edits in these evaluations, the performance between one-shot and iterative image generation increases sharply with additional steps. We include more results in figures 11, 12, 13, 14, and 15.

3.7 Dataset Ablation

Setup Finally, we perform ablation studies on our dataset generation pipeline to evaluate its contribution to the results. First, we use our end-to-end data generation pipeline to annotate 15,000 images, which we call 'EraseDraw-15k'. This pipeline could easily be scaled to orders of magnitude more images given enough compute. In addition, we convert 50,000 images from GQA [62] with annotated scene graphs into edit instructions by segmentation and erasing. We train EraseDraw on a combination of these two datasets.

Results Despite being autonomously generated and smaller size compared to GQA-Insert, the model trained on EraseDraw-15k competes well with its GQA-Insert counterpart, with the combination of both datasets yielding superior results (see Table 2). This performance differential can be attributed to the complementary qualities of the datasets: GQA-Insert's human-labeled accuracy versus EraseDraw-15k's broader diversity, albeit with more errors and limited volume.

4 Discussion

4.1 Conclusion and Limitations

We introduce EraseDraw, a framework for autonomously generating object insertion data by leveraging the fact that erasing is an easier task than drawing. We train a diffusion model using our data, and we show that it achieves state-of-the-art object insertion results. Additionally, we introduce a new paradigm for iteratively composing

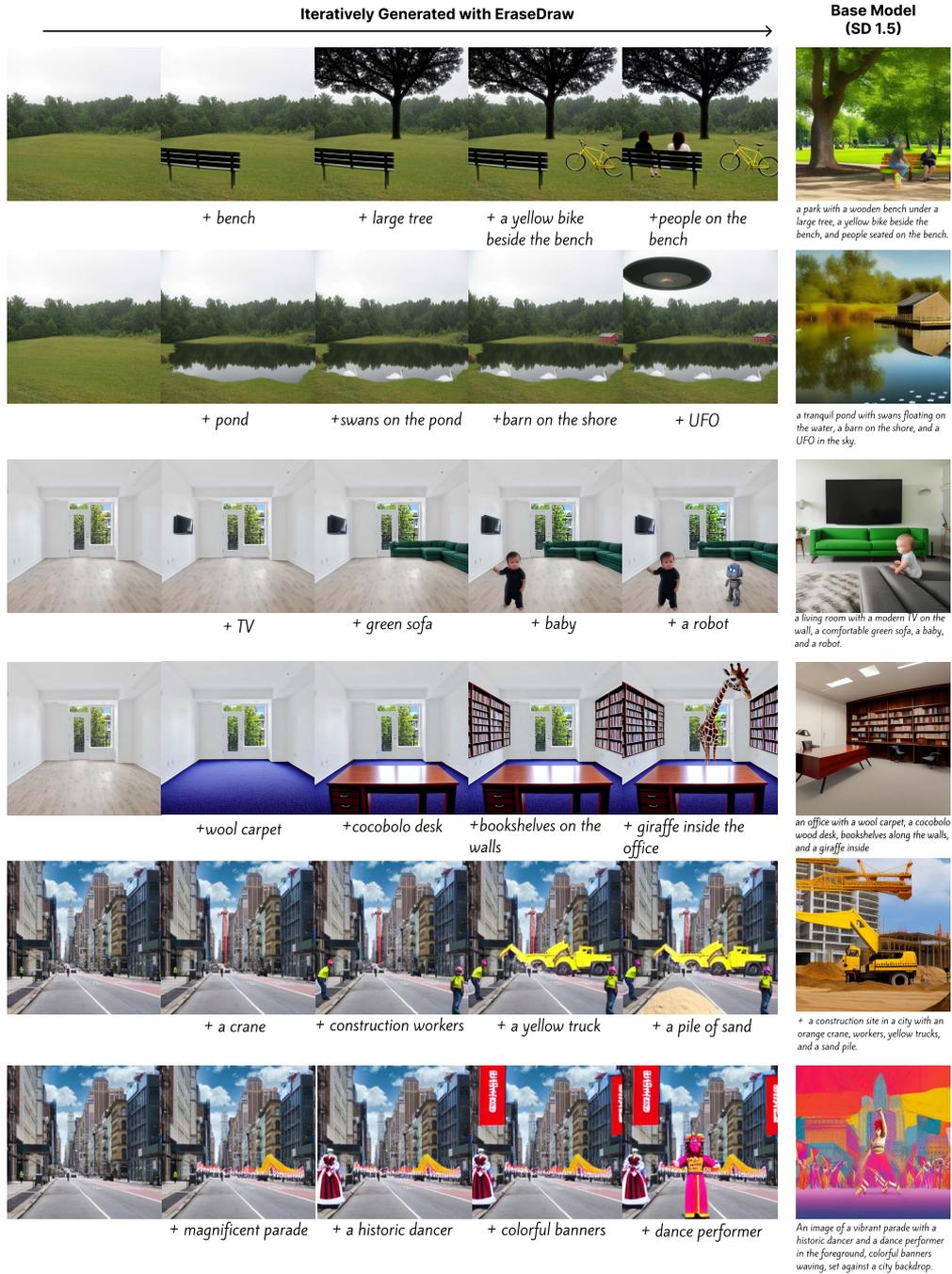


Fig. 9: Iterative Generation

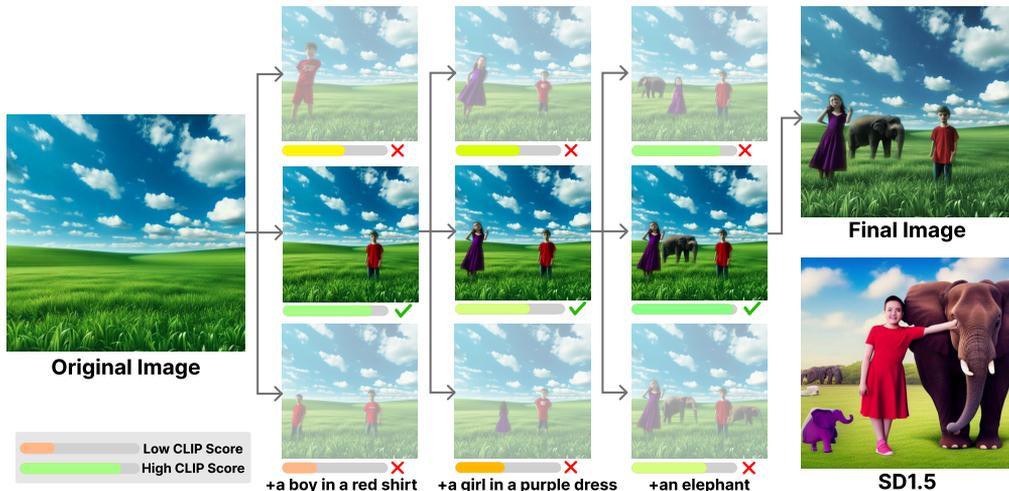


Fig. 10: Beam Search Given the original image on the left, we run beam search using CLIP distance between the edit instruction and the output image as the score heuristic. In this figure, we display the results from the top beam. To compare with one-shot diffusion using SD1.5, we use the prompt "A boy in a red shirt, girl in a purple dress, and an elephant on a grassy field."

an image using object insertion steps. We emphasize that our approach is model-agnostic; all general-purpose image editing models can benefit from including object insertion examples from EraseDraw in their training data. Besides image generation and editing, EraseDraw can unlock downstream applications such as giving robots a commonsense understanding of object placements and augmenting existing image datasets with novel objects.

Our work is limited by the relatively small created dataset and the base model. In particular, we note that the quality of the inserted objects is limited by the capabilities of the base model we finetune from, consequently our model is unable to generate anatomically accurate hands or faces, and it has limited ability to follow instructions involving precise object placements. Furthermore, while our framework can be used for processing arbitrary unlabeled images, our dataset is derived from OpenImages and COCO, which may have biases for certain object categories.

Dataset	Dataset Size	Autonomously Generated	$\text{CLIP}_{\text{out}} \uparrow$	$\text{CLIP}_{\text{dir}} \uparrow$
GQA-Insert	50,000	No	0.0907	0.1279
EraseDraw-15k	15,000	Yes	0.0905	0.1254
Mixed	65,000	No	0.0930	0.1383

Table 2: EraseDraw Dataset Ablation Studies

We believe scaling up our autonomous data pipeline and fine-tuning stronger base models will yield superior performance.

4.2 Societal Impact

Image editing tools such EraseDraw could potentially be misused for creating misleading or harmful content. It is crucial for the research community to develop techniques for mitigating misuse, and to promote responsible use of these technologies.

Acknowledgements: The authors would like to thank for Huy Ha, Samir Gadre, and Zeyi Liu valuable feedback and discussions.

References

- [1] Amit, T., Shaharbany, T., Nachmani, E., Wolf, L.: Segdiff: Image segmentation with diffusion probabilistic models. arXiv preprint arXiv:2112.00390 (2021)
- [2] Avrahami, O., Lischinski, D., Fried, O.: Blended diffusion for text-driven editing of natural images. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18208–18218 (2022)
- [3] Avrahami, O., Tamir, G., Benaim, S., Fried, O.: Emu edit: Multi-task instruction-guided image editing with learned task embeddings. arXiv preprint arXiv:2304.04384 (2023)
- [4] Azadi, S., Pathak, D., Ebrahimi, S., Darrell, T.: Compositional gan: Learning image-conditional binary composition. In: International Journal of Computer Vision. vol. 128, pp. 2629–2642. Springer (2020)
- [5] Bar-Tal, O., Ofri-Amar, D., Fridman, R., Kasten, Y., Ullman, S., Dekel, T.: Text2live: Text-driven layered image and video editing. arXiv preprint arXiv:2204.02491 (2022)
- [6] Blattmann, A., Dockhorn, T., Kulal, S., Mendelevitch, D., Kilian, M., Lorenz, D., Levi, Y., English, Z., Voleti, V., Letts, A., et al.: Stable video diffusion: Scaling latent video diffusion models to large datasets. arXiv preprint arXiv:2311.15127 (2023)
- [7] Brooks, T., Holynski, A., Efros, A.A.: Instructpix2pix: Learning to follow image editing instructions. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18392–18402 (2023)
- [8] Brooks, T., Holynski, A., Efros, A.A.: Instructpix2pix: Learning to follow image editing instructions. In: CVPR (2023)
- [9] Brooks, T., Peebles, B., Holmes, C., DePue, W., Guo, Y., Jing, L., Schnurr, D., Taylor, J., Luhman, T., Luhman, E., Ng, C., Wang, R., Ramesh, A.: Video generation models as world simulators (2024), <https://openai.com/research/video-generation-models-as-world-simulators>
- [10] Chen, S., Sun, P., Song, Y., Luo, P.: Diffusiondet: Diffusion model for object detection. arXiv preprint arXiv:2211.09788 (2022)
- [11] Choi, W., Chao, Y.W., Pantofaru, C., Savarese, S.: Context-driven 3d scene understanding from a single image. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 2688–2695 (2012)

- [12] Deitke, M., Liu, R., Wallingford, M., Ngo, H., Michel, O., Kusupati, A., Fan, A., Laforte, C., Voleti, V., Gadre, S.Y., et al.: Objaverse-xl: A universe of 10m+ 3d objects. *Advances in Neural Information Processing Systems* **36** (2024)
- [13] Dhariwal, P., Nichol, A.: Diffusion models beat gans on image synthesis. *NeurIPS* (2021)
- [14] Fisher, M., Ritchie, D., Savva, M., Funkhouser, T., Hanrahan, P.: Example-based synthesis of 3d object arrangements. *ACM Transactions on Graphics (TOG)* **31**(6), 1–11 (2012)
- [15] Gal, R., Alaluf, Y., Atzmon, Y., Patashnik, O., Bermano, A.H., Chechik, G., Cohen-Or, D.: An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618* (2022)
- [16] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. *NeurIPS* (2014)
- [17] Hertz, A., Mokady, R., Tenenbaum, J., Aberman, K., Pritch, Y., Cohen-Or, D.: Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626* (2022)
- [18] Hertz, A., Mokady, R., Tenenbaum, J., Aberman, K., Pritch, Y., Cohen-Or, D.: Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626* (2022)
- [19] Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. *Advances in neural information processing systems* **33**, 6840–6851 (2020)
- [20] Huang, X., Belongie, S.: Arbitrary style transfer in real-time with adaptive instance normalization. In: *Proceedings of the IEEE international conference on computer vision*. pp. 1501–1510 (2017)
- [21] Hudson, D.A., Manning, C.D.: Gqa: A new dataset for real-world visual reasoning and compositional question answering. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. pp. 6700–6709 (2019)
- [22] Isola, P., Zhu, J.Y., Zhou, T., Efros, A.A.: Image-to-image translation with conditional adversarial networks. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 1125–1134 (2017)
- [23] Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. pp. 4401–4410 (2019)
- [24] Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., Aila, T.: Analyzing and improving the image quality of stylegan. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. pp. 8110–8119 (2020)
- [25] Kawar, B., Lang, O., Tov, O., Irani, M.: Imagen editor: Text-based selection and editing of images. *arXiv preprint arXiv:2211.15481* (2022)
- [26] Kawar, B., Zada, S., Lang, O., Tov, O., Chang, H., Dekel, T., Mosseri, I., Irani, M.: Imagic: Text-based real image editing with diffusion models. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 18215–18224 (2022)
- [27] Kingma, D.P., Welling, M.: Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* (2013)
- [28] Lee, J.Y., Tseng, Z., Abbeel, P.: Relaxed placement: Learning to synthesize compositional scene layouts with object relations. In: *Computer Vision and Pattern Recognition (CVPR)* (2022)

- [29] Li, A.C., Prabhudesai, M., Duggal, S., Brown, E., Pathak, D.: Your diffusion model is secretly a zero-shot classifier. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 2206–2217 (October 2023)
- [30] Lin, D., Fidler, S., Urtasun, R.: Holistic scene understanding for 3d object detection with rgb-d cameras. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1417–1424 (2013)
- [31] Liu, R., Wu, R., Van Hoorick, B., Tokmakov, P., Zakharov, S., Vondrick, C.: Zero-1-to-3: Zero-shot one image to 3d object. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 9298–9309 (2023)
- [32] Liu, Y., Lin, C., Zeng, Z., Long, X., Liu, L., Komura, T., Wang, W.: Syncdreamer: Learning to generate multiview-consistent images from a single-view image. arXiv preprint arXiv:2309.03453 (2023)
- [33] Meng, C., Song, Y., Song, J., Wu, J., Zhu, J.Y., Ermon, S.: Sdedit: Image synthesis and editing with stochastic differential equations. In: International Conference on Learning Representations (2021)
- [34] Michel, O., Bhattad, A., VanderBilt, E., Krishna, R., Kembhavi, A., Gupta, T.: Object 3dit: Language-guided 3d-aware image editing. *Advances in Neural Information Processing Systems* **36** (2024)
- [35] Mokady, R., Hertz, A., Aberman, K., Pritch, Y., Cohen-Or, D.: Null-text inversion for editing real images using guided diffusion models. arXiv preprint arXiv:2211.09794 (2022)
- [36] Nichol, A., Dhariwal, P., Ramesh, A., Shyam, P., Mishkin, P., McGrew, B., Sutskever, I., Chen, M.: Glide: Towards photorealistic image generation and editing with text-guided diffusion models. arXiv preprint arXiv:2112.10741 (2021)
- [37] Nichol, A.Q., Dhariwal, P.: Improved denoising diffusion probabilistic models. In: International Conference on Machine Learning. pp. 8162–8171. PMLR (2021)
- [38] Ozguroglu, E., Liu, R., Surís, D., Chen, D., Dave, A., Tokmakov, P., Vondrick, C.: pix2gestalt: Amodal segmentation by synthesizing wholes. In: CVPR (2024)
- [39] Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., Chen, M.: Hierarchical text-conditional image generation with clip latents. arXiv preprint arXiv:2204.06125 (2022)
- [40] Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: CVPR (2022)
- [41] Ruiz, N., Li, Y., Jampani, V., Pritch, Y., Rubinstein, M., Aberman, K.: Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In: CVPR (2023)
- [42] Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E., Ghasemipour, S.K.S., Ayan, B.K., Mahdavi, S.S., Lopes, R.G., et al.: Photorealistic text-to-image diffusion models with deep language understanding. arXiv preprint arXiv:2205.11487 (2022)
- [43] Sang, Y., Park, J., Moon, S.G., Mu, S., Mu, Q., Guo, J., Zhang, Y., Lee, S., Kim, M., Lee, J., et al.: Image sculpting: Interactive image editing using 3d geometric operations. arXiv preprint arXiv:2303.13786 (2023)
- [44] Sargent, K., Li, Z., Shah, T., Herrmann, C., Yu, H.X., Zhang, Y., Chan, E.R., Lagun, D., Fei-Fei, L., Sun, D., Wu, J.: ZeroNVS: Zero-shot 360-degree view synthesis from a single real image. arXiv preprint arXiv:2310.17994 (2023)

- [45] Schuhmann, C., Beaumont, R., Vencu, R., Gordon, C., Wightman, R., Cherti, M., Coombes, T., Katta, A., Mullis, C., Wortsman, M., et al.: Laion-5B: An open large-scale dataset for training next generation image-text models. *NeurIPS* (2022)
- [46] Shen, Y., Gu, J., Tang, X., Zhou, B.: Interpreting the latent space of gans for semantic face editing. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 9243–9252 (2020)
- [47] Sheynin, S., Polyak, A., Singer, U., Kirstain, Y., Zohar, A., Ashual, O., Parikh, D., Taigman, Y.: Emu edit: Precise image editing via recognition and generation tasks. *arXiv preprint arXiv:2311.10089* (2023)
- [48] Su, D., Yu, C., Frank, B., Guibas, L., Welling, M., Tu, Z.: Edict: Exact diffusion inversion via coupled transformations. *arXiv preprint arXiv:2211.12446* (2022)
- [49] Su, Z., Chou, B., Yang, B., Zhang, R., Qiao, Y., Xu, J., Wang, Z., Cheng, X., Chen, X., Wu, H., et al.: Magicbrush: Text-to-image editing with a human in the loop. *arXiv preprint arXiv:2302.04754* (2023)
- [50] Voynov, A., Babenko, A.: Unsupervised discovery of interpretable directions in the gan latent space. In: *International Conference on Machine Learning*. pp. 9786–9796. PMLR (2020)
- [51] Wang, T., Zhang, T., Yang, B., Lu, H., Dai, D., Cai, J., Zhang, D., Van Gool, L.: Pretraining is all you need for image-to-image translation. In: *European conference on computer vision*. pp. 30–48. Springer (2022)
- [52] Wu, R., Liu, R., Vondrick, C., Zheng, C.: Sin3dm: Learning a diffusion model from a single 3d textured shape. *arXiv preprint arXiv:2305.15399* (2023)
- [53] Xu, J., Liu, S., Vahdat, A., Byeon, W., Wang, X., De Mello, S.: Open-Vocabulary Panoptic Segmentation with Text-to-Image Diffusion Models. *arXiv preprint arXiv:2303.04803* (2023)
- [54] Zhan, G., Zheng, C., Xie, W., Zisserman, A.: Amodal ground truth and completion in the wild. *CVPR* (2024)
- [55] Zhang, C., Wang, L., Yang, R.: Scene design by integrating geometry and physics for realistic image synthesis. In: *Computer Graphics Forum*. vol. 33, pp. 61–70. Wiley Online Library (2014)
- [56] Zhao, W.H., Jiang, J., Weng, J., He, J., Lim, E.P., Yan, H., Li, X.: Image-based contextual advertisement recommendation. In: *Proceedings of the 34th international ACM SIGIR conference on Research and development in Information Retrieval*. pp. 821–830 (2011)
- [57] Zhu, J.Y., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: *Proceedings of the IEEE international conference on computer vision*. pp. 2223–2232 (2017)
- [58] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [59] Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18208–18218, 2022.
- [60] Samaneh Azadi, Deepak Pathak, Sayna Ebrahimi, and Trevor Darrell. Compositional gan: Learning image-conditional binary composition. In *International Journal of Computer Vision*, pages 2629–2642. Springer, 2020.

- [61] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18392–18402, 2023.
- [62] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6700–6709, 2019.
- [63] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *Advances in Neural Information Processing Systems*, 35:26565–26577, 2022.
- [64] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18215–18224, 2022.
- [65] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023.
- [66] Joon-Young Lee, Zach Tseng, and Pieter Abbeel. Relaxed placement: Learning to synthesize compositional scene layouts with object relations. In *Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [67] Simian Luo, Yiqin Tan, Suraj Patil, Daniel Gu, Patrick von Platen, Apolinário Passos, Longbo Huang, Jian Li, and Hang Zhao. Lcm-lora: A universal stable-diffusion acceleration module. *arXiv preprint arXiv:2311.05556*, 2023.
- [68] Chenlin Meng, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Image synthesis and editing with stochastic differential equations. In *International Conference on Learning Representations*, 2021.
- [69] Oscar Michel, Anand Bhattad, Eli VanderBilt, Ranjay Krishna, Aniruddha Kembhavi, and Tanmay Gupta. Object 3dit: Language-guided 3d-aware image editing. *Advances in Neural Information Processing Systems*, 36, 2024.
- [70] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- [71]] Hive-magicbrush model checkpoint. <https://huggingface.co/osunlp/HIVE-MagicBrush/resolve/main/MagicBrush-epoch-000130.ckpt>, 2023. Accessed: 2024-03-21.
- [72] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- [73] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022.
- [74] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed K Sajjadi Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022.

- [75] Yuyang Sang, Junwoo Park, Seong-Gyun Moon, Shan Mu, Qi Mu, Jia Guo, Yahan Zhang, Seunghwan Lee, Minhyuk Kim, Jaegul Lee, et al. Image sculpting: Interactive image editing using 3d geometric operations. *arXiv preprint arXiv:2303.13786*, 2023.
- [76] Shelly Sheynin, Adam Polyak, Uriel Singer, Yuval Kirstain, Amit Zohar, Oron Ashual, Devi Parikh, and Yaniv Taigman. Emu edit: Precise image editing via recognition and generation tasks. *arXiv preprint arXiv:2311.10089*, 2023.
- [77] Daniel Su, Chenglin Yu, Berthy Frank, Leonidas Guibas, Max Welling, and Zhiwei Tu. Edict: Exact diffusion inversion via coupled transformations. *arXiv preprint arXiv:2211.12446*, 2022.
- [78] Zhengyi Su, Brian Chou, Bowen Yang, Richard Zhang, Yinghao Qiao, Jiannan Xu, Ziwei Wang, Xu Cheng, Xiaohui Chen, Hao Wu, et al. Magicbrush: Text-to-image editing with a human in the loop. *arXiv preprint arXiv:2302.04754*, 2023.
- [79] 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 *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 2149–2159, 2022.
- [80] Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. Cogvlm: Visual expert for pretrained language models. *arXiv preprint arXiv:2311.03079*, 2023.
- [81] Ahmet Burak Yildirim, Vedat Baday, Erkut Erdem, Aykut Erdem, and Aysegul Dundar. Inst-inpaint: Instructing to remove objects with diffusion models. *arXiv preprint arXiv:2304.03246*, 2023.
- [82] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11975–11986, 2023.

EraseDraw: Learning to Insert Objects by Erasing Them from Images

Supplementary Material

A Training

We train our model for 1000 steps on $2 \times 48\text{GB}$ A6000 GPUs over 12 hours at 256×256 resolution. With the help of gradient accumulation, we train at an effective batch size of 1536 images. We adopt all other model related training hyperparameters from [61]. We find that our model can generalize to non-square images and resolutions as high as 512×512

B Accelerating Inference

While the inference duration of diffusion models is prohibitive, we accelerate our model using a pre-trained LCM-LoRA [67] for Stable Diffusion 1.5. Surprisingly, despite being trained for the text-to-image task, the LCM-LoRA preserves the insertion abilities of our model while cutting down the number of inference steps to as low as 4 to 16 steps depending on the desired image quality.

All results in our paper except our quantitative evaluations and qualitative comparison figures have been generated using LCM LoRA with 16 steps using 512×512 resolution.

Our qualitative comparison figures are generated with 50 steps of Euler ancestral sampling outlined in Karras et al. [63] with 512×512 resolution.

C Prompts

Extracting Captions from Images using Large Vision Language Model To extract captions from an image using GPT-4 Vision, we use the following prompt

- 1 1. Describe the objects in this image along with their attributes and their spatial relations with respect to the other objects.
- 2 2. For every individual object,
- 3 a) Come up with a "subject identification" for that object. The subject identification should be a simple way to identify the object in the image. Color and shape may be helpful to include.
- 4 Examples:
- 5 If two fish are in the picture, one is red and the other is blue, then you could identify their subject identificaitons as "the red fish" and "the blue fish".

```

6       If three men are in the picture, and they're spread out
       into a row, then you could identify their subject
       identifications as "man on the left","man on the right",man in
       the middle".
7       If only a single cat is in the picture, then the subject
       identification "cat" is sufficient.
8       The subject identification MUST NOT include a noun other than
       the subject.
9       b) Come up with simple captions of the form [adjective] [
       subject] [prepositional phrase] that describe the location of
       the object with respect to other objects or the image.
10      (e.g. a man with a blue shirt standing in front of the wall,
       an elephant next to the tree, a bat held by a player, the dog
       on the right of the image, etc.).
11      Make sure that all of the captions refer to exact the same
       subject.
12 3. Exclude large background elements from your captions, such as
       the sky, the ground, the walls, etc.
13
14 Finally, return your final response as a JSON in the form
15 {
16     "[subject identification]":
17     [
18         "[caption 1]", "[caption 2]", "[caption 3]", ...
19     ],
20     ...
21 }
22
23 You may now begin!

```

Detecting Bounding Boxes To detect an object named [object] using CogVLM [80], we prompt the model with

```
1 Where is [object]?
```

D Additional Examples of Beam Search

In figures 11, 12, 13, 14 15, 16, we showcase results from employing the EraseDraw method combined with beam search for gradual image composition. The complete image prompt, detailed in each figure’s caption, is divided into five sub-prompts, corresponding to five steps in the beam search process. At each step t , images are ranked by their CLIP similarity to the combined sub-prompts from 0 through t .

Each figure’s top six rows display images with the highest and lowest four CLIP scores. The last row shows a comparison with images generated by Stable Diffusion 1.5. To ensure fairness in terms of computation budget, we allow Stable Diffusion to generate as many images as beam search at each step, and apply CLIP filtering to get the best image. More specifically, for each step t in the process, Stable Diffusion 1.5

generates $(\text{beam width}) \times (\text{branching factor}) \times t$ images. These images are generated using the cumulative prompt, which is the combination of sub-prompts from 0 through t . We then select the image with the highest CLIP score relative to the cumulative prompt for presentation.

We compute all CLIP scores using siglip-base-patch16-224 [82] model from Huggingface.

E Dataset Generation Algorithm

We outline our precise dataset generation procedure below in Algorithm 1.

Algorithm 1 Autonomous Data Generation

- 1: **Input:** Input dataset D of unlabeled images
 - 2: **Output:** Dataset D' of $(C_{Image}, C_{Text}, \mathbf{x})$ tuples.
 - 3: Initialize $D' \leftarrow \emptyset$
 - 4: **for** each image \mathbf{x} in D **do**
 - 5: Use image captioner to propose objects \mathbf{o}_i , simple captions T_i for each object, and a set of complex captions $\{T_1^{(i)} \dots, T_K^{(i)}\}$ for each object
 - 6: $\mathbf{b}_i = \text{DetectBoundingBox}(\mathbf{x}, T_i)$
 - 7: $\mathbf{m}_i = \text{BoundingBoxToSegmentation}(\mathbf{x}, \mathbf{b}_i)$
 - 8: **for** each complex caption $T_j^{(i)}$ **do**
 - 9: $C_{Image} = \text{Erase}(\mathbf{x}, \mathbf{m}_i)$
 - 10: append $(C_{Image}, C_{Text} = T_j^{(i)}, \mathbf{x})$ to D'
 - 11: **end for**
 - 12: **end for**
-

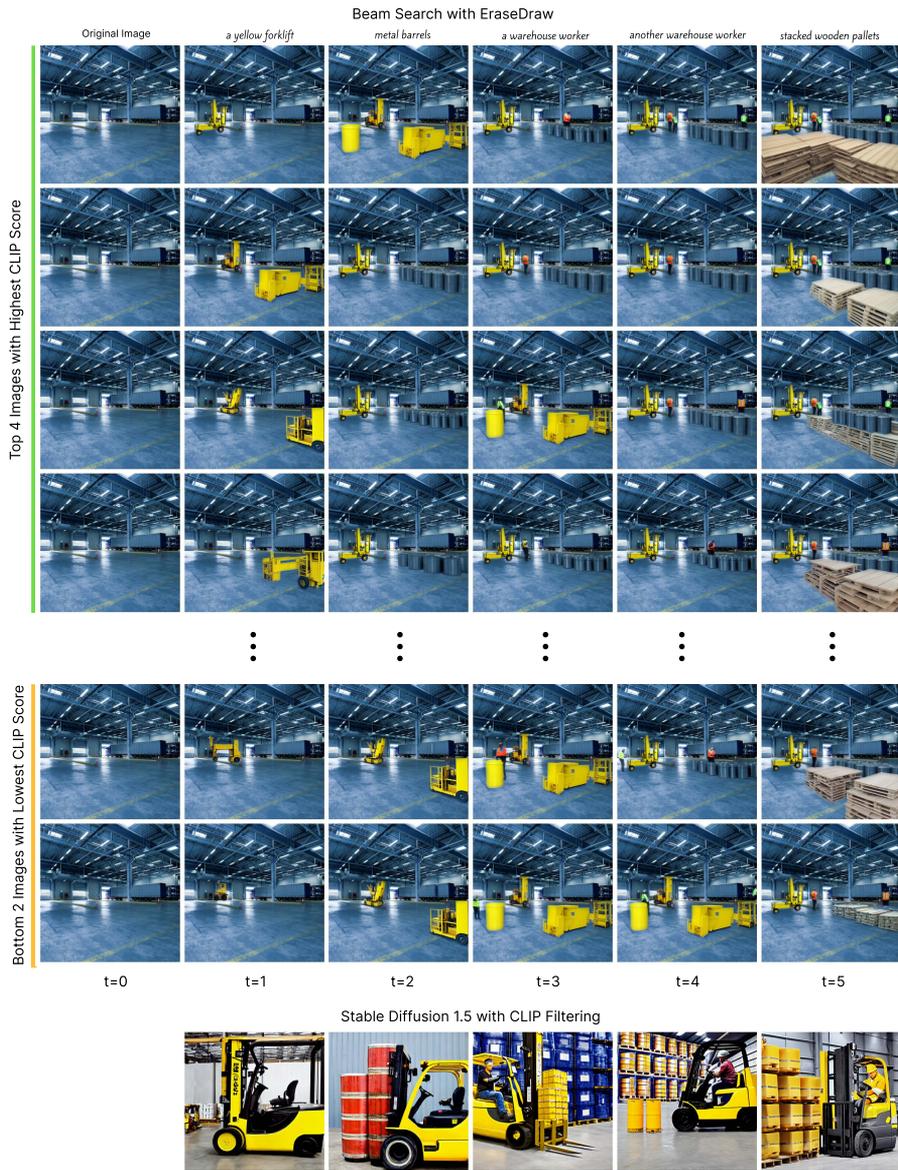


Fig. 11: Beam search with beam width $k = 3$ and branching factor $N = 4$ on the prompt "a yellow forklift, metal barrels, a warehouse worker, another warehouse worker, stacked wooden pallets"



Fig. 12: Beam search with beam width $k = 5$ and branching factor $N = 4$ on the prompt "a bright green potted plant, a sofa, a bookshelf, an elegant rug on the floor, a painting on the wall"

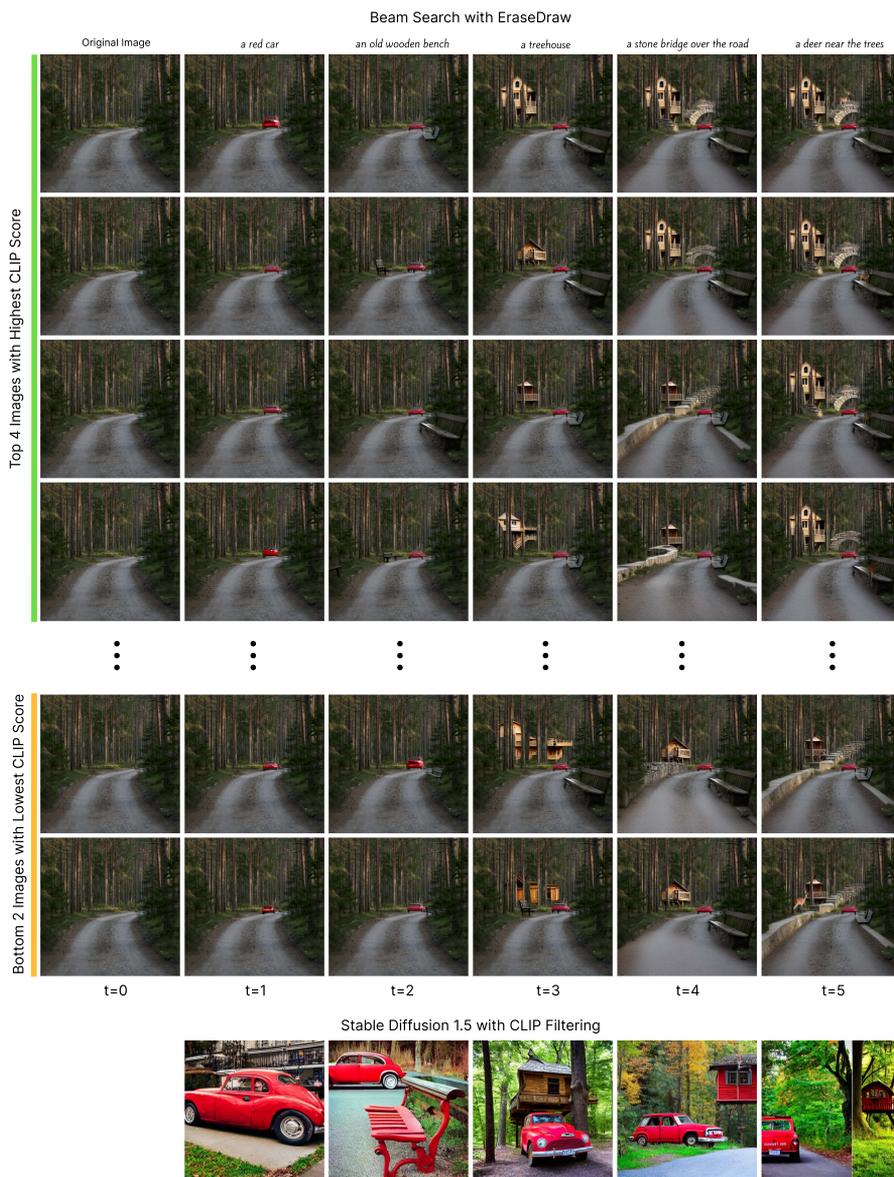


Fig. 13: Beam search with beam width $k = 3$ and branching factor $N = 4$ on the prompt "a red car, an old wooden bench, a tree house, a stone bridge over the road, a deer near the trees"

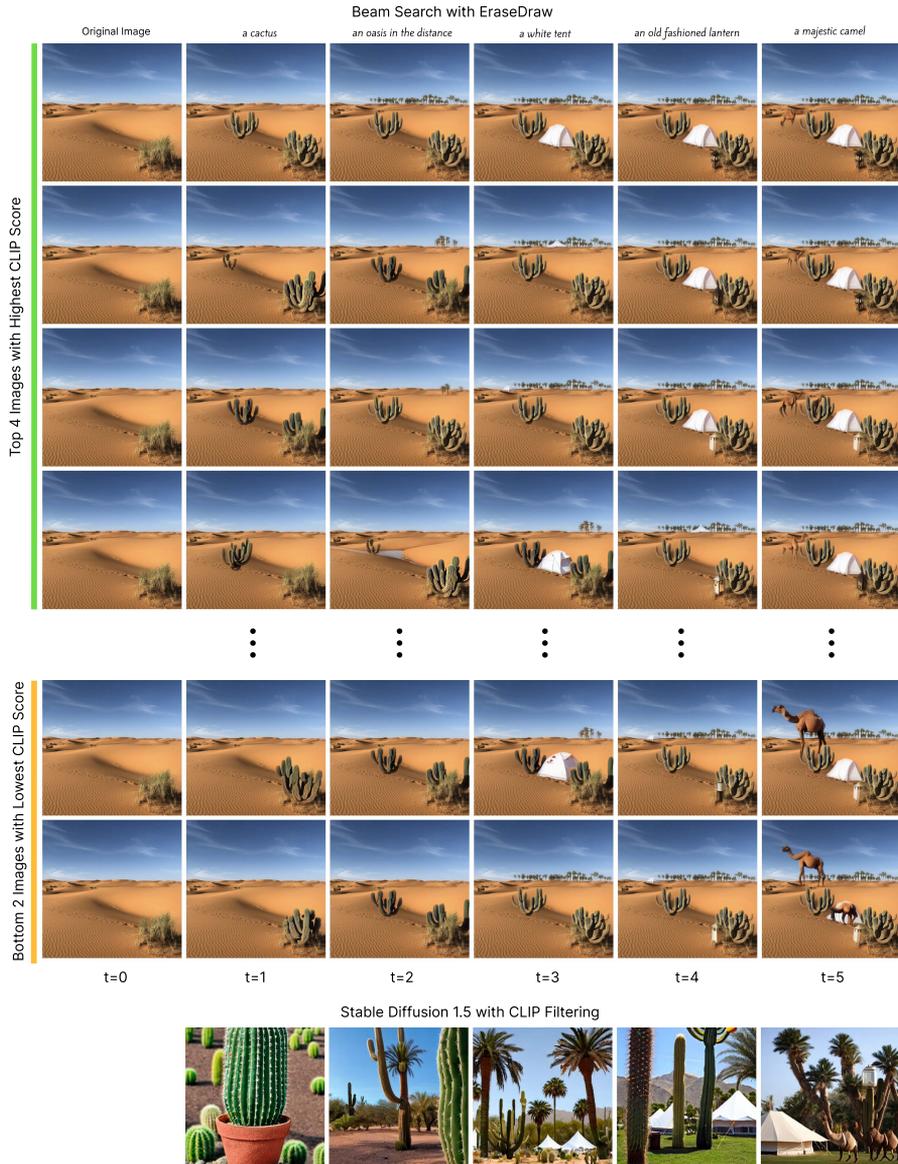


Fig. 14: Beam search with beam width $k = 3$ and branching factor $N = 4$ on the prompt "a cactus, an oasis in the distance, a white ten, an old fashioned lantern, a majestic camel"

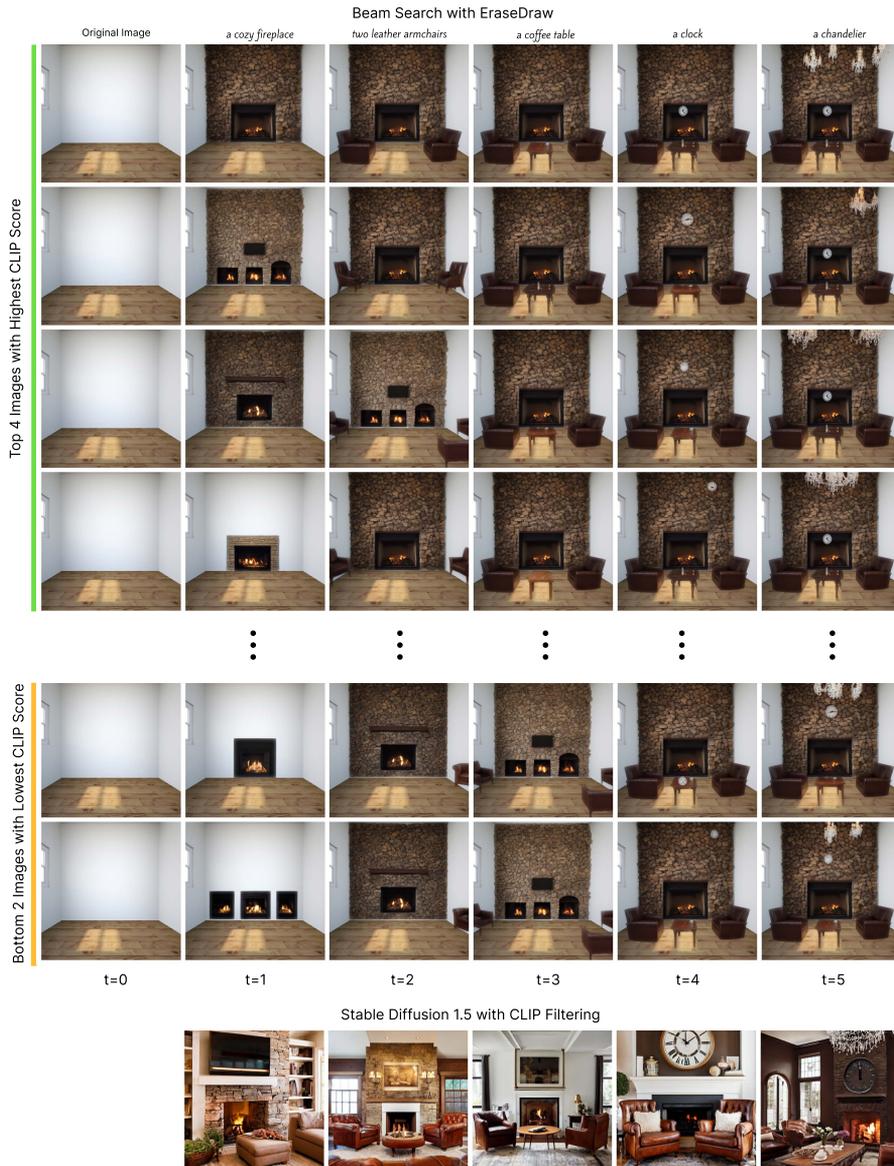


Fig. 15: Beam search with beam width $k = 3$ and branching factor $N = 4$ on the prompt "a cozy fireplace, two leather armchairs, a coffee table, a clock, a chandelier"

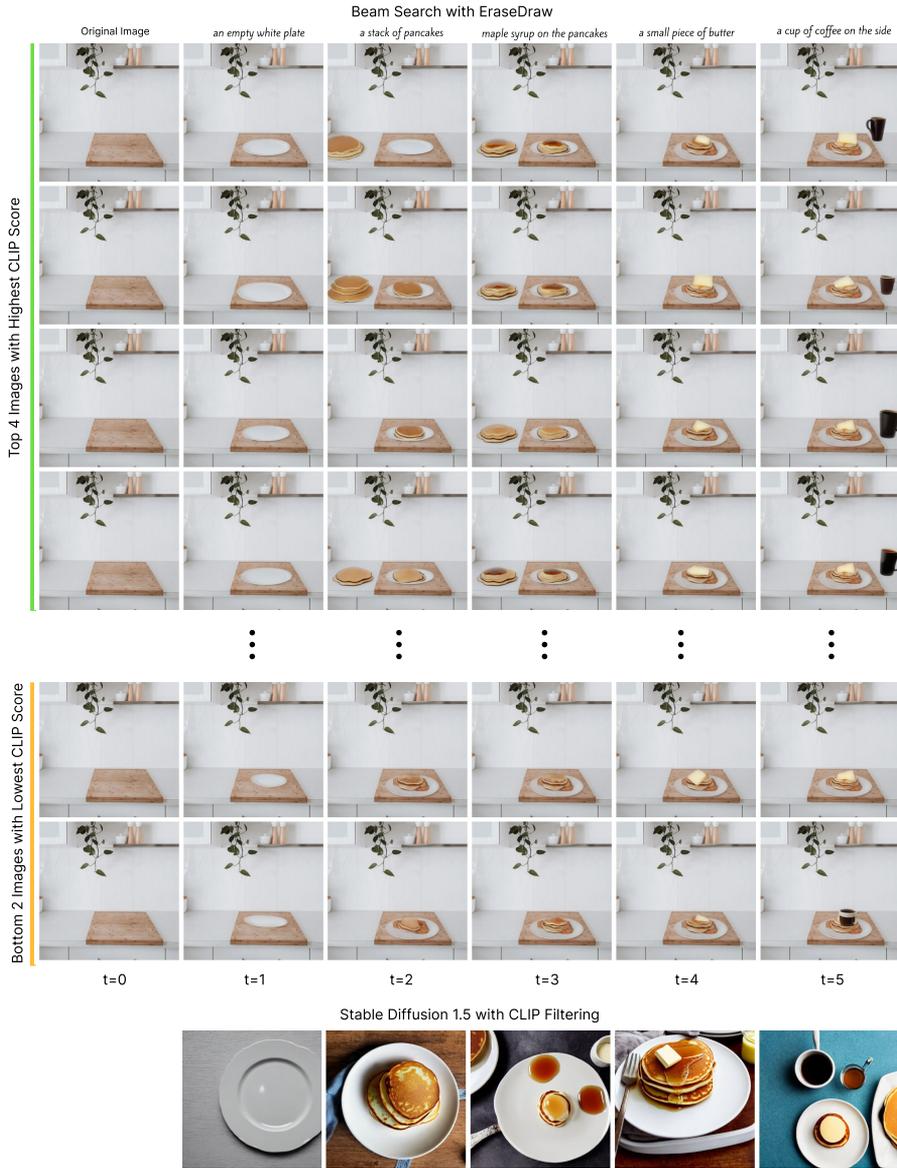


Fig. 16: Beam search with beam width $k = 3$ and branching factor $N = 4$ on the prompt "an empty white plate, a stack of pancakes, maple syrup on the pancakes, a small piece of butter, a cup of coffee on the side"



Fig. 17: EraseDraw's results on insertion tasks given in Figure 2 of the paper