# SPICE: A Synergistic, Precise, Iterative, and Customizable Image Editing Workflow

Kenan Tang [* 1]    Yanhong Li [* 2]    Yao Qin [1]

## Abstract

Recent prompt-based image editing models have demonstrated impressive prompt-following capability at structural editing tasks. However, existing models still fail to perform local edits, follow detailed editing prompts, or maintain global image quality beyond a single editing step. To address these challenges, we introduce SPICE, a *training-free* workflow that accepts arbitrary resolutions and aspect ratios, accurately follows user requirements, and improves image quality consistently during more than 100 editing steps. By synergizing the strengths of a base diffusion model and a Canny edge ControlNet model, SPICE robustly handles free-form editing instructions from the user. SPICE outperforms state-of-the-art baselines on a challenging realistic image-editing dataset consisting of semantic editing (object addition, removal, replacement, and background change), stylistic editing (texture changes), and structural editing (action change) tasks. Not only does SPICE achieve the highest quantitative performance according to standard evaluation metrics, but it is also consistently preferred by users over existing image-editing methods. We release the workflow implementation for popular diffusion model Web UIs to support further research and artistic exploration[1].

## 1. Introduction

Image editing is the task of changing the content of an image according to a user's requirements. An example is adding an apple of a specific size at a specific location on an image

---
[*]Equal contribution  [1]University of California, Santa Barbara, Santa Barbara, CA, United States [2]University of Chicago, Chicago, IL, United States. Correspondence to: Yao Qin <yao-qin@ucsb.edu>.

[1]https://github.com/kenantang/spice

(Figure 1)[2]. A powerful image editing tool is vital for many applications from creative design to scientific research. In artistic workflows, photographers modify or remove objects in an image to enhance storytelling and visual impact. In machine learning research, image editing has been used in video editing (Fan et al., 2025), data augmentation (Hirota et al., 2024), and benchmark construction (Wu et al., 2024).

Existing vision language models have achieved initial success on the image editing task (Brooks et al., 2023; Zhang et al., 2024). Such models usually take in an original image and user's editing prompt as the input and output the edited image. In addition, some methods use a binary mask to improve the editing results (Wang et al., 2023; Zhao et al., 2024). However, for advanced artistic or research purposes that require more than one editing step, existing methods are disqualified by the following limitations. First, pixels outside the mask often deteriorate after editing. Second, the user cannot specify the precise size and location of an object by the mask. Third, models struggle with unusual editing tasks, such as adding a backpack to a bench facing away from the viewer (Figure 3). Taken together, these limitations make it extremely difficult for users to refine an image across multiple edits, as image degradation inevitably accumulates during iterative editing.

To overcome these limitations, we propose **SPICE**, a *training-free* workflow that can be easily adopted in existing text-to-image diffusion models. Our workflow follows a structured three-step process: (1) mask generation, (2) color and edge hint generation, and (3) two-stage denoising, as shown in Figure 2. First, in the mask generation step, the user provides a mask to define the editing region and the context size, localizing the modifications while using essential contextual information from the original image. By constraining the editing region to the user-provided mask, SPICE strictly preserves image quality and avoid deterioration outside the mask.

Second, during color and edge hint generation, SPICE allows the user to provide arbitrarily detailed or simplified image-space information to the model. This goes beyond

---
[2]Figures in this paper are provided in high resolution. Readers are encouraged to zoom in to examine the details.

*Figure 1.* **SPICE enables a user to edit the image exactly as they want, and image details improve instead of deteriorating after many editing steps.** The first row shows the full image of a 3000×2000 resolution. The second row shows a 900×600 region enlarged for better visibility. In this example, a user uses 9 editing steps to perform various edits, including structure change, object removal, object addition, object replacement, text addition, color change, and detail fixes. Steps 7, 8, and 9 together fix the fridge structure. The labels above the first row are not the true editing instructions but their abbreviated version. For example, in the "Add a Word" column, the user wants to add the specific word "suspicious" to the white bowl. The prompt is "An open fridge with food in it. A bowl with a word 'suspicious' on it." The result aligns with the user's requirement.

conventional image editing methods (Brooks et al., 2023; Zhang et al., 2024) that rely solely on textual prompts which can hardly provide the precise size and location according to the user's requirement. In contrast to existing methods, SPICE integrates structured visual guidance from a hinted image, offering a level of control over size and positioning that exceeds what textual prompts can achieve.

Finally, SPICE uses a two-stage denoising process, where a Canny edge ControlNet model (Zhang et al., 2023) in the *early* denoising steps integrates image-space hints, and a base diffusion model refines the output in the *later* steps to generate diverse variations, as shown in Figure 2(c). By synergizing the complementary strengths of both models, the two-stage denoising process enables *precise* and *customizable* editing, in which the model faithfully follows the user's requirements on the location, size, and other properties of the edited object. This not only enhances the quality of photo-realistic images in complex tasks but also allows AI-generated artwork to move beyond stereotypical and fundamental limitations, such as predominantly portraying a single character or repeatedly erring on details.

By overcoming the limitations, SPICE consistently achieves success in the iterative editing tasks of *more than 100 editing steps*. The editing steps are a mixture of semantic and structural editing steps, where the image quality is consistently improved. For the first time, our workflow enables scaling test-time compute (Snell et al., 2024) in image generation, where the user can decide how each unit of additional compute time should contribute to the final image quality. This user-friendly design distinguishes our workflow from contemporary test-time scaling approaches

that rely on automatic search algorithms and rewards (Ma et al., 2025; Singhal et al., 2025), over which the user has no granular control.

Despite the strong capability, our workflow can be easily integrated into all popular diffusion model Web UIs, such as Stable Diffusion Web UI Automatic1111, ComfyUI, and Stable Diffusion Web UI Forge. Unlike existing tools such as ADetailer[3] or Regional Prompter[4], which offer an overwhelming number of hyperparameters that can be difficult to navigate, SPICE provides a smaller set of hyperparameters for ease-of-use. Furthermore, SPICE introduces minimal computational overhead in each editing step, making editing as efficient as generating an image from text with the same model and sampling hyperparameters. This extremely low overhead allows our workflow to run on consumer GPUs, unlocking high-quality and iterative image editing for more users. Combining strong capabilities with low implementation and computational cost, our workflow provides a powerful yet accessible image editing tool for researchers and non-researchers alike with huge convenience and flexibility.

## 2. Methods

SPICE is based on inpainting (Lugmayr et al., 2022), an operation that uses a diffusion model to replace a masked region on an image, conditioned on a prompt that describes the new image. Note that this prompt differs from the editing prompt, which uses a verb (e.g., add or remove) to describe
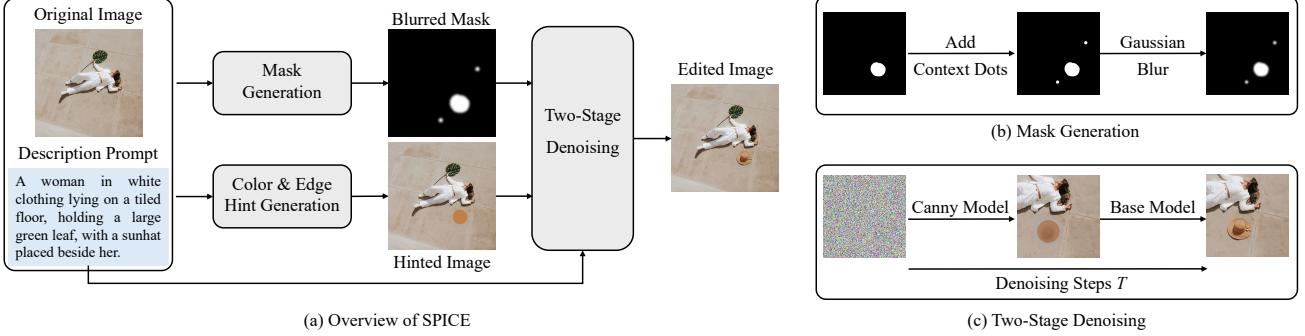
---

[3] https://github.com/Bing-su/adetailer
[4] https://github.com/hako-mikan/sd-webui-regional-prompter

(a) Overview of SPICE

(b) Mask Generation

(c) Two-Stage Denoising

*Figure 2.* **By sketching a binary mask with context dots and a color & edge hint, users can effortlessly achieve realistic edits with SPICE.** Subfigure (a) shows the overview of our workflow, while Subfigures (b) and (c) show the internal steps. In this example, the user requires a sunhat to be added next to the woman. First, the user sketches both a mask with context dots and a hinted image containing color and edge hints. The mask is automatically blurred after being sketched. Then, during the two-stage denoising step, the Canny and base models perform the early and late denoising steps, respectively. See Figure 5 for more examples of masks and hints.

the editing operation (Zhang et al., 2024). For example, when the editing prompt is "add an apple in the fridge", the description prompt will be "an apple in a fridge". In this section, we introduce the three key steps of our workflow, namely mask generation, color and edge hint generation and two-stage denoising. An overview is shown in Figure 2.

## 2.1. Mask Generation

**Context Selection.** We denote the original image to be edited as $I_T \in [0,1]^{H \times W \times C}$, where $T$ is the index of the editing step, $H$ and $W$ are the image height and width, and $C = 3$ represents the RGB color channels. Other than a prompt $p$, traditional inpainting methods[5] require the user to provide a binary mask $M \in \{0,1\}^{H \times W}$, where 1 indicates the region to be edited. Given the mask and the prompt, the inpainting operation uses the following steps to generate an edited image. First, a bounding box of the region to be edited is calculated, and the bounding box is extended either vertically or horizontally to ensure its aspect ratio matches a user-specified resolution supported by the diffusion model, such as $1216 \times 832$ (Podell et al., 2024). The user-specified resolution is usually larger than the extended bounding box. Next, pixels on $I$ and $M$ within the extended bounding box are upsampled to the user-specified resolution. Then, the diffusion model generates a new image from latent noise. The generation process is conditioned on the description prompt and the upsampled pixels from $I$ and $M$. Finally, the output is downsampled to the resolution of the extended bounding box, and the inpainted region on $I_T$ is replaced by the downsampled output, resulting in $I_{T+1}$.

However, since these traditional methods generate the inpainted region without contextual information outside the extended bounding box, the output often appears unnatural, with inconsistent lighting or color compared to the rest of

the image. A naive solution is the "whole image" mode, where the entire image is used as context, and the extended bounding box is not used. However, this introduces two major problems: (1) poor performance on small objects (e.g., deformed fingers or scrambled patterns), and (2) distortion caused by resizing the whole image to a model-supported resolution when aspect ratios of the two differ.

To address these issues, we introduce *context dots*, a pair of dots at opposite corners of the desired bounding box. Context dots ensure that the extended bounding box includes sufficient context for generation. The user directly adds context dots to the original mask, resulting in a context mask $M_{\text{context}} \in \{0,1\}^{H \times W}$, as shown in Figure 2(b). This simple enhancement offers three key advantages: (1) the user can exclude image areas that interfere with the inpainting process, (2) the user can specify a resolution between that of the inpainted region and that of the full image, balancing local details and global context (Section 3.4), and (3) context dots minimally affect surrounding pixels, limiting changes to only the desired editing region.

**Soft Inpainting.** On inpainted images, there is usually an unwanted sharp boundary between the inpainted region and the remaining parts of the image. In more severe cases, an inpainted object can be incomplete. The reason is that even with a context, diffusion models are not robust enough to generate pixels that seamlessly blend with existing ones. To mitigate these artifacts, we adopt Differential Diffusion (Levin & Fried, 2023), a method that allows the diffusion model to be conditioned by continuous mask values in [0, 1] during generation. We provide the continuous values as a soft mask $M_{\text{soft}} \in [0,1]^{H \times W}$ by applying a simple Gaussian blur to $M_{\text{context}}$, as shown in Figure 2(b). This procedure, together with thresholds for blending the original and inpainted image, is named as Soft Inpainting in popular Web UIs. We use Soft Inpainting when it is available in an Web UI or implement our own simplified version (without thresholds) when it is not available.

---

[5]https://github.com/AUTOMATIC1111/
stable-diffusion-webui

## 2.2. Color and Edge Hint Generation

In many existing image editing systems (Croitoru et al., 2023; Yang et al., 2023; Huang et al., 2024), users typically specify desired content through a textual prompt. However, words alone cannot easily describe certain details, such as asymmetric apparel designs or intricate color patterns. Instead, these subtle details can be effectively conveyed with an additional visual hint in the image space. To this end, SPICE allows a user to first create a rough sketch or color layout using standard editing software (e.g., Krita or Adobe Photoshop). The resultant image, denoted $I_{\text{hinted}} \in [0, 1]^{H \times W \times C}$, then replaces the original image as the inpainting input. To incorporate information from this hinted image, a denoising strength hyperparameter in [0, 1] specifies how much the original pixels on the image should be changed[6]. At a moderate denoising strength (e.g., 0.5 to 0.7), the inpainting model produces pixels that remain close to the user's sketched hints—preserving intended colors, shapes, or patterns—while still diverging enough from the sketch to form realistic objects.

Note that the user can provide any form of color and edge hints, including sketches, reference images pasted in a collage style, or even another region of the original image (using the Clone Stamp Tool in PhotoShop). This flexibility is an improvement over existing methods, such as Magic-Quill (Liu et al., 2024b) that uses downsampled $32 \times 32$ color blocks for guidance and thus loses high-frequency details. Hence, SPICE allows a user to perform image editing from any point on the human-model collaboration spectrum. At one end, the user fully edit the image by drawing out every detail, without the help of diffusion models. At the other end, the user fully delegates the model to edit the image. Usually, the user can get decent editing results by staying on the end where the user input is minimal (Section 3.1). We will also discuss how a small set of hyperparameters enable the user to move freely along this spectrum (Section 3.4).

## 2.3. Two-Stage Denoising

Image editing methods (Lugmayr et al., 2022) typically rely on a single diffusion model to perform all denoising steps when generating an image from latent noise. However, different diffusion models have complementary strengths. On the one hand, a general-purpose text-to-image base model, such as Flux.1 [dev][7], excels at generating rich variations in its output but can only be conditioned on textual prompts.

On the other hand, Flux.1 [dev] Canny[8], a model derived from Flux.1 [dev] that contains a Canny edge ControlNet (i.e., a Canny model), can be conditioned on Canny edge information (Zhang et al., 2023) but sacrifices variability.

To synergize the strengths, we propose a two-stage denoising process. Specifically, we use a Canny model $f_{\text{Canny}}$ during *early* denoising steps to incorporate image-space hints. Starting from latent noise $z_0$, the Canny model can condition its generation on the Canny edge information $E_{\text{hinted}}$ (extracted from the hinted image $I_{\text{hinted}}$) within the extended bounding box, so that the edge hints can be followed. After the early steps, we get an intermediate latent image

$$z_{\text{Canny}} = f_{\text{Canny}}(I_{\text{hinted}}, p, E_{\text{hinted}}, M_{\text{soft}}, z_0). \quad (1)$$

In the *late* denoising steps, we use a base model $f_{\text{Base}}$ to generate diverse content from $z_{\text{Canny}}$, achieving realism and sophistication despite the simplicity of color and edge hints. This can be formulated as

$$z_{\text{base}} = f_{\text{base}}(I_{\text{hinted}}, p, M_{\text{soft}}, z_{\text{Canny}}). \quad (2)$$

The final latent image $z_{\text{base}}$ will be decoded into the edited RGB image. This two-stage denoising process ensures that the denoising process benefits from the strengths of both models. Meanwhile, by simply adjusting the proportion of denoising steps assigned to each model, users can intuitively balance variability and controllability.

Empirically, we find that for a wide range of image editing tasks, Canny Edge ControlNet models outperforms other ControlNet variants (e.g., depth or pose) in the first denoising stage. Furthermore, Canny models are more widely available than other ControlNet models. Therefore, we only use Canny models but not other ControlNet models.

## 3. Experiments

In this section, we evaluate SPICE under various image-editing scenarios. First, we evaluate our workflow on a standard benchmark of single-step editing. Then, we systematically verify the effectiveness of the three features of SPICE, namely precise, iterative, and customizable editing.

### 3.1. Benchmark Results

**Evaluation Benchmark.** We use the second version of the EditEval (Huang et al., 2024) benchmark to evaluate our workflow and baseline methods. EditEval is a single-step editing benchmark consisting of original images, source captions, target captions, and editing prompts. The original images have resolutions ranging from $1863 \times 1863$ to $8742 \times 8742$. EditEval covers semantic editing (object addition, object removal, object replacement, and background

---

[6]The hyperparameter does not work as a simple linear blending coefficient at the post-processing stage, and its exact implementation differs for various algorithms. To avoid ambiguity, we refrain from providing a general equation here. Interested readers can refer to the source code of popular Web UIs for the equations.

[7]https://huggingface.co/black-forest-labs/FLUX.1-dev

[8]https://huggingface.co/black-forest-labs/FLUX.1-Canny-dev

**Object Addition:** Add a backpack placed on the bench.

**Object Replacement:** Replace the road sign with a mailbox.

**Object Removal:** Remove the four women.

**Background Change:** Change the riverside to a desert.

**Texture Change:** Turn the handbag into the glass.
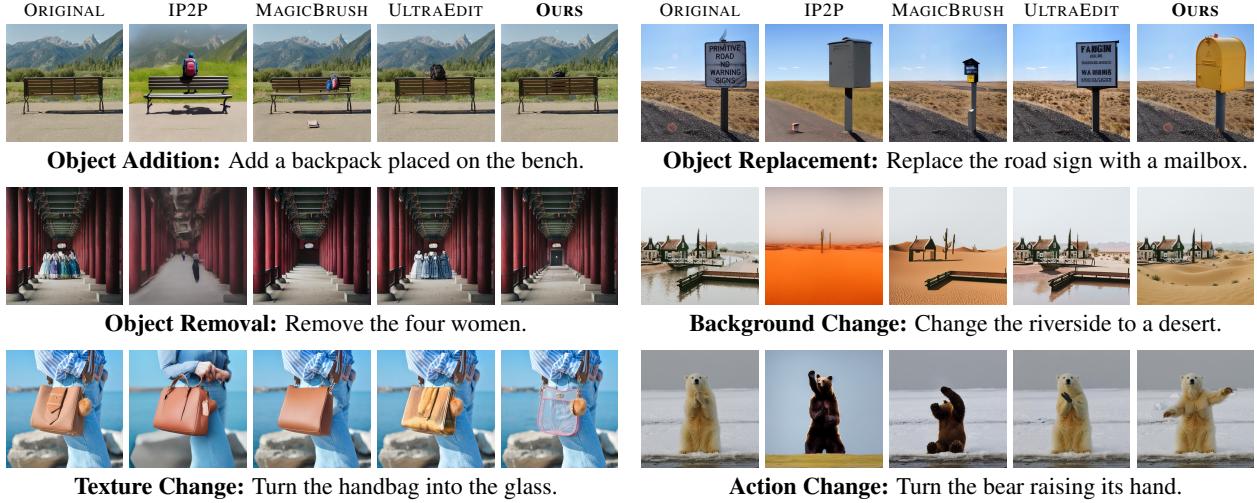
**Action Change:** Turn the bear raising its hand.

*Figure 3.* **Our workflow outperforms baseline methods in 6 editing categories from EditEval.** Each group of five images shows an example from an editing task. Each group starts from original image, followed by IP2P, MagicBrush, UltraEdit, and our results.

change), stylistic editing (style change and texture change), and structural editing tasks (action change). Due to the extremely high resolutions of images, many editing tasks in EditEval require challenging fine-grained edits. An example is removing a tiny insect from a bird's beak, which requires editing fewer than 1% of all pixels. EditEval allows a fair comparison of all methods, as none of them were trained on a dataset with the same distribution as EditEval.

**Baseline Methods.** The baseline methods include Instruct-Pix2Pix (IP2P) (Brooks et al., 2023), IP2P trained on MagicBrush (Zhang et al., 2024), and UltraEdit (Zhao et al., 2024). For these methods, we use the original editing prompts from EditEval. The inference hyperparameters are set to recommended values (Appendix A). For UltraEdit, we use the mask-based checkpoint and the same binary masks (with context dots) as in our workflow. Hence, UltraEdit similarly benefits from the extra user input, ensuring a fair comparison between our workflow and this strong baseline.

**Evaluation Metrics.** We use the CLIP (Radford et al., 2021) text-image direction similarity ($\text{CLIP}_{dir}$) and CLIP output similarity ($\text{CLIP}_{out}$) metrics from the Emu Edit benchmark (Sheynin et al., 2024). These two metrics are suitable for reference-image-free evaluation. We do not calculate the L1 distance between edited and original images, as the L1 distance does not monotonically increase with the editing quality (higher distance can simultaneously indicate better global editing performance and worse local editing performance). In EditEval, we exclude the style change task, as the task is better handled by style LoRAs or specialized checkpoints (e.g., Flux.1 [dev] Redux[9]). After excluding 24

images, there are $n = 126$ images remaining.

**Implementation Details.** For our workflow, we use 0.9 denoising strength, 5 Canny model steps, and 25 base model steps. We use Flux.1 [dev] Canny as the Canny model and Flux.1 [dev] as the base model. We also use a LoRA (Midjourney Dreamlike Fantasy FLUX LoRA[10]) with 1.0 strength on the base model to stabilize the output style. We manually draw color patches and masks with a hard round brush, which is a coarse-grained brush that limits the complexity of user input. We also used Photoshop's selection tools whenever necessary, such as the Rectangular Marquee tool and Quick Selection tool. For each original image, we spend less than one minute to do the selection and add the color patch. The color patch opacity is 0.8. To prevent cherry-picking, we fix the random seed at 0. We use target captions from EditEval as description prompts.

**Quantitative and Qualitative Evaluation.** Following Sheynin et al. (2024), we display $\text{CLIP}_{dir}$ and $\text{CLIP}_{out}$ in Table 1 to compare CLIP text-image direction similarity and CLIP output similarity. We can see that our workflow achieves the highest scores on the EditEval benchmark. In addition, we compare the edited images by baseline methods and our workflow on all six different editing categories for visual comparison. Figure 3 shows the failure patterns of baseline methods including undesired global color shift of the grass (Object Addition), removing the lens flare outside the masked region (Object Replacement), inability to remove multiple humans from the image (Object Removal), and wrong anatomy of the polar bear (Action Change). In contrast, our workflow maintains global color, keeps details

---

[9]https://huggingface.co/black-forest-labs/FLUX.1-Redux-dev

[10]https://civitai.com/models/679736/midjourney-dreamlike-fantasy-flux-lora

outside the mask untouched, removes objects as requested, and generates animals with correct anatomy.

| METHOD | $\text{CLIP}_{dir}$ | $\text{CLIP}_{out}$ |
|---|---|---|
| IP2P | $0.193 \pm 0.120$ | $0.282 \pm 0.042$ |
| MAGICBRUSH | $0.210 \pm 0.147$ | $0.302 \pm 0.038$ |
| ULTRAEDIT | $0.193 \pm 0.148$ | $0.301 \pm 0.039$ |
| OURS | $\mathbf{0.221 \pm 0.150}$ | $\mathbf{0.305 \pm 0.033}$ |

**Human Study.** To provide a comprehensive qualitative evaluation, we further conduct a human evaluation to compare results from different models. We ask a total of 9 annotators to select their preferred image from an image pair, where one image comes from our workflow and another image comes from one of the baseline methods. To avoid introducing any bias as the workflow designers, we simply ask the annotators to select the image that better follows the editing instruction (more details in Appendix B).

Figure 4 shows the human evaluation results. Our workflow is predominantly preferred. Although at a much lower chance than baseline methods, our workflow does sometimes fail to follow instructions, as indicated by the "both bad" cases. In Section 3.3, we show how the flexibility of our workflow allows us to overcome this limitation easily.
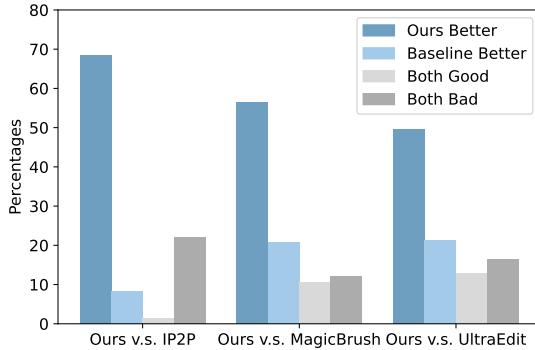


*Figure 4.* **Our workflow is preferred by the annotators over any baseline method.** In this evaluation, we use only a single editing step and fixed hyperparameters for our workflow, and thus our workflow can still fail ("both bad" cases). In Section 3.3, we show how relaxing these constraints can dramatically improve results.

**Handling Challenging Edits.** In Appendix C, we show that our workflow outperforms baselines on challenging examples from two other popular benchmarks (Emu Edit and MagicBrush). On the challenging examples, our workflow also outperforms SeedEdit (Shi et al., 2024a), a latest commercial baseline from Doubao AI.

**Ablation Study.** While FLUX.1 [dev] is a strong image generation backbone, we demonstrate that our workflow out-

performs the backbone model FLUX.1 [dev] by an ablation study. Also, we designed our workflow such that it mitigates many issues of the specialized inpainting model FLUX.1 [dev] Fill[11]. Hence, the better performance of our workflow than baseline methods does not result from simply selecting a stronger backbone. Results can be found in Appendix D.

**Minimal Burden on Users.** Despite the superior performance, our workflow imposes minimal burden on the users. While the users need to provide both the masks and the hints, the two inputs can be casually sketched, because our workflow is robust to missing details and inaccurate shapes. Examples of excellent editing results from simple user inputs are shown in Figure 5.
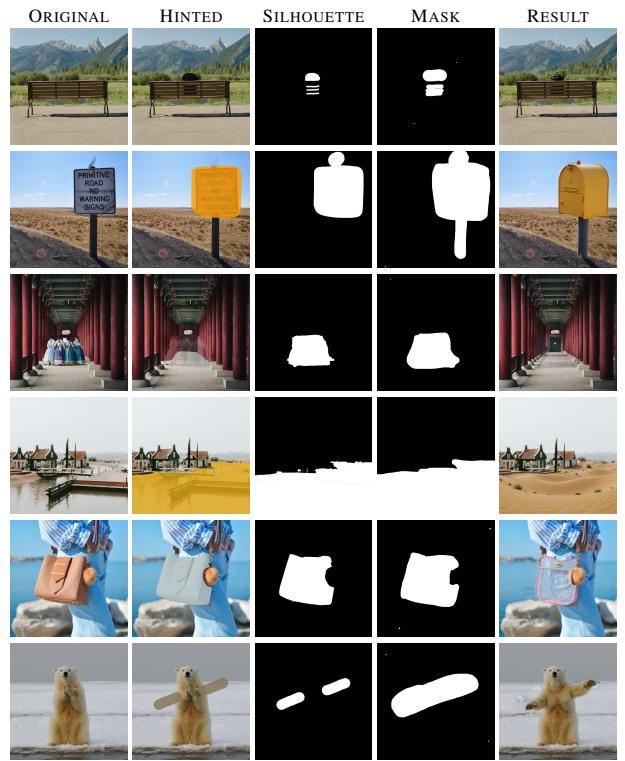


*Figure 5.* **SPICE can achieve excellent editing results from simple user inputs.** To better visualize the hint, we further show the silhouette of the hint. Silhouettes are not used in the generation process, but are manually extracted after generation for visualization purpose only. Neither the hints nor the masks need to precisely match the desired shape of the edited object.

### 3.2. Precise Editing

To quantify the precision of our workflow, we add objects while specifying 5 properties, including size, location, rotation, color, and aspect ratio (Figure 6). To specify the value of each property, we use a hint in the form of a color patch. In the edited image, we measure the property value of
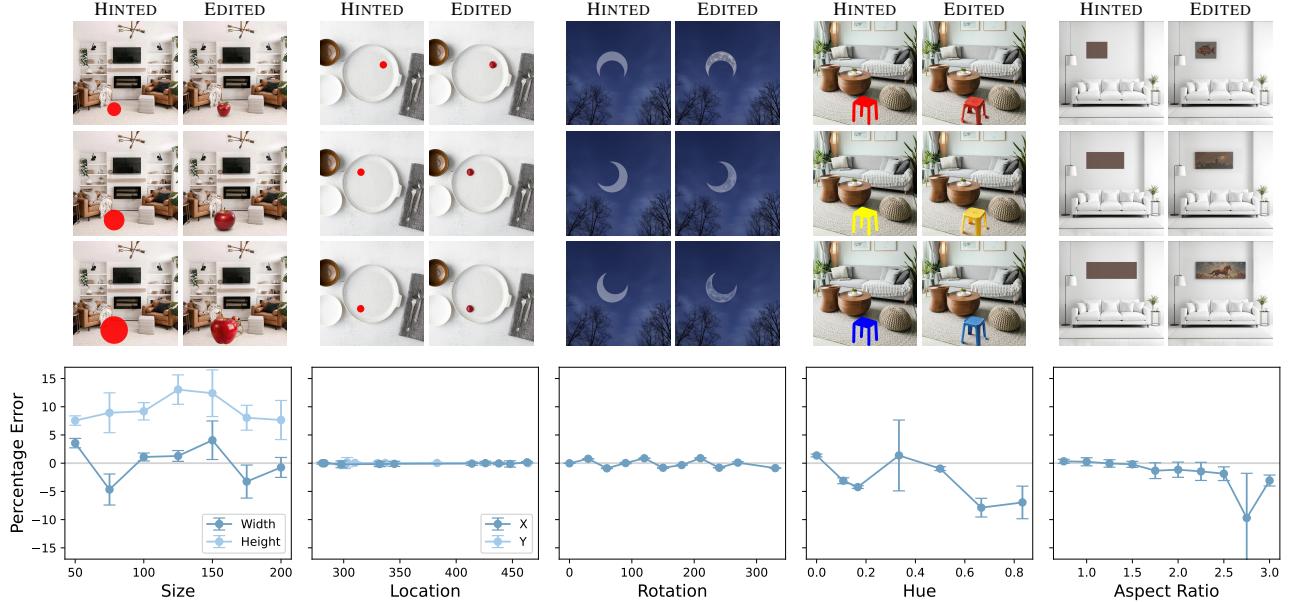
---

[11]https://huggingface.co/black-forest-labs/FLUX.1-Fill-dev

*Figure 6.* **Our workflow delivers precise and repeatable edits, where the generated objects align with user specifications.** The top images compare the user-specified color hints with the generated objects, and the bottom plots show the percentage errors. On each plot, we show the mean and standard deviation of percentage errors across 10 random seeds. All mean percentage errors are close to 0.

the added object and calculate the percentage error against the user-specified property value (Appendix E). We use 10 random seeds to account for generative variability.

Figure 6 shows that the mean of percentage errors for each specified property value is close to 0, with a small standard deviation across random seeds. There are still systematic errors, such as the height of the apple being consistently 10% larger than specified. However, for the apple, the height is larger because the apple's stem is unspecified by the hint. Overall, the results confirm that our workflow delivers precise, repeatable edits critical for real-world scenarios.

### 3.3. Iterative Editing

In this subsection, we consider two iterative editing regimes, one using fewer than 5 editing steps, another more than 100. Multiple editing steps are often necessary, as one single editing or generation step usually fails to align with complicated user requirements (further discussed in Section 3.4).

In the low-steps regime, if we are not satisfied with results from the first editing step, we are able to improve the result using another editing step with our workflow. In Section 3.1, we identify that our workflow does not work well for a few cases in single-step editing. In these cases, a large proportion of the edited region is satisfactory, and we want to fix the small remaining proportion. If we are using one of the baseline methods, we can only run the model multiple times with the same source image and editing prompt, hoping to get a better result with a different random seed (Appendix F). However, with our workflow, we can simply improve the result by recycling the output from the first step (Figure 7).

In the high-steps regime, while many previous methods support multi-step editing, heavy artifacts are generated and propagated in each step, disqualifying these methods for long editing tasks. For example, when Emu Edit is prompted to edit an image multiple times without changing a red bench in the scene[12], the image area containing the bench is still edited, and the quality of this area becomes worse with each edit. The edges of planks that constitute the bench are increasingly blurred and deformed. In contrast, our workflow steadily improves the image after multiple steps. In Figure 8(a), we show a result after 40 diverse editing steps, including outpainting, structural edits, upscaling, relighting, and composition adjustment (forming a heart shape using the sky, the crescent, the claw, and the branches). Interestingly, to relight characters according to the environment, we do not need to provide any color hints. We also do not need a specialized relighting model, such as IC-Light (Zhang et al., 2025). Full steps are shown in Appendix G. In Figure 8(b), we start with a model-generated image with various errors in shadows, objects, and anatomy. Using the same model as the backbone in our workflow, we fix these errors one by one across more than 100 editing steps. This image is produced by BoleroMix (Pony) v1.41, a checkpoint derived from SDXL. More examples of final edited results in different styles can be found in Appendix H, where usually 100 to 200 editing steps are performed to achieve the desired level of correctness in details. Each image is created from scratch within 4 hours in one sitting.
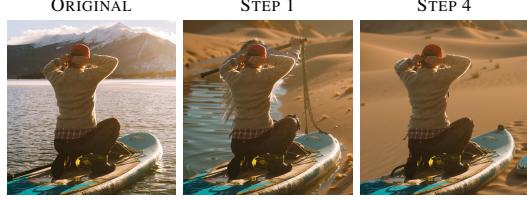
---

[12]https://emu-edit.metademolab.com/assets/videos/dog.mp4

*Figure 7.* **We use SPICE to recycle failures.** In this example, the editing instruction is "change the lake to a desert". After Step 1, SPICE fails to remove the water from the background. However, 3 more editing steps greatly improve the outcome by removing obvious artifacts. In each step, the mask and hint are redrawn to emphasize different parts of the image for editing, but the description prompt is fixed, imposing low burden on the user.
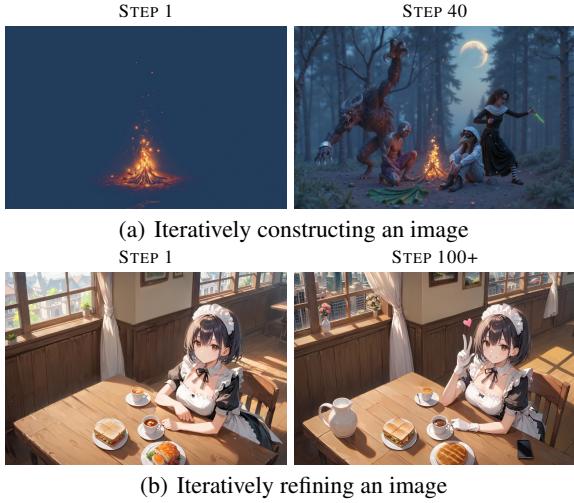


(a) Iteratively constructing an image



(b) Iteratively refining an image

*Figure 8.* **We use SPICE to consistently improve image quality during a large number of editing steps.** The editing steps can either change the structure or refine the details. The top images are generated using Flux.1 [dev]. The bottom images are generated using BoleroMix (Pony) v1.41, a checkpoint derived from SDXL.

### 3.4. Customizable Editing

Prompt-based image generation models frequently fail to fulfill their promise that users can create content according to their own will. Users are frustrated for two main reasons.

First, prompts are hard to design and often cannot instruct the model to produce complex output. A user assumes that the model follows the prompt, but the model does not interpret the prompt as humans do. One example is that DALL·E 3 (Betker et al., 2023) in November 2024 still fails to generate a rabbit with four ears or a violin without a bridge[13], as shown in Figure 9(a). Another example is that DALL·E 3 fails to interpret relations or numbers, such as a potato under a spoon or five fish (Conwell et al., 2024).

---

[13] https://cs.nyu.edu/~davise/
papers/DALL-E-Parts/PartsNovember24/
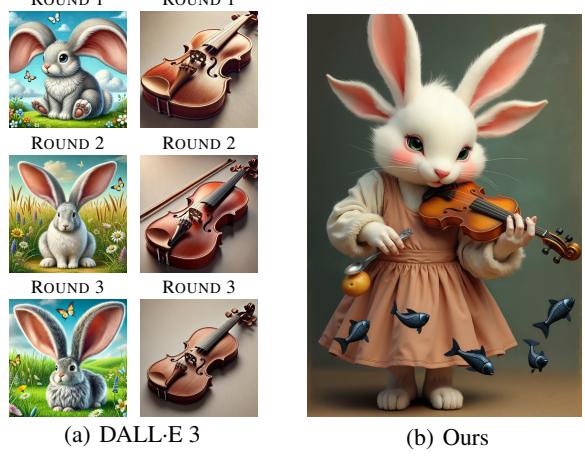DALL-E-Parts-November24.html



(a) DALL·E 3          (b) Ours

*Figure 9.* **SPICE can generate content that DALL·E 3 cannot.** Examples are a rabbit with 4 ears and a violin without a bridge. By January 2025, DALL·E 3 fails even after the user points out the error and asks the model to edit the errors multiple times.

Secondly, the innate randomness of the model dramatically increases the generation cost or even forbids the generation of complicated images. Suppose a user wants to generate an image with 10 different objects at various locations on this image. If each object has a 50% chance to be generated perfectly, the success rate of all objects being perfect is below 0.1%. In real scenarios, the perfection rate for each object will only be lower, and the number of objects larger. Then, generating such an image is practically impossible.

With SPICE, instead of desperately engineering prompts or varying the random seed, users handle these difficulties by tuning 3 hyperparameters, namely context size, denoising strength, and Canny model steps. We first demonstrate the effect of each hyperparameter (Figure 10). Then, we show how these hyperparameters can help the user reliably generate an image that DALL·E 3 cannot (Figure 9).

**Context Size.** During inpainting, a diffusion model generates content consistent with what is already in the context. In Figure 10(a), the user requires another yellow vegetable to be added to the center of the image. The prompt is "yellow vegetables on an orange background". The color hint is the yellow circle in the center. The model's interpretation of the color hint depends on the selected context. When a small context includes none of the other vegetables, the model generates multiple vegetables from the single yellow circle. When a medium context includes corners of the other vegetables, the model generates one vegetable. When a large context includes all other vegetables, the model fails to generate a new vegetable. When the context includes one vegetable either on the top left or the top right corner, the model generates a vegetable with the same direction of shade. While we do not theoretically explain the failure when using large context sizes, we observe it frequently

in other editing tasks. We also observe the failure to add objects using the "whole image" mode (Section 2.1), an observation that motivated our design of context dots.

**Denoising Strength.** The model balances hinted colors and realistic shadows at a medium denoising strength. In Figure 10(b), the user requires a Rubik's cube with certain face colors to be added to the image. Because the colors cannot be succinctly described by words, the user inserts a reference image as the hint. With a low denoising strength, the colors are preserved, but the shadows are not generated. We observe the opposite with a high denoising strength.

**Canny Model Steps.** The model balances realistic details and hinted edges at a medium number of Canny model steps. In Figure 10(c), the user requires a crescent to be added to the sky, but the prompt is simply "a moon". With few Canny model steps, the moon looks realistic, but its shape is wrong. We observe the opposite with many Canny model steps.

By slightly tuning the 3 hyperparameters (Appendix I), we are able to generate Figure 9(b), an image with a rabbit of four ears, holding a violin without a bridge with the left hand and a spoon with the right hand. There is a potato under the spoon and five fish in the air. In contrast, DALL·E 3 cannot generate even one of the components on this image.
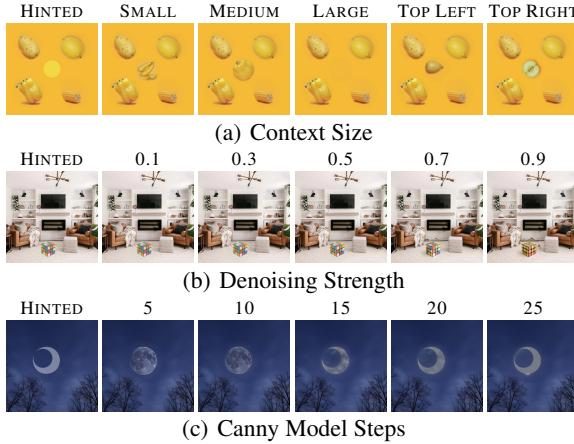


(a) Context Size

(b) Denoising Strength

(c) Canny Model Steps

*Figure 10.* **Users can achieve desired effects by customizing the value of 3 hyperparameters.** (a) To improve the consistency of edited region with a certain part of the image, the user can cover this part in the context. (b) Lower denoising strength preserves the color, whereas higher denoising strength inserts realistic shadows. (c) Fewer Canny model steps produce high-frequency details, whereas more steps allows the generated image to faithfully follow the hinted shape. Some examples look sub-optimal, because they show the results of only varying one hyperparameter. In practice, users can adjust all three parameters together for better edits.

## 4. Related Work

**Image Editing Methods.** Diffusion model-based image editing techniques can be categorized based on the type of inputs used to guide image generation (Croitoru et al., 2023; Yang et al., 2023; Huang et al., 2024). Most existing methods rely primarily on textual prompts (Brooks et al., 2023; Fu et al., 2024), while many also incorporate masks to specify the editing region (Wang et al., 2023; Zhao et al., 2024). However, due to the inherent imprecision of both prompts and masks, some methods integrate additional inputs to provide finer control over the generation process (Yang et al., 2024; Shi et al., 2024b; Liu et al., 2024a).

Among these, two methods are most closely related to our workflow. First, MagicQuill (Liu et al., 2024b) enables user-guided editing using color and edge information. However, its color guidance is constrained to low-resolution $32 \times 32$ blocks, regardless of image size, preventing fine-grained edits such as those required in EditEval. More discussion and a case study to compare SPICE and MagicQuill is included in Appendix J. Secondly, self-guidance (Epstein et al., 2023) allows users to specify object placement and size through textual prompts, but this approach remains inherently imprecise due to the limitations of language-based descriptions. SPICE overcomes the constraints of the two methods by accepting full-resolution color and edge hints and synergizing two complementary diffusion models. Hence, SPICE effectively interprets and applies the provided guidance, leading to precise and high-quality edits.

**Image Editing Benchmarks** Several datasets have been introduced for training and benchmarking image editing models, including EditBench (Wang et al., 2023), MagicBrush (Zhang et al., 2024), HQ-Edit (Hui et al., 2024), InstructPix2Pix (Brooks et al., 2023), UltraEdit (Zhao et al., 2024), Seed-Data-Edit (Ge et al., 2024), and EditEval (Huang et al., 2024). While dataset sizes have grown over time, the number of challenging test cases remains limited. Future benchmarks should prioritize structural editing tasks, such as modifying object actions or layouts, rather than focusing primarily on simpler semantic edits like color or texture changes. Simple semantic edits can already be handled almost perfectly by our workflow.

## 5. Conclusion

Existing prompt-based image editing models fail in performing local edits, working under different resolutions, following user instructions, and maintaining image quality during multiple editing steps. We propose SPICE, a training-free workflow that addressees all challenges above. We release the workflow to facilitate future research.

## Impact Statement

For the first time, our workflow unlocks the ability to perform photo-realistic edits with perfect details using diffusion

models. Since the workflow is fully open source, malicious users may misuse the workflow, leading to negative societal consequences. However, as there exist various techniques for detecting AI-generated content, we do not believe the edited images will pose an immediate threat to the public.

# References

Betker, J., Goh, G., Jing, L., Brooks, T., Wang, J., Li, L., Ouyang, L., Zhuang, J., Lee, J., Guo, Y., et al. Improving image generation with better captions. *Computer Science. https://cdn. openai. com/papers/dall-e-3. pdf*, 2(3): 8, 2023.

Brooks, T., Holynski, A., and 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.

Conwell, C., Tawiah-Quashie, R., and Ullman, T. Relations, negations, and numbers: Looking for logic in generative text-to-image models. *arXiv preprint arXiv:2411.17066*, 2024.

Croitoru, F.-A., Hondru, V., Ionescu, R. T., and Shah, M. Diffusion models in vision: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45 (9):10850–10869, 2023.

Epstein, D., Jabri, A., Poole, B., Efros, A., and Holynski, A. Diffusion self-guidance for controllable image generation. *Advances in Neural Information Processing Systems*, 36: 16222–16239, 2023.

Fan, X., Bhattad, A., and Krishna, R. Videoshop: Localized semantic video editing with noise-extrapolated diffusion inversion. In Leonardis, A., Ricci, E., Roth, S., Russakovsky, O., Sattler, T., and Varol, G. (eds.), *Computer Vision – ECCV 2024*, pp. 232–250, Cham, 2025. Springer Nature Switzerland. ISBN 978-3-031-73254-6.

Fu, T.-J., Hu, W., Du, X., Wang, W. Y., Yang, Y., and Gan, Z. Guiding instruction-based image editing via multimodal large language models. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=S1RKWSyZ2Y.

Ge, Y., Zhao, S., Li, C., Ge, Y., and Shan, Y. Seed-data-edit technical report: A hybrid dataset for instructional image editing, 2024. URL https://arxiv.org/abs/2405.04007.

Hirota, Y., Andrews, J., Zhao, D., Papakyriakopoulos, O., Modas, A., Nakashima, Y., and Xiang, A. Resampled datasets are not enough: Mitigating societal bias beyond single attributes. In Al-Onaizan, Y., Bansal, M., and Chen, Y.-N. (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 8249–8267, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main. 471. URL https://aclanthology.org/2024.emnlp-main.471/.

Huang, Y., Huang, J., Liu, Y., Yan, M., Lv, J., Liu, J., Xiong, W., Zhang, H., Chen, S., and Cao, L. Diffusion model-based image editing: A survey. *arXiv preprint arXiv:2402.17525*, 2024.

Hui, M., Yang, S., Zhao, B., Shi, Y., Wang, H., Wang, P., Zhou, Y., and Xie, C. Hq-edit: A high-quality dataset for instruction-based image editing, 2024. URL https://arxiv.org/abs/2404.09990.

Levin, E. and Fried, O. Differential diffusion: Giving each pixel its strength. *arXiv preprint arXiv:2306.00950*, 2023.

Liu, H., Xu, C., Yang, Y., Zeng, L., and He, S. Drag your noise: Interactive point-based editing via diffusion semantic propagation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 6743–6752, 2024a.

Liu, Z., Yu, Y., Ouyang, H., Wang, Q., Cheng, K. L., Wang, W., Liu, Z., Chen, Q., and Shen, Y. Magicquill: An intelligent interactive image editing system. *arXiv preprint arXiv:2411.09703*, 2024b.

Lugmayr, A., Danelljan, M., Romero, A., Yu, F., Timofte, R., and Van Gool, L. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11461–11471, 2022.

Ma, N., Tong, S., Jia, H., Hu, H., Su, Y.-C., Zhang, M., Yang, X., Li, Y., Jaakkola, T., Jia, X., et al. Inference-time scaling for diffusion models beyond scaling denoising steps. *arXiv preprint arXiv:2501.09732*, 2025.

Podell, D., English, Z., Lacey, K., Blattmann, A., Dockhorn, T., Müller, J., Penna, J., and Rombach, R. SDXL: Improving latent diffusion models for high-resolution image synthesis. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=di52zR8xgf.

Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.

Sheynin, S., Polyak, A., Singer, U., Kirstain, Y., Zohar, A., Ashual, O., Parikh, D., and Taigman, Y. Emu edit: Precise image editing via recognition and generation tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8871–8879, 2024.

Shi, Y., Wang, P., and Huang, W. Seededit: Align image re-generation to image editing, 2024a. URL https://arxiv.org/abs/2411.06686.

Shi, Y., Xue, C., Liew, J. H., Pan, J., Yan, H., Zhang, W., Tan, V. Y., and Bai, S. Dragdiffusion: Harnessing diffusion models for interactive point-based image editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8839–8849, 2024b.

Singhal, R., Horvitz, Z., Teehan, R., Ren, M., Yu, Z., McKeown, K., and Ranganath, R. A general framework for inference-time scaling and steering of diffusion models. *arXiv preprint arXiv:2501.06848*, 2025.

Snell, C., Lee, J., Xu, K., and Kumar, A. Scaling llm test-time compute optimally can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*, 2024.

Wang, S., Saharia, C., Montgomery, C., Pont-Tuset, J., Noy, S., Pellegrini, S., Onoe, Y., Laszlo, S., Fleet, D. J., Soricut, R., et al. Imagen editor and editbench: Advancing and evaluating text-guided image inpainting. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 18359–18369, 2023.

Wu, X., Guan, T., Li, D., Huang, S., Liu, X., Wang, X., Xian, R., Shrivastava, A., Huang, F., Boyd-Graber, J. L., Zhou, T., and Manocha, D. AutoHallusion: Automatic generation of hallucination benchmarks for vision-language models. In Al-Onaizan, Y., Bansal, M., and Chen, Y.-N. (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 8395–8419, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-emnlp. 493. URL https://aclanthology.org/2024.findings-emnlp.493/.

Yang, L., Zhang, Z., Song, Y., Hong, S., Xu, R., Zhao, Y., Zhang, W., Cui, B., and Yang, M.-H. Diffusion models: A comprehensive survey of methods and applications. *ACM Computing Surveys*, 56(4):1–39, 2023.

Yang, Y., Peng, H., Shen, Y., Yang, Y., Hu, H., Qiu, L., Koike, H., et al. Imagebrush: Learning visual in-context instructions for exemplar-based image manipulation. *Advances in Neural Information Processing Systems*, 36, 2024.

Zhang, K., Mo, L., Chen, W., Sun, H., and Su, Y. Magicbrush: A manually annotated dataset for instruction-guided image editing. *Advances in Neural Information Processing Systems*, 36, 2024.

Zhang, L., Rao, A., and Agrawala, M. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3836–3847, 2023.

Zhang, L., Rao, A., and Agrawala, M. Scaling in-the-wild training for diffusion-based illumination harmonization and editing by imposing consistent light transport. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=u1cQYxRI1H.

Zhao, H., Ma, X., Chen, L., Si, S., Wu, R., An, K., Yu, P., Zhang, M., Li, Q., and Chang, B. Ultraedit: Instruction-based fine-grained image editing at scale. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024. URL https://openreview.net/forum?id=9ZDdlgH6O8.

## A. Baseline Model Hyperparameters

For InstructPix2Pix, we use the model and hyperparameters according the Hugging Face official page[14]. The number of inference steps is set to 10, and the image guidance scale is set to 1.

For MagicBrush, we also refer to the Hugging Face page[15]. Specifically, we use the recommended "recent best checkpoint", which is `MagicBrush-epoch-52-step-4999.ckpt`. We do not provide extra hyperparameters when calling the inference script from the command line, except that we provide a random seed. Hence, all parameters are fixed at their default values, presumably recommended by the model developers.

For UltraEdit, we first use the model version that supports masks[16]. We use the recommended hyperparameter values. We also use the free-form version[17], since binary masks are not available in Emu Edit and MagicBrush test sets. Since there is not a recommended set of hyperparameters on the official page for the free-form version, we use the same hyperparameters for the free-form version and the masked version.

## B. Human Evaluation Details

For human evaluation, we adopt a preference setting. An annotator is asked to choose a better image out of a pair or consider both images to be "good" or "bad". We construct an *image pair* by using one image from our workflow and another image from one of the baselines. In total, there are $3 \times 126 = 378$ image pairs for 3 baselines. For each of the 3 subsets (126 image pairs), we ask 3 annotators for evaluation, totaling 9 annotators. For each image pair, the annotator is shown 3 images, including the original image, an edited image from one method, and an edited image from another method. The order of our workflow and the baseline is randomly chosen for each image pair. The order is hidden from the annotators. We show the editing category and the editing prompt, but not the source prompt or the target prompt. Before the annotation process begins, each annotator is given the same evaluation instructions, where we do not specify detailed criteria for preference other than the adherence to the editing instruction.

An example of an image pair the annotator sees is shown in Figure 11. To ensure fairness, we send the same instructions to the annotators. To ensure minimal bias, we do not explain the relative strengths and weaknesses of each method. We do not further communicate with annotators or provide clarifications in written or spoken form, until the annotators finishes the task. The only message we send to annotators is shown below (from "Hi" to "winner").

Hi [NAME],

We are working on a project about using diffusion models to edit images. We would like you to help evaluate the results. This evaluation requires you to choose the better image from a pair. There are a total of 126 pairs, so the evaluation would not take long.

The image pairs can be found at

[LINK TO THE FOLDER CONTAINING IMAGE PAIRS]

Please submit your evaluation using this spreadsheet

[LINK TO THE SPREADSHEET]

Put a 1 in the option you choose. Please feel free to discuss anything you observe with us. Thanks for your help!

The evaluation instruction:

Please choose from A and B the image that better follows the editing instruction. If the instruction is followed by both A and B, choose the image that looks better. Only choose "both good" or "both bad" if there is no clear winner.

---

[14]https://huggingface.co/timbrooks/instruct-pix2pix
[15]https://huggingface.co/osunlp/InstructPix2Pix-MagicBrush
[16]https://huggingface.co/BleachNick/SD3_UltraEdit_w_mask
[17]https://huggingface.co/BleachNick/SD3_UltraEdit_freeform

**Object Addition 15**

**Add a backpack placed on the bench.**
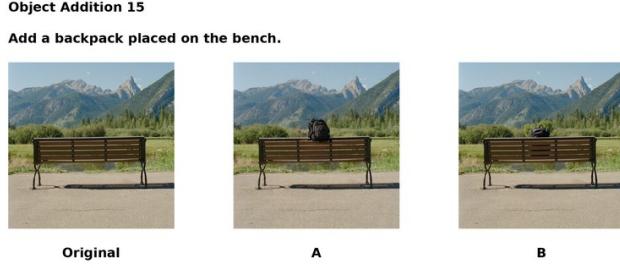
Original        A        B

*Figure 11.* **An example of an image pair that the annotator sees.** The annotator does not know which option comes from which model. While they can see the name of the baseline method they are evaluating, none of the annotators know about the baseline methods.

## C. Challenging Examples from More Benchmarks

We further compare the performance of our workflow against baseline methods on selected challenging examples. A challenging example is an example where all baseline methods fail. We select challenging examples from the Emu Edit and MagicBrush test sets. For the Emu Edit test set, the baseline methods are Emu Edit, IP2P, MagicBrush, and UltraEdit (Figure 12 and Figure 13). For the MagicBrush test set, the baseline methods are Reference (the best DALL·E 2 generations), IP2P, MagicBrush, and UltraEdit (Figure 14 and Figure 15). After we exclude the style change task, both test sets cover 7 different editing categories, as reported in the original papers. We select 2 challenging examples from each category, totaling 14 for each test set. Our workflow succeeds on all challenging examples. Our workflow also outperforms SeedEdit from Doubao AI on these examples. The authors translate the original prompts into Chinese when using SeedEdit.

## D. Ablation Studies

We qualitatively demonstrate the effect of removing each component from our workflow (Figure 16). Then, we demonstrate that our workflow mitigates many issues of the specialized inpainting model FLUX.1 [dev] Fill (Figure 17). Together, the analysis in this section shows that our workflow improves upon the backbone image generation model.

### D.1. Necessity of Components

Removing any component from our workflow leads to the following negative consequences (Figure 16).

**Context Selection.** Without an appropriate context, the generated region might show a different style than the overall image. Even when the style is correct, the inpainted region could still show a very different texture, an artifact that distinguishes the region from surrounding pixels.

**Soft Inpainting.** Without soft-inpainting, both major and minor artifacts appear. Here, the major artifact is the missing shoes from the doll. The minor artifact is the abrupt change of texture around the mask edges (on the bench backrest).

**Color and Edge Hinting.** Without color and edge hints, the object sometimes can still be added. However, the size of the object does not follow the user's requirement.

**Two-Stage Denoising.** Without two-stage denoising, the generated object and the color patch have different sizes.

### D.2. Flux.1 [dev] Fill Failures

Our workflow mitigates many issues of the Flux.1 [dev] Fill, a specialized inpainting model trained on the same backbone of Flux.1 [dev]. We discuss the issues below and demonstrate them in Figure 17.

**Global Editing Failures.** While Flux.1 [dev] Fill performs comparably with our method in some local editing tasks, it has a low success rate in global (background) editing. This is probably because Flux.1 [dev] Fill was trained with more images
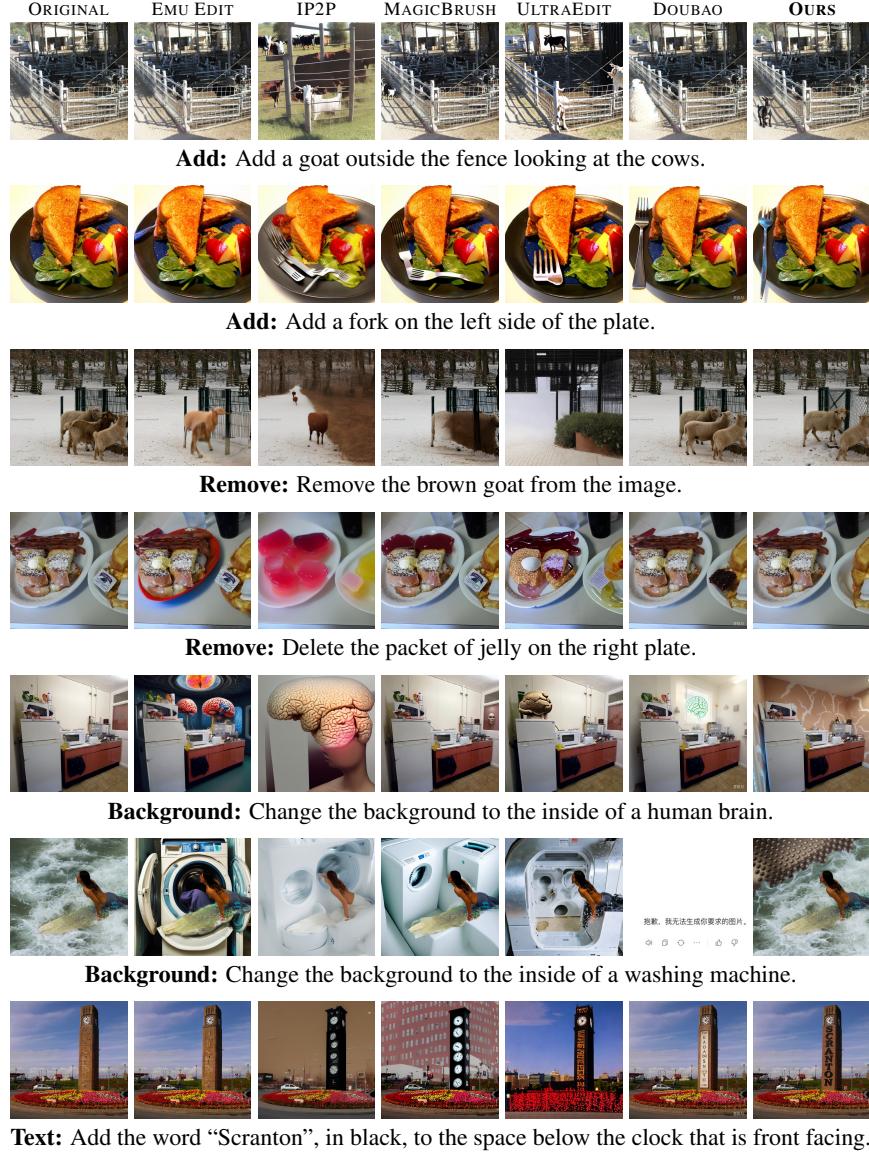
**Add:** Add a goat outside the fence looking at the cows.

**Add:** Add a fork on the left side of the plate.

**Remove:** Remove the brown goat from the image.

**Remove:** Delete the packet of jelly on the right plate.

**Background:** Change the background to the inside of a human brain.

**Background:** Change the background to the inside of a washing machine.

**Text:** Add the word "Scranton", in black, to the space below the clock that is front facing.

*Figure 12.* **Our workflow performs the best on the first half of 14 challenging examples from Emu Edit.** Doubao AI refuses to generate an image for the washing machine example, showing an error message of "Sorry, I cannot generate the image you requested".
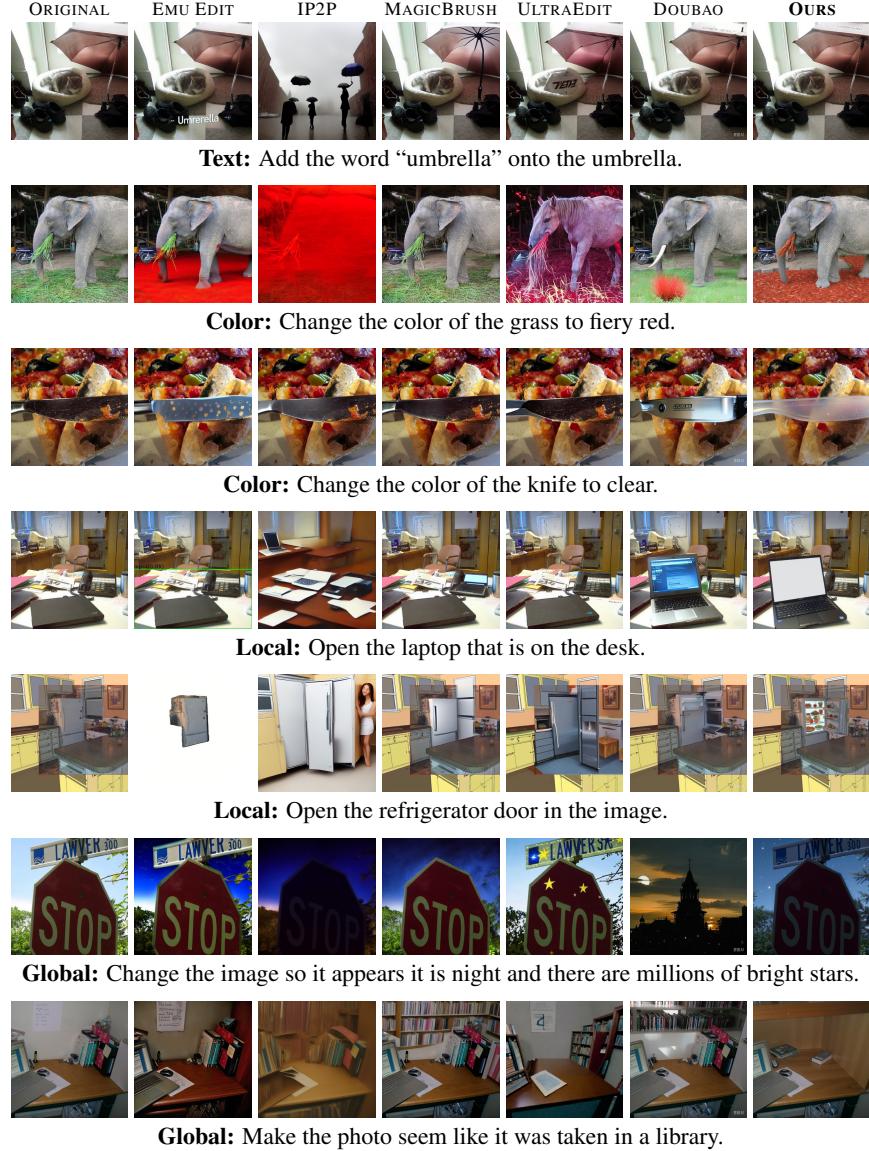
*Figure 13.* **Our workflow performs the best on the second half of 14 challenging examples from Emu Edit.**
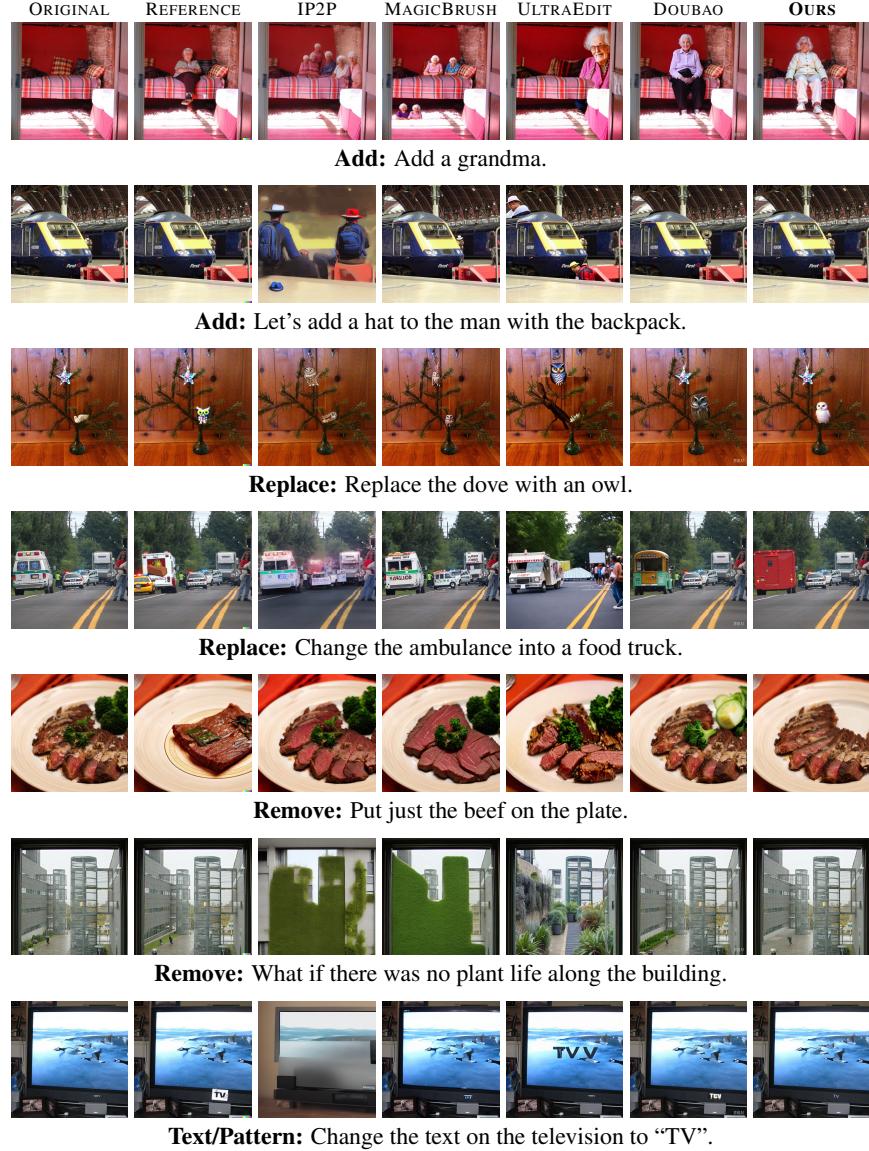
| ORIGINAL | REFERENCE | IP2P | MAGICBRUSH | ULTRAEDIT | DOUBAO | **OURS** |
|----------|-----------|------|------------|-----------|--------|----------|

**Add:** Add a grandma.

**Add:** Let's add a hat to the man with the backpack.

**Replace:** Replace the dove with an owl.

**Replace:** Change the ambulance into a food truck.

**Remove:** Put just the beef on the plate.

**Remove:** What if there was no plant life along the building.

**Text/Pattern:** Change the text on the television to "TV".

*Figure 14.* **Our workflow performs the best on the first half of 14 challenging examples from MagicBrush.**

16

| ORIGINAL | REFERENCE | IP2P | MAGICBRUSH | ULTRAEDIT | DOUBAO | **OURS** |

**Text/Pattern:** Change the stop sign into a no entry sign.

**Color:** Make one of the sheep a black sheep.

**Color:** Change the fire hydrant from red to yellow.

**Action:** Let the cat look shocked.

**Action:** Make the man smile.

**Counting:** Turn the two windows into a single window.

**Counting:** Let there be a stop sign and only one road sign.
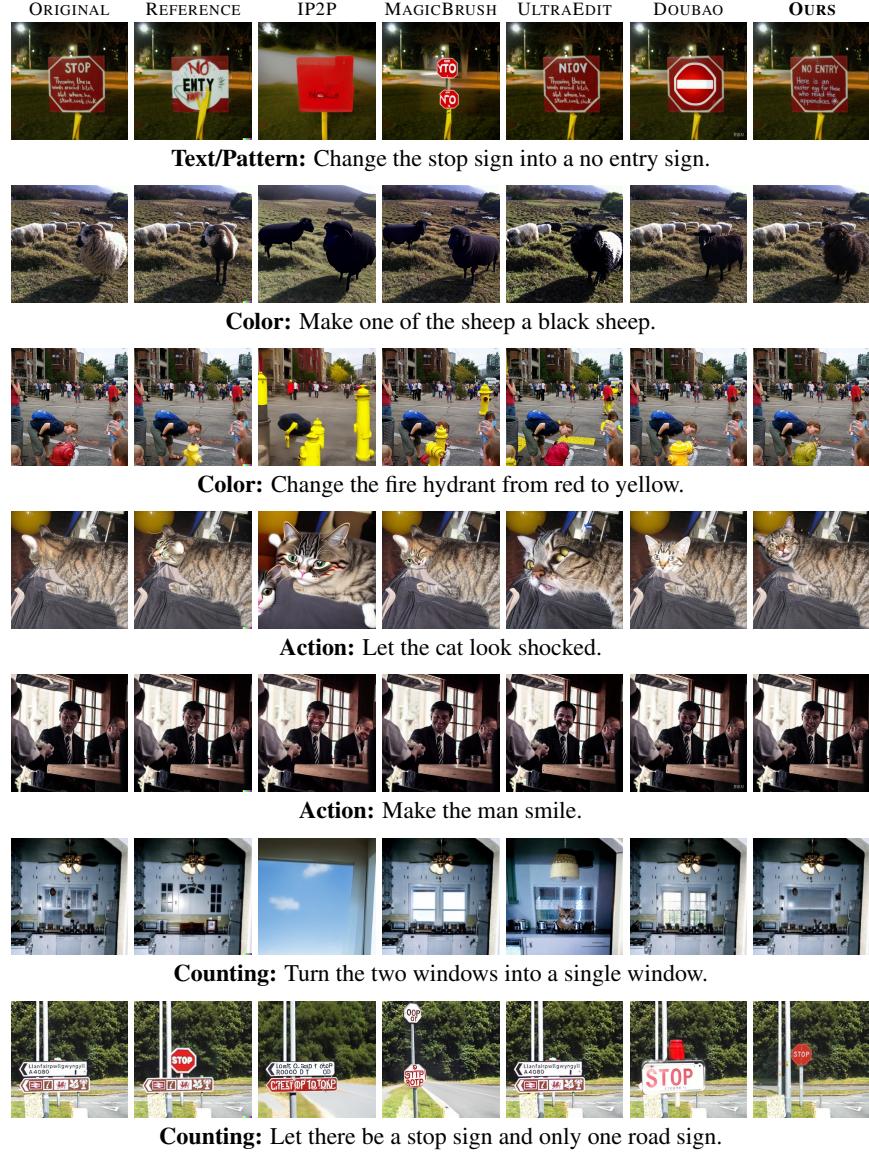
*Figure 15.* **Our workflow performs the best on the second half of 14 challenging examples from MagicBrush.**

with local edits than ones with global edits. Our method requires no further training, so no additional bias is imposed on either local or global editing.

**Color Drifts.** Flux.1 [dev] Fill sometimes results in an obvious color drift outside the mask (saturation of the grass color drops), a weakness acknowledged by the official developers[18]. This is because unlike our workflow, Flux.1 [dev] Fill is not strictly a local editing model. By design, our method strictly preserves pixel values outside the mask.

**Prompt Bleeding.** When multiple objects are mentioned in the prompt, Flux.1 [dev] Fill suffers from the well-known prompt bleeding issue, where characteristics of objects are mixed instead of being independent (a flamingo that looks like a camel). Empirically, by providing color and edge hints, our method effectively handles the prompt bleeding issue, despite using the same prompt containing multiple objects. Moreover, including multiple objects in the prompt improves model's understanding of the context in our workflow, which sometimes leads to even better results.
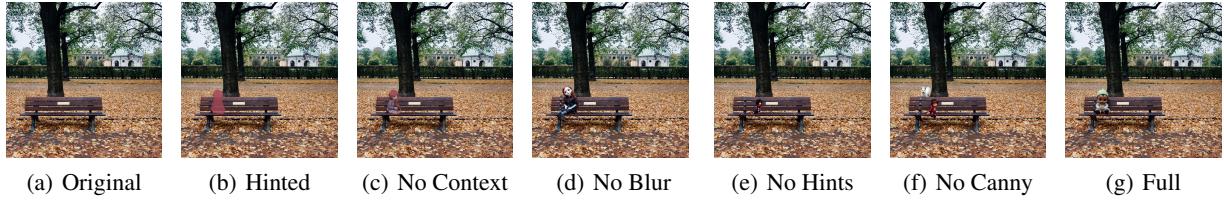
| (a) Original | (b) Hinted | (c) No Context | (d) No Blur | (e) No Hints | (f) No Canny | (g) Full |

*Figure 16.* **Ablation studies show that each component of our workflow is necessary.** Removing any one leads to suboptimal quality.

## E. Measuring Properties of Generated Objects

To quantify the precision of our workflow, we need to measure the properties of the generated objects. The properties include size, location, rotation, color, and aspect ratio. We use the Language Segment-Anything[19] tool to identify the segmentation mask of a generated object. Here, we use segmentation mask to refer to the pixels that cover the object. We use the default model (`sam2.1_hiera_small`) and settings provided by the developers, which empirically performs well for our purpose. For an object, its width and height are defined as the width and height of the bounding box of the segmentation mask. The location is defined as the center of the bounding box. The rotation (for the crescent) is calculated as the direction from the center of mass of the segmentation mask to the center of the bounding box. The color is defined as the average RGB value of pixels on the generated image covered by the segmentation mask. To find a single value to represent the color, we convert the average RGB value into HSV representation and use the hue. The aspect ratio is the width divided by height.

After measuring these properties of the generated object, we calculate a percentage error between the generated property value and specified property value. For size, we vary the diameter of a red circle and compare the height and width of the apple's bounding box with the diameter. For location, we vary the center coordinates of a red circle and compare the center coordinates of the cherry's bounding box with the center coordinates of the circle. For rotation, we vary the rotation angle of a crescent-shaped color patch and compare the orientation of the crescent with the orientation of the color patch. For color, we vary the hue of a color patch and compare the average hue of the plastic chair to the hue of the patch. For aspect ratio, we vary the width-to-height ratio of a rectangle and compare the width-to-height ratio of the painting with the width-to-height
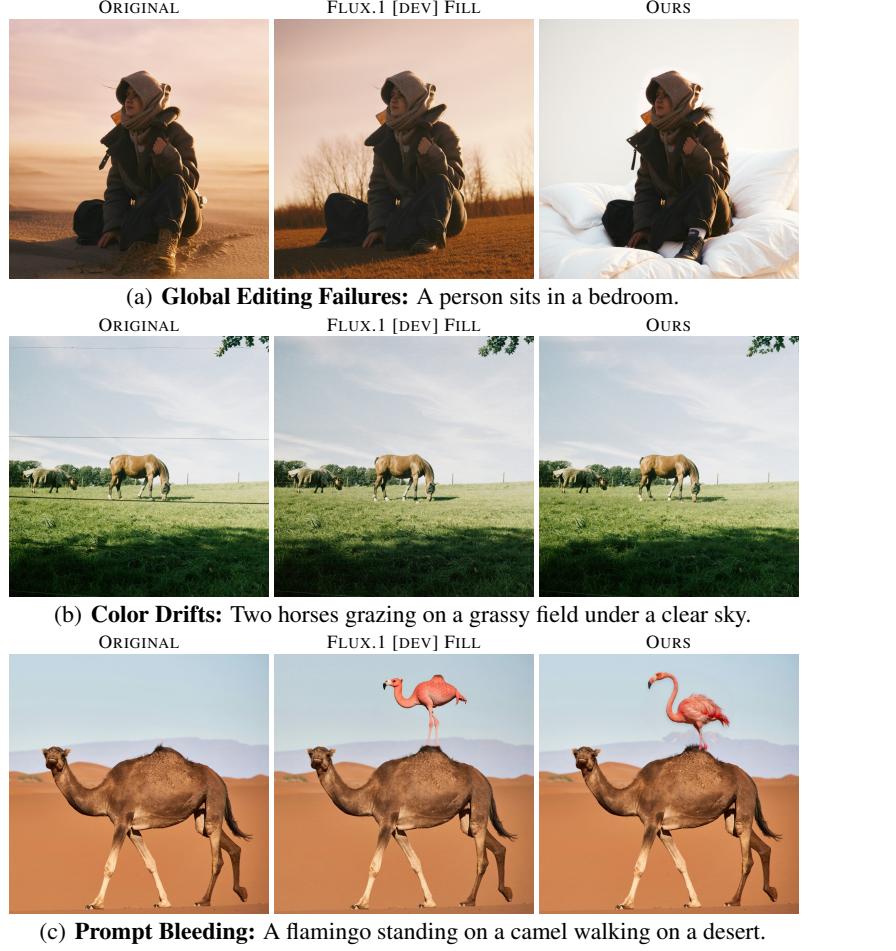
---

[18] https://huggingface.co/black-forest-labs/FLUX.1-Fill-dev#limitations
[19] https://github.com/luca-medeiros/lang-segment-anything

(a) **Global Editing Failures:** A person sits in a bedroom.



(b) **Color Drifts:** Two horses grazing on a grassy field under a clear sky.



(c) **Prompt Bleeding:** A flamingo standing on a camel walking on a desert.

*Figure 17.* **Our workflow mitigates issues of the strong Flux.1 [dev] Fill model.** While Flux.1 [dev] Fill is a specialized inpainting checkpoint, it suffers intrinsic limitations including failures in global editing, color drifts, and prompt bleeding. In contrast, our training-free workflow is free from these issues. In the captions, we show the description prompt for each editing task. Note that we do not use the editing prompt ("add a flamingo on top of the camel"), because these two methods are both designed to accept description prompts instead of editing prompts.

ratio of the rectangle. The definition of percentage errors are listed below:

$$\text{Percentage Error of Width} = \frac{\text{Generated Width} - \text{Specified Width}}{\text{Specified Width}} \times 100\%,$$

$$\text{Percentage Error of Height} = \frac{\text{Generated Height} - \text{Specified Height}}{\text{Specified Height}} \times 100\%,$$

$$\text{Percentage Error of X} = \frac{\text{Generated X} - \text{Specified X}}{\text{Specified X}} \times 100\%,$$

$$\text{Percentage Error of Y} = \frac{\text{Generated Y} - \text{Specified Y}}{\text{Specified Y}} \times 100\%,$$

$$\text{Percentage Error of Rotation} = \frac{\text{Generated Rotation} - \text{Specified Rotation}}{360°} \times 100\%,$$

$$\text{Percentage Error of Color} = \frac{\text{Generated Hue} - \text{Specified Hue}}{1.0} \times 100\%,$$

$$\text{Percentage Error of Aspect Ratio} = \frac{\text{Generated Aspect Ratio} - \text{Specified Aspect Ratio}}{\text{Specified Aspect Ratio}} \times 100\%.$$

To account for the periodic nature of rotation and hue, we use 360° for rotation and 1.0 for hue instead of the specified property value.

## F. Changing the Random Seed

With baseline methods, a user can change the random seed to get multiple results and select the best one. However, even with more compute, baseline models will consistently fail on the same challenging example (Figure 18). In contrast, with our workflow, a user can better utilize the increased compute if the user performs editing step by step.
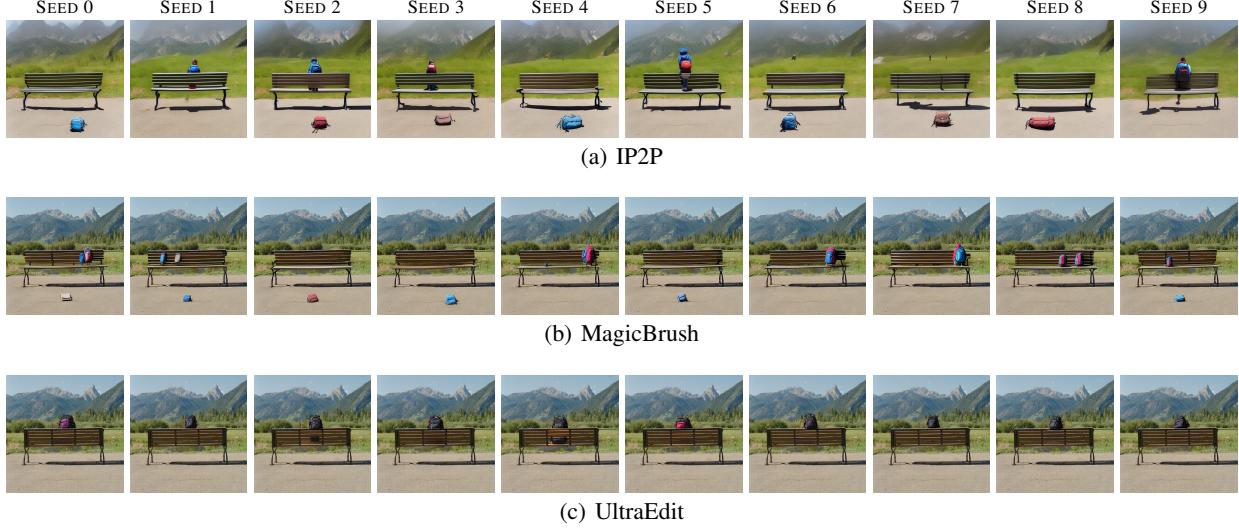


(a) IP2P

(b) MagicBrush

(c) UltraEdit

Figure 18. **For baseline methods, changing the random seed does not help.** The random seeds are shown on the top. With 10 times the original compute, baseline methods still fail on the challenging editing task from EditEval of adding a bag to the bench. Our workflow succeeds in the first edit and will succeed in fewer than 10 edits should the first edit fail. The original image and our result are in Figure 3.

## G. Iterative Construction

Figure 19 shows all 40 steps we use to iteratively construct Figure 8(a). Some steps involve subtle but important edits that are best visualized at a high resolution. We will release a video to better illustrate the editing operation done at each step.

## H. Results of Different Art Styles

Figure 20 shows two additional results. One is generated with BoleroMix (Pony) v1.41 as the backbone, and the other is generated with the same Flux.1 [dev] as the backbone. These additional results illustrate the flexibility of our workflow in editing and iteratively generating images of different art styles. Note that these are fan-made images, where details may differ from those of the original characters.

## I. Hyperparameter Recommendations

To generate Figure 9(b), the recommended hyperparameter ranges are shown in Table 2. Note that several consecutive editing steps are necessary in a certain editing region to achieve the shown effect.

## J. Comparison with MagicQuill

MagicQuill (Liu et al., 2024b) is a method very similar to ours, because it also accepts edge and color information from the user as its input. However, the following differences in design choices lead to a better performance of our method.

*Figure 19.* **We use 40 steps to iteratively construct the image and fix detail errors.** With our workflow, high frequency details do not deteriorate over the steps. With baseline methods, the deterioration happens at the first step.

BoleroMix (Pony) v1.41 (Characters from Touhou Project)    Flux.1 [dev] (Characters from Hades 2)



*Figure 20.* **Our workflow excels at editing and iteratively generating images of different art styles.** Both images are generated using around 100 to 200 editing steps, within 4 hours in one sitting. Character LoRAs are not used, so many editing steps are spent on correcting details of the character's clothes. Note that due to the complex composition and high resolution (2048×2048), it is extremely hard, if not impossible, to generate these images with a single text-to-image step.

**Fewer ControlNets.** MagicQuill uses three ControlNets, namely inpainting, edge, and color ControlNets. The control branch also needs to be further trained. Meanwhile, our workflow uses only one ControlNet, which is the Canny-edge ControlNet, and no further training is required.

**Higher Flexibility.** As our workflow uses different models for different stages, it can be adapted even when the ControlNets are not released as a module on the base model but as a full model (for example, the official Flux.1 [dev] Canny model by

*Table 2.* **We recommend different hyperparameter combinations according to specific editing tasks.** Empirically, these combinations have a higher success rate than others.

| OBJECT | CONTEXT | DENOISING | CANNY |
|---|---|---|---|
| Ears | Head | Low | High |
| Bridge | Violin | Moderate | Moderate |
| Potato | Arm | High | Moderate |
| Fish | Dress | Moderate | Low |

Black Forest Labs). Hence, the higher flexibility allows us to adapt our workflow to the Flux.1 [dev] base model within one week of the release of Flux.1 [dev] Canny in November 2024, whereas MagicQuill still supports neither SDXL or FLUX models by April 2025[20].

**Better Detail Preservation.** Our workflow does not downsample the color information. Hence, high frequency details are preserved from user input. In contrast, MagicQuill downsamples user color input to a $32 \times 32$ resolution during pre-processing, an information bottleneck that limits the ability of users to provide detailed color hints.

**Disentangled Mask and Colors** MagicQuill assumes that the user only wants to edit the region around locations where the color patch is provided. The mask is uniformly expanded in all directions from the color patch. However, consider the example of adding a backpack on a bench (Figure 21). In this case, the mask should be expanded more in the horizontal direction than the vertical direction, because the mask should not fully cover the backrest of the bench. To address this difficulty, the disentangled design in our workflow returns the freedom to the user.

**Streamlined Components.** First, our workflow does not uses an MLLM to guess user inputs. While MLLM provides convenience in MagicQuill for common objects, it heavily hallucinates when adding an object that is partially occluded by other objects, because the provided color patch does not have the shape of the full un-occluded object. The hallucination leads to negative user experience when challenging edits are required. Second, our workflow does not use an additional CNN to extract edges. Empirically, the Canny edges provides satisfactory performance and can be quickly extracted in our workflow.

**Unlimited Resolution.** By default, MagicQuill generates an image with the resolution of the whole image and replaces the masked region with generated content. This is troublesome when the user wants to iteratively perform global and local edits on the same image. Our workflow mitigates the multi-scale editing issue by proposing context dots, an intuitive and user-friendly method for specifying generation resolutions. The existing "Resolution Adjustement" functionality in MagicQuill (updated December 6, 2024) is not designed for multi-scale editing.

**Multi-Purpose.** By changing parameters and resolutions, our workflow smoothly transitions from an image editor to a detail enhancer, or from an inpainter to an outpainter. Examples of both transitions can be found in Figure 19. Neither detail enhancing nor outpainting is supported by MagicQuill.

Due to the above advantages, our method outperforms MagicQuill on a challenging example (Figure 21), with the same Realistic Vision V6.0 B1 (based on SD 1.5) model[21] as the backbone. To ensure fair comparison, we take the first-shot, single-step result from our method, but try our best to tune the hyperparameters and to draw accurate shapes for MagicQuill. After over 20 attempts, we are still unable to correctly add the backpack onto the bench using MagicQuill.

Here, we additionally provide comparison with GPT-4o (March 2025) and Gemini 2.0 Flash (March 2025). The fact that an SD-1.5-based model (December 2023) outperforms latest commercial VLMs suggests intrinsic limitations of text-only image editing frameworks.

We encourage readers with more experience in the baseline methods to independently verify their performance upper-bound.

---

[20]https://github.com/ant-research/MagicQuill/issues/89
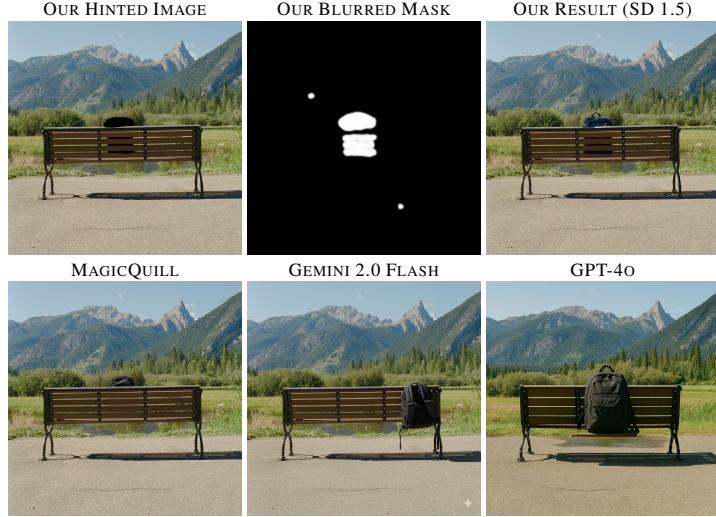[21]https://civitai.com/models/4201?modelVersionId=245598

*Figure 21.* **Using Realistic Vision V6.0 B1, an SD-1.5-based model released in December 2023, as the backbone, our workflow outperforms both MagicQuill (with the same backbone) and latest commercial VLMs (released in March 2025).** For our workflow and MagicQuill, we used the same description prompt "black backpack on a bench". For Gemini 2.0 Flash and GPT-4o, we used the same editing prompt "add a black backpack onto this bench". Only our method results in almost perfect spatial relationships, with only a minor error at the bottom of the backpack. MagicQuill is unable to generate the lower part of the backpack, which is supposed to be visible through slits on the backrest. Gemini 2.0 Flash and GPT-4o both generate a floating backpack in front of the backrest. The original image and our better result of using Flux.1 [dev] as the backbone are shown in Figure 3.

## K. Limitations

Our workflow has the following limitations.

**Texture.** Our workflow performs seamless inpainting when the original image is either (1) the output of the workflow or (2) the text-to-image output of the same base model used in the workflow. If the goal is to edit a real photo or an image generated by another base model, the distribution mismatch will cause minor texture mismatches that are impossible to fix perfectly. While such artifacts are barely noticeable with human eyes, they would not escape heuristics-based editing artifact detectors.

**Reference Images.** Our workflow does not allow copying all details of a reference image in a single editing step. If the user requires editing complicated apparel, including larger pieces of clothing and smaller accessories, we recommend using multiple editing steps, varying the context size according to the item size, and adding one item at a time.

**Local Minima.** If the base model has some strong local minima, our workflow is unable to help the user generate content out of those local minima. Two examples from the base model (Flux.1 [dev] with a LoRA) we use are (1) a fixed hairstyle regardless of the edge hints and (2) two patellas on one knee regardless of the color hints. To mitigate these artifacts, a convenient and successful workaround is to increase the Canny model steps (correspondingly reducing the base model steps), since the Canny model and the base model seldom suffer from the same local minima. The resulting image after successful edits is shown in the right panel of Figure 20.

**User Input.** To achieve a high editing quality, a user needs to provide mask and color guidance and to tune hyperparameters. This is harder for the user than baseline methods, where a user only needs to select a better-looking image, if there exists such an image in the outputs. We have not thoroughly tested our workflow in fully automated settings.

## L. Workflow Implementation

The steps in our workflow can be easily implemented (or have already been separately implemented) in the most popular diffusion model Web UIs (Stable Diffusion Web UI Automatic1111, ComfyUI, and Stable Diffusion Web UI Forge). While there are myriad other methods which perform similar functionalities, we narrow down the necessary components to the set

of our choice. The novelty of our workflow is that we achieve greatly improved image editing quality with this small set of simple, user-friendly steps.

For Flux.1 [dev], our workflow is best implemented in ComfyUI. To implement our workflow in ComfyUI, we refer to the following YouTube and Reddit tutorials:

1. `https://www.youtube.com/watch?v=mI0UWm7BNtQ`

2. `https://www.youtube.com/watch?v=GvZXK4phJmE`

3. `https://www.reddit.com/r/comfyui/comments/1d1x5ek/dual_prompting_with_split_sigmas`

Other than Flux.1 [dev], we thoroughly test our workflow with BoleroMix (Pony) v1.41[22], a Japanese animation-style base model derived from SDXL (Podell et al., 2024). For this one and other SDXL-derived models, the user can implement the workflow by simply enabling relevant add-ons in Stable Diffusion Web UI Automatic1111. We also test our workflow in other models at a smaller scale, where all steps in our workflow similarly work.

We have tested our workflow on one NVIDIA RTX A6000, NVIDIA A100 (80 GB), or NVIDIA GeForce RTX 4090. On A6000, each editing step using Flux.1 [dev] (1024×1024 resolution and 30 diffusion steps) only takes around 30 seconds. The price of a server with a single A6000 can be as low as 0.5\$ per hour[23], equivalent to 0.004\$ per editing step. This is one tenth of the cost of DALL·E 3, which is 0.04\$ per image[24]. Due to the low success rate of DALL·E 3 on hard editing tasks (Section 3.4), the cost gap is only larger in real scenarios.

The workflow, together with the installation instructions for popular Web UIs, has been released in a GitHub repository. Please see the abstract of the paper for the link.

---

[22]`https://civitai.com/models/448716?modelVersionId=629179`
[23]`https://cloud.vast.ai/` (January 2025)
[24]`https://openai.com/api/pricing/`