

Transparent Image Layer Diffusion using Latent Transparency

LVMIN ZHANG, Stanford University, USA
 MANEESH AGRAWALA, Stanford University, USA

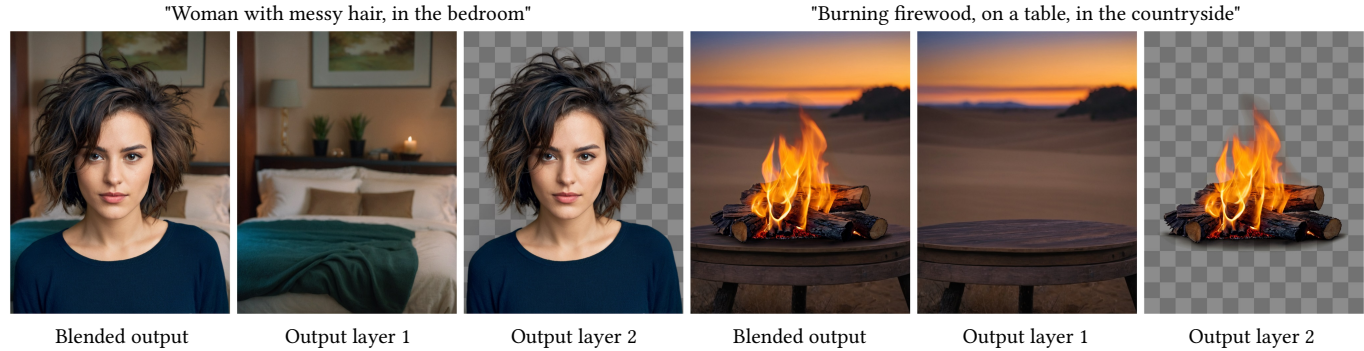


Fig. 1. Generating transparent images and layers. For the given text prompts (top), our framework is capable of generating multiple layers with transparency. These layers can be blended to produce images corresponding to the prompts. Zoom in to see details including messy hair and semi-transparent fire.

We present an approach enabling large-scale pretrained latent diffusion models to generate transparent images. The method allows generation of single transparent images or of multiple transparent layers. The method learns a “latent transparency” that encodes alpha channel transparency into the latent manifold of a pretrained latent diffusion model. It preserves the production-ready quality of the large diffusion model by regulating the added transparency as a latent offset with minimal changes to the original latent distribution of the pretrained model. In this way, any latent diffusion model can be converted into a transparent image generator by finetuning it with the adjusted latent space. We train the model with 1M transparent image layer pairs collected using a human-in-the-loop collection scheme. We show that latent transparency can be applied to different open source image generators, or be adapted to various conditional control systems to achieve applications like foreground/background-conditioned layer generation, joint layer generation, structural control of layer contents, *etc.* A user study finds that in most cases (97%) users prefer our natively generated transparent content over previous ad-hoc solutions such as generating and then matting. Users also report the quality of our generated transparent images is comparable to real commercial transparent assets like Adobe Stock.

CCS Concepts: • **Applied computing** → **Fine arts; Media arts.**

Additional Key Words and Phrases: Transparent images, image editing, image layer, text-to-image diffusion

Authors’ addresses: Lvmin Zhang, lvmin@stanford.edu, Stanford University, USA; Maneesh Agrawala, maneesh@cs.stanford.edu, Stanford University, USA.

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

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.
 ACM 0730-0301/2024/7-ART100
<https://doi.org/10.1145/3658150>

ACM Reference Format:

Lvmin Zhang and Maneesh Agrawala. 2024. Transparent Image Layer Diffusion using Latent Transparency. *ACM Trans. Graph.* 43, 4, Article 100 (July 2024), 38 pages. <https://doi.org/10.1145/3658150>

1 INTRODUCTION

While large-scale models for generating images have become foundational in computer vision and graphics, surprisingly little research attention has been given to layered content generation or transparent image generation. This situation is in stark contrast to substantial market demand. The vast majority of visual content editing software and workflows are layer-based, relying heavily on transparent or layered elements to compose and create content.

The primary factors contributing to this research gap are the lack of training data and the difficulty in manipulating the data representation of existing large-scale image generators. High-quality transparent image elements on the Internet are typically hosted by commercial image stocks with limited (and costly) access, in contrast to text-image datasets that already include billions of images (*e.g.*, LAION [Schuhmann et al. 2022]). The largest open-source transparent image datasets are often less than 50K in size (*e.g.*, DIM [Xu et al. 2017] includes 45,500 transparent images). Meanwhile, most open-source image generation models, *e.g.*, Stable Diffusion, are latent diffusion models that are sensitive to their latent space data representations. Even minor changes to the latent distribution could severely degrade inference or finetuning. For instance, Stable Diffusion 1.5 and XL use different latent spaces, and finetuning with mismatched latents can cause significant degradation in output image quality [Stability 2022b]. This adds to the challenge of manipulating the data representation of existing models to support additional formats like transparent images.

We present a “latent transparency” approach that enables large-scale pretrained latent diffusion models to generate transparent images as well as multiple transparent layers. This method encodes image transparency into a latent offset that is explicitly regulated

to avoid disrupting the latent distribution. The latent transparency is encoded and decoded by external independent models, ensuring that the original pretrained latent encoder/decoder is preserved, so as to maintain high-quality results of state-of-the-art diffusion models. To generate multiple layers together, we use a shared attention mechanism that ensures consistency and harmonious blending between image layers, and we train LoRAs to adapt the models to different layer conditions.

We employ a human-in-the-loop scheme to train our framework and collect data simultaneously. We finalize the scale of our dataset at 1M transparent images, covering a diversity of content topics and styles. We then use state-of-the-art methods to extend the dataset to multi-layer samples. This dataset not only enables the training of transparent image generators but can also be used in different applications like background/foreground-conditioned generation, structure-guided generation, style transfer, *etc.*

Experiments show that in a majority of cases (97%), users prefer the transparent content generated natively by our method over previous ad-hoc solutions like generating-then-matting. When we compare the quality of our generated results with the search results from commercial transparent assets sites like Adobe Stock, user preference rates suggest that quality is comparable.

In summary, we (1) propose “latent transparency”, an approach to enable large-scale pretrained latent diffusion models to generate single transparent images or multiple transparent layers, (2) we present a shared attention mechanism to generate layers with consistent and harmonious blending, and (3) we present a pretrained model for transparent image generation, two pretrained LoRAs for multiple layer generation, as well as several additional ablative architectures for multi-layer generation.

2 RELATED WORK

2.1 Hiding Images inside Perturbations

Research in multiple fields point out a phenomenon: neural networks have the ability to “hide” features in perturbations inside existing features without changing the overall feature distributions, *e.g.*, hiding an image inside another image through small, invisible pixel perturbations. A typical CycleGAN [Zhu et al. 2017] experiment showcases *face-to-ramen*, where the human face identity could be hidden in a picture of ramen. Similarly, invertible downscaling [Xiao et al. 2020] and invertible grayscale [Xia et al. 2018] indicate that neural networks can hide a large image inside a smaller one, or hide a colorful image inside a grayscale one, and then reconstruct the original image. In another widely verified experiment Goodfellow et al. [2015] show that adversarial example signals can be hidden inside feature perturbations to influence the behaviors of other neural networks. In this paper, our proposed “latent transparency” utilizes similar principles: hiding image transparency features inside a small perturbation added to the latent space of Stable Diffusion [Stability 2022a], while at the same time avoiding changes to the overall distribution of the latent space.

2.2 Diffusion Probabilistic Models and Latent Diffusion

Diffusion Probabilistic Model [Sohl-Dickstein et al. 2015] and related training and sampling methods like Denoising Diffusion Probabilistic Model (DDPM) [Ho et al. 2020], Denoising Diffusion Implicit Model (DDIM) [Song et al. 2021], and score-based diffusion [Song et al. 2020] contribute to the foundations of recent large-scale image generators. Early image diffusion methods usually directly use pixel colors as training data [Kong and Ping 2021; San-Roman et al. 2021; Song et al. 2021]. In contrast, the Latent Diffusion Model (LDM) [Rombach et al. 2022] operates in latent space and has been shown to enable easier training while lowering computation requirements. This method has been further extended to create Stable Diffusion [Stability 2022a]. Recently, eDiff-I [Balaji et al. 2022] has used an ensemble of multiple conditions including a T5 text encoder [Raffel et al. 2019], a CLIP text and image embedding encoder [Ilharco et al. 2021]. Versatile Diffusion [Xu et al. 2022] adopts a multi-purpose diffusion framework to process text, an image, and variations within a single model.

2.3 Customized Diffusion Models and Image Editing

Early methods to customize diffusion models have focused on text-guidance [Avrahami et al. 2022; Kim et al. 2022a; Nichol et al. 2021]. Image diffusion algorithms also naturally support inpainting [Avrahami et al. 2022; Ramesh et al. 2022]. Textual Inversion [Gal et al. 2022] and DreamBooth [Ruiz et al. 2022] can personalize the contents of generated results based on a small set of exemplar images of the same topic or object. Recently, control models have also been used to add additional conditions for the generation of text-to-image models, *e.g.*, ControlNet [Zhang and Agrawala 2023], lightweight T2I-adapter [Mou et al. 2023], *etc.* IP-Adapter [Ye et al. 2023] uses a cross-attention mechanism to separate text and image features, allowing for the control signals of the reference image as a visual prompt. [Li et al. 2023a] uses masks in neural network features to achieve semantic region control. Inversion-based methods are also popular in editing images. The DDPM [Ho et al. 2020] theory indicates that a diffusion algorithm constructs data with accumulated small variations and those variations, conditioned on noise, can be manipulated with inverted optimization. Mokady *et al.* [2023] shows that DDIM inversion can optimize images without requiring inputs to be generated by a previously known diffusion process (null-text embedding). Cao *et al.* [2023] and Narek *et al.* [2023] manipulate spatial cross-attention features of Stable Diffusion layers together with DDPM inversion. Hertz *et al.* [2023] edit attention activations of the input images with user-given text prompts and feed them back to the diffusion models. DiffEdit [2023] generates region masks for image editing, given input images and user prompts. Diffusion-CLIP [Kim et al. 2022b] finetunes diffusion models with CLIP loss against prompts. Imagic [Kawar et al. 2023] jointly optimizes text embedding of user prompts and the model gradients to reconstruct the image for image editing applications.

2.4 Transparent Image Layer Processing

Transparent image processing is closely related to image decomposition, layer extraction, color palette processing, as well as image

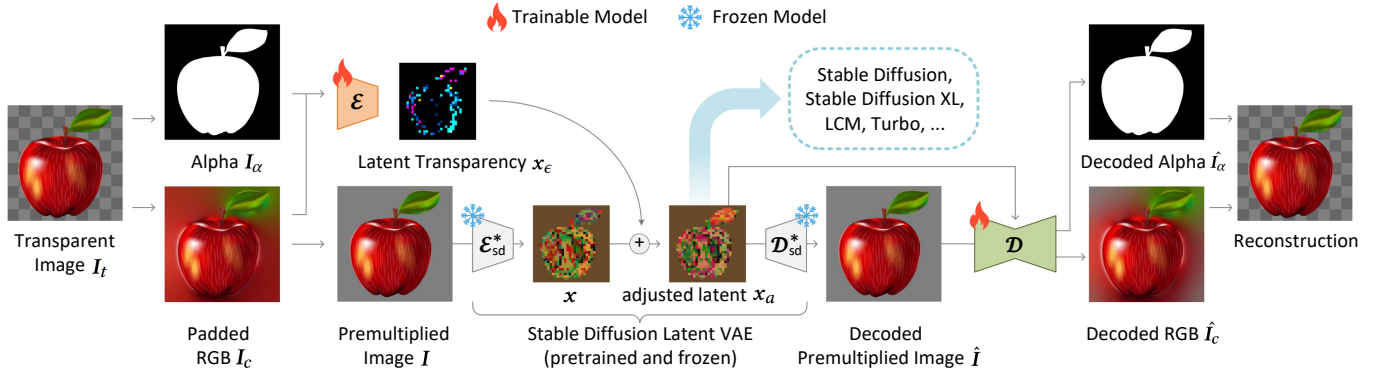


Fig. 2. **Latent Transparency.** Given an input transparent image, our framework encode a “latent transparency” to adjust the latent space of Stable Diffusion. The adjusted latent images can be decoded to reconstruct the color and alpha. This latent space with transparency can be further used in training or fine-tuning pretrained image diffusion models.

matting [Aksoy et al. 2017a, 2016; Tang et al. 2019]. Typical color-based decomposition can be viewed as a RGB color space geometry problem [Du et al. 2023; Tan et al. 2019, 2015, 2018, 2016]. These ideas have also been extended to more advanced blending of image layers [Koyama and Goto 2018]. Unmixing-based color separation also contributes to image decomposition [Aksoy et al. 2017b], and semantic features can be used in image soft segmentation [Aksoy et al. 2018]. We compare our approach to several state-of-the-art deep-learning based matting methods in our experiments and discussion. *PPMatting* [Chen et al. 2022] is a neural network image matting model trained from scratch using standard matting datasets. *Matting Anything* [Li et al. 2023b] is a image matting model using the Segment Anything Model (SAM) [Kirillov et al. 2023] as a backbone. *VitMatte* [Yao et al. 2024] is a tri-map-based matting method using a Vision Transformer (ViT). Text2Layer [Zhang et al. 2023] attempts to use foreground segmentation guidance to achieved layered effects in diffusion models, and indicates that its main bottleneck is the quality of foreground matting method since its learning objective is constructed from the image segmentation of matting models. Our approach starts from native generation of transparent images rather than post-processing of image matting, and is fundamentally different from previous approaches that use matting as post-processing of model outputs or use matting for dataset synthesizing.

2.5 Image Harmonization

Harmonious blending of transparent image layers is closely related to image harmonization research. Achieving “harmony” is usually seen as a problem of correlating color, contrast, and style constituents between foreground and background to ensure natural appearance and consistent composition. Deep learning approaches [Chen et al. 2023; Guerreiro et al. 2023; Guo et al. 2021; Tan et al. 2023; Tsai et al. 2017; Zhu et al. 2015] have been proposed to harmonize images, using annotated datasets [Cong et al. 2020; Niu et al. 2023]. These works utilize the learning capabilities of neural networks to acquire the prior knowledge of harmonization.

3 METHOD

Our approach enables a Latent Diffusion Model (LDM), like Stable Diffusion, to generate transparent images, and then extends the model to jointly generate multiple transparent layers together. In section 3.1, we introduce the method to adjust the LDM latent space to support transparent image encoding/decoding. In section 3.2, we adapt pretrained latent diffusion models with the adjusted latent space to generate transparent images. In section 3.3, we describe the method for joint or conditional layer generating. Finally, we detail the dataset preparation and implementation details for neural network training in section 3.4.

Definitions. To clarify the presentation we first define some terms. For any transparent image $I_t \in \mathbb{R}^{h \times w \times 4}$ with RGBA channels, we denote the first 3 RGB color channels as $I_c \in \mathbb{R}^{h \times w \times 3}$ and the alpha channel as $I_\alpha \in \mathbb{R}^{h \times w \times 1}$. Since the colors are physically undefined at pixels where the alpha value is strictly zero, in this paper, all undefined areas in I_c are always padded by an iterative Gaussian filter (see also supplementary material) to avoid aliasing and unnecessary edge patterns. We call I_c the “padded RGB image” (Fig. 2). The I_t can be converted to a “premultiplied image” as $I = I_c * I_\alpha$ where $*$ denotes pixelwise multiplication. In this paper, all RGB values are in range $[-1, 1]$ (consistent with Stable Diffusion) while all alpha values are in range $[0, 1]$. The premultiplied image I can be seen as a common non-transparent RGB image that can be processed by any RGB-formatted neural networks. Visualizations of these images are shown in Fig. 2.

3.1 Latent Transparency

Our goal is to add transparency support to large-scale latent diffusion models, like Stable Diffusion (SD), that typically uses a latent encoder (VAE) to convert RGB images to latent images before feeding it to a diffusion model. Herein, the VAE and the diffusion model should share the same latent distribution, as any major mismatch can significantly degrade the inference/training/fine-tuning of the latent diffusion framework. When we adjust the latent space to support transparency, the original latent distribution must be preserved

as much as possible. These seemingly conflicting goals (adding transparency support while preserving the original latent distribution) can be handled with a straight-forward measurement: we can check how well the modified latent distribution can be decoded by the original pretrained frozen latent decoder — if decoding a modified latent image creates severe artifacts, the latent distribution is misaligned or broken.

We can write this "harmfulness" measurement mathematically as follows. Given an RGB image I , the pretrained and frozen Stable Diffusion latent encoder $\mathcal{E}_{sd}^*(\cdot)$ and decoder $\mathcal{D}_{sd}^*(\cdot)$, where the $*$ indicates frozen models, we denote the latent image as $\mathbf{x} = \mathcal{E}_{sd}^*(I)$. Assuming this latent image \mathbf{x} is modified by any offset \mathbf{x}_ϵ , produces an adjusted latent $\mathbf{x}_a = \mathbf{x} + \mathbf{x}_\epsilon$. The decoded RGB reconstruction can then be written as $\hat{I} = \mathcal{D}_{sd}^*(\mathbf{x}_a)$ and we can evaluate how "harmful" the offset \mathbf{x}_ϵ is as

$$\mathcal{L}_{\text{identity}} = \|\mathbf{I} - \hat{I}\|_2 = \|\mathbf{I} - \mathcal{D}_{sd}^*(\mathcal{E}_{sd}^*(\mathbf{I}) + \mathbf{x}_\epsilon)\|_2, \quad (1)$$

where $\|\cdot\|_2$ is the L2 norm distance (mean squared error). Intuitively, if $\mathcal{L}_{\text{identity}}$ is relatively high, the \mathbf{x}_ϵ could be harmful and may have destroyed the reconstruction functionality of SD encoder-decoder, otherwise if $\mathcal{L}_{\text{identity}}$ is relatively low, the offset \mathbf{x}_ϵ does not break the latent reconstruction and the modified latent can still be handled by the pretrained Stable Diffusion.

Besides, since most mainstream VAE for diffusion models are KL-Divergence or Diagonal Gaussian Distribution models, these models often has a naively trained parameter for the standard deviation as an offset in the latent space. Considering such deviation denoted as \mathbf{x}_{std} , we can make use of this pretrained parameter to construct $\mathbf{x}_\epsilon = \lambda_{\text{offset}} \mathbf{x}_{\text{std}} \mathbf{x}_{\text{offset}}$ where $\mathbf{x}_{\text{offset}}$ is the raw output from newly added encoder, \mathbf{x}_{std} is the deviation output of pretrained VAE, and λ_{offset} is a weighting parameter with default $\lambda_{\text{offset}} = 1e2$.

We make use of the latent offset \mathbf{x}_ϵ to establish "latent transparency" for encoding/decoding transparent images. More specifically, we train from scratch a latent transparency encoder $\mathcal{E}(\cdot, \cdot)$ that takes the RGB channels \mathbf{I}_c and alpha channel \mathbf{I}_α as input to convert pixel-space transparency into a latent offset

$$\mathbf{x}_\epsilon = \mathcal{E}(\mathbf{I}_c, \mathbf{I}_\alpha). \quad (2)$$

We then train from scratch another latent transparency decoder $\mathcal{D}(\cdot, \cdot)$ that takes the adjusted latent $\mathbf{x}_a = \mathbf{x} + \mathbf{x}_\epsilon$ and the aforementioned RGB reconstruction $\hat{I} = \mathcal{D}_{sd}^*(\mathbf{x}_a)$ to extract the transparent image from the adjusted latent space

$$[\hat{\mathbf{I}}_c, \hat{\mathbf{I}}_\alpha] = \mathcal{D}(\hat{I}, \mathbf{x}_a), \quad (3)$$

where $\hat{\mathbf{I}}_c, \hat{\mathbf{I}}_\alpha$ are the reconstructed color and alpha channels. The neural network layer architecture of $\mathcal{E}(\cdot, \cdot)$ and $\mathcal{D}(\cdot, \cdot)$ is in the supplementary material. We evaluate the reconstruction with

$$\mathcal{L}_{\text{recon}} = \|\mathbf{I}_c - \hat{\mathbf{I}}_c\|_2 + \|\mathbf{I}_\alpha - \hat{\mathbf{I}}_\alpha\|_2, \quad (4)$$

and we experimentally find that the result quality can be further improved by introducing a PatchGAN discriminator loss

$$\mathcal{L}_{\text{disc}} = \mathbb{L}_{\text{disc}}([\hat{\mathbf{I}}_c, \hat{\mathbf{I}}_\alpha]), \quad (5)$$

where $\mathbb{L}_{\text{disc}}(\cdot, \cdot)$ is a GAN objective from a 5-layer patch discriminator (details in supplementary material). The final objective can be jointly written as

$$\mathcal{L}_{\text{vae}} = \lambda_{\text{recon}} \mathcal{L}_{\text{recon}} + \lambda_{\text{identity}} \mathcal{L}_{\text{identity}} + \lambda_{\text{disc}} \mathcal{L}_{\text{disc}}, \quad (6)$$

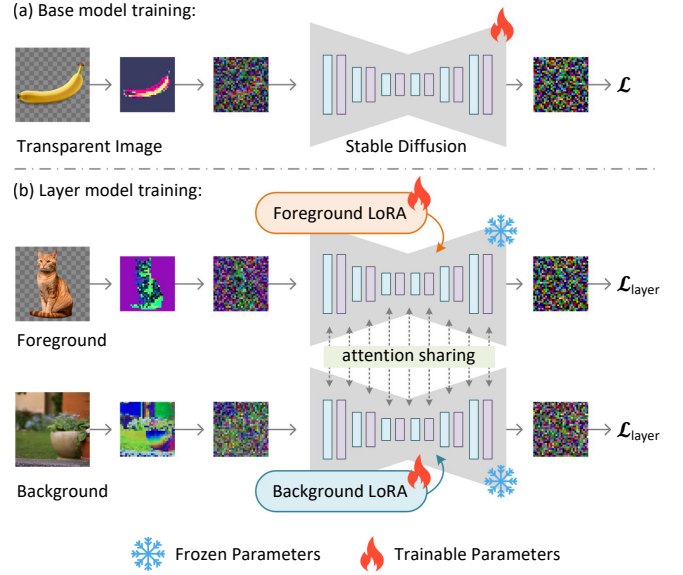


Fig. 3. **Model Training.** We visualize the training of the base model to generate transparent images, and the training of the multi-layer model to generate multiple layers together. When training the base diffusion model (a), all model weights are trainable, whereas for training the multi-layer model (b), only two LoRAs are trainable (the foreground LoRA and background LoRA).

where λ_{\dots} are weighting parameters: by default we use $\lambda_{\text{recon}} = 1$, $\lambda_{\text{identity}} = 1$, $\lambda_{\text{disc}} = 0.01$. By training this framework with \mathcal{L}_{vae} , the adjusted latent \mathbf{x}_a can be encoded from transparent images or vice versa, and those latent images can be used in fine-tuning Stable Diffusion. We visualize the pipeline in Fig. 2.

3.2 Diffusion Model with Latent Transparency

Since the altered latent space with latent transparency is explicitly regulated to align with the original pretrained latent distribution (Eq. 1), Stable Diffusion can be directly fine-tuned on the altered latent space. Given the adjusted latent \mathbf{x}_a , diffusion algorithms progressively add noise to the image and produce a noisy image \mathbf{x}_t , with t denoting how many times noise is added. When t is large enough, the latent image approximates pure noise. Given a set of conditions including the time step t and text prompt \mathbf{c}_t , image diffusion algorithms learn a network ϵ_θ that predicts the noise added to the noisy latent image \mathbf{x}_t with

$$\mathcal{L} = \mathbb{E}_{\mathbf{x}_t, t, \mathbf{c}_t, \epsilon \sim \mathcal{N}(0,1)} \left[\|\epsilon - \epsilon_\theta(\mathbf{x}_t, t, \mathbf{c}_t)\|_2^2 \right] \quad (7)$$

where \mathcal{L} is the overall learning objective of the entire diffusion model. This training is visualized in Fig. 3-(a).

3.3 Generating Multiple Layers

We further extend the base model to a multi-layer model using attention sharing and LoRAs [Hu et al. 2021], as shown in Fig. 3-(b). We denote the foreground noisy latent as \mathbf{x}_f and background as \mathbf{x}_b , and train two LoRAs, a foreground LoRA parameterized by θ_f and a background LoRA by θ_b , to denoise the latent images. If the two

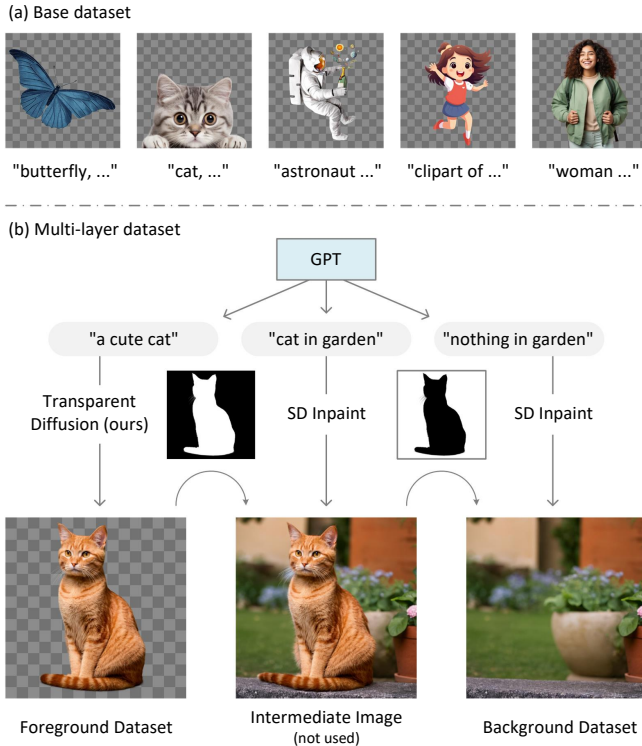


Fig. 4. **Dataset Preparation.** We demonstrate the preparation of the two datasets: the transparent image dataset (base dataset) and multi-layer dataset. The base dataset is collected by downloading online transparent images and a human-in-the-loop training method. The multi-layer dataset is synthesized with our transparent diffusion model and several state-of-the-art models including ChatGPT, SDXL inpaint model, etc. The final scale of each dataset is around 1M.

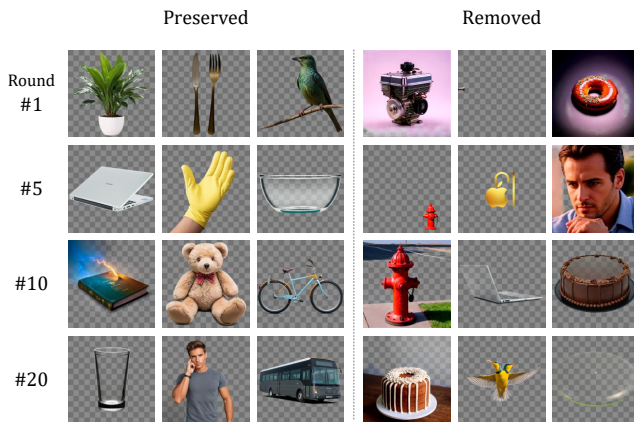


Fig. 5. **Human-in-the-loop data screening.** We visualize sample examples that are preserved versus removed in each round during the dataset collection process. We show examples from the round 1, 5, 10, and 20. The prompts are randomly sampled during the collecting process.

models independently denoise the two images, we have the two objectives with

$$\begin{cases} \mathbb{E}_{\mathbf{x}_f, t, c_t, \epsilon_f \sim \mathcal{N}(0,1)} \left[\|\epsilon_f - \epsilon_{\theta, \theta_f}(\mathbf{x}_f, t, c_t)\|_2^2 \right] \\ \mathbb{E}_{\mathbf{x}_b, t, c_t, \epsilon_b \sim \mathcal{N}(0,1)} \left[\|\epsilon_b - \epsilon_{\theta, \theta_b}(\mathbf{x}_b, t, c_t)\|_2^2 \right] \end{cases} \quad (8)$$

where ϵ_f, ϵ_b are latent noise for the foreground and background. We then merge the two independent diffusion process to achieve coherent generation. For each attention layer in the diffusion model, we concatenate all {key, query, value} vectors activated by the two images, so that the two passes can be merged into a jointly optimized big model $\epsilon_{\theta, \theta_f, \theta_b}(\cdot)$. We denote the merged noise as concatenated $\epsilon_m = [\epsilon_f, \epsilon_b]$, and we have the final objective

$$\mathcal{L}_{\text{layer}} = \mathbb{E}_{\mathbf{x}_f, \mathbf{x}_b, t, c_t, \epsilon_m \sim \mathcal{N}(0,1)} \left[\|\epsilon_m - \epsilon_{\theta, \theta_f, \theta_b}(\mathbf{x}_f, \mathbf{x}_b, t, c_t)\|_2^2 \right] \quad (9)$$

to coherently generate multiple layers together. We can also make simple modifications to this objective to support conditional layer generation (e.g., foreground-conditioned background generation or background-conditioned foreground generation). More specifically, by using a clean latent for the foreground instead of noisy latent (i.e., by always setting $\epsilon_f = 0$), the model will not denoise foreground, and the framework becomes a foreground-conditioned generator. Similarly, by setting $\epsilon_b = 0$, the framework becomes a background-conditioned generator. We implement all these conditional variations in experiments.

3.4 Dataset Preparation and Training Details

Base Dataset. We use a human-in-the-loop method to collect a dataset of transparent images and train our models. The dataset initially contains 20k high-quality transparent PNG images purchased or downloaded free from 5 online image stocks (all images include commercial use permission (examples in Fig. 4-(a)). We then train the SDXL VAE with latent transparency using randomly sampled images with equal probability (at batch size 8), and then train the SDXL diffusion model using the same data with adjusted latents. Next we repeat the following steps for a total of 25 rounds. At the beginning of each round, we generate 10k random samples using the last model in the previous round. and the random prompts from LAIONPOP [Schuhmann and Bevan 2023]. We then manually pick 1000 samples to add back to the training dataset. The newly added samples are given a 2x higher probability of appearing in training batches in the next round. We then train the latent transparency encoder-decoder and diffusion models again. After 25 rounds, the size of the dataset increases to 45K. Afterwards, we generate 5M sample pairs without human interaction and use the LAION Aesthetic threshold [Schuhmann et al. 2022] setting of 5.5 and clip score sorting to obtain 1M sample pairs. We automatically remove samples that do not contain any transparent pixels as well as those that do not contain any visible pixels. Finally, all images are captioned with LLaVA [Liu et al. 2023] (an open-source multi-modal GPT similar to GPT4v) to get detailed text prompts. The training of both the VAE and the diffusion model is finalized with another 15k iterations using the final 1M dataset.

Statistical Analysis. We briefly analyze here how human data selection improves the quality of the dataset as well as the model



Fig. 6. **Qualitative Results.** We showcase various examples of transparent images generated by our model. The prompts for each group is given at the top of the examples. These examples only use our base single-layer model.

Table 1. **Statistical Record of Human-in-the-loop Collection.** We report the Defective Sample Count (DSC) per 100 sampling during the rounds of human data selection. We resume the model checkpoints recorded after each round of data collection and generate 100 samples for each checkpoint. Users find how many samples are of obvious defects (like fully empty image, or fully non-transparent image, or obvious errors like opaque glass, etc) and report the number as *DSC* (lower is better ↓).

Round	#0	#1	#2	#3	#4	#5	#6	#7	#8	#9	#10	#12	#14	#16	#18	#20
DSC ↓	61	62	37	53	41	25	17	15	23	21	25	11	6	9	3	5

capabilities. As shown in Fig. 5, we visualize samples that are preserved or removed in each round of the human-in-the-loop selection. We can see that human efforts removes some obvious flaws (*e.g.*, empty images, fully opaque colors for glass, *etc.*) and enhances the diversity of the dataset content (*e.g.*, the glowing effects on magic books, *etc.*). In Table 1, we resume the checkpoint from each round of data collection to sample images, and ask the users to review the images and count images with obvious defects. We can see that as the number of rounds increases, the rate of defective outputs gradually decreases.

Multi-layer Dataset. We further extend our $\{text, transparent\}$ image dataset into a $\{text, foreground\}$ layer, $\{background\}$ layer dataset, so as to train the multi-layer models. As shown in Fig. 4-(b), we ask GPTs (we used ChatGPT for 100k requests and then moved to LLAMA2 for 900k requests) to generate structured prompts pairs for foreground like “a cute cat”, entire image like “cat in garden”, and background like “nothing in garden” (we ask GPT to add the word “nothing” to the background prompt). The foreground prompt

is processed by our trained transparent image generator (Section 3.2) to obtain the transparent images. Then, we use Diffusers Stable Diffusion XL Inpaint model [diffusers 2024] to inpaint all pixels with alpha less than one to obtain intermediate images using the prompt for the entire images. Finally, we invert the alpha mask, erode $k = 8$ pixels and inpaint again with the background prompt to get the background layer. We repeat this process 1M times to generate 1M layer pairs.

Training Details. We use the AdamW optimizer at learning rate 1e-5 for both VAE and diffusion model. The pretrained Stable Diffusion model is SDXL [Podell et al. 2023]. For the the LoRA [Hu et al. 2021] training, we always use rank 256 for all layers. We use the Diffusers’ standard for naming and extracting LoRA keys. In the human-in-the-loop data collection, each round contains 10k iterations at batch size 16. The training devices are 4x A100 80G NV-link, and the entire training takes one week (to reduce budget, the training is paused when human are collection data for the next round of optimization) and the real GPU time is about 350 A100 hours. Our approach is training friendly for personal-scale or lab-scale research as the 350 GPU hours can often be processed within 1K USD.

4 EXPERIMENTS

We detail qualitative and quantitative experiments with our system. We first present qualitative results with single images (Section 4.1), multiple layers (Section 4.2), as well as iterative generation (Section 4.3), and then show that our framework can also be combined with control modules for wider applications (Section 4.4). We then analysis the importance of each component with ablative study

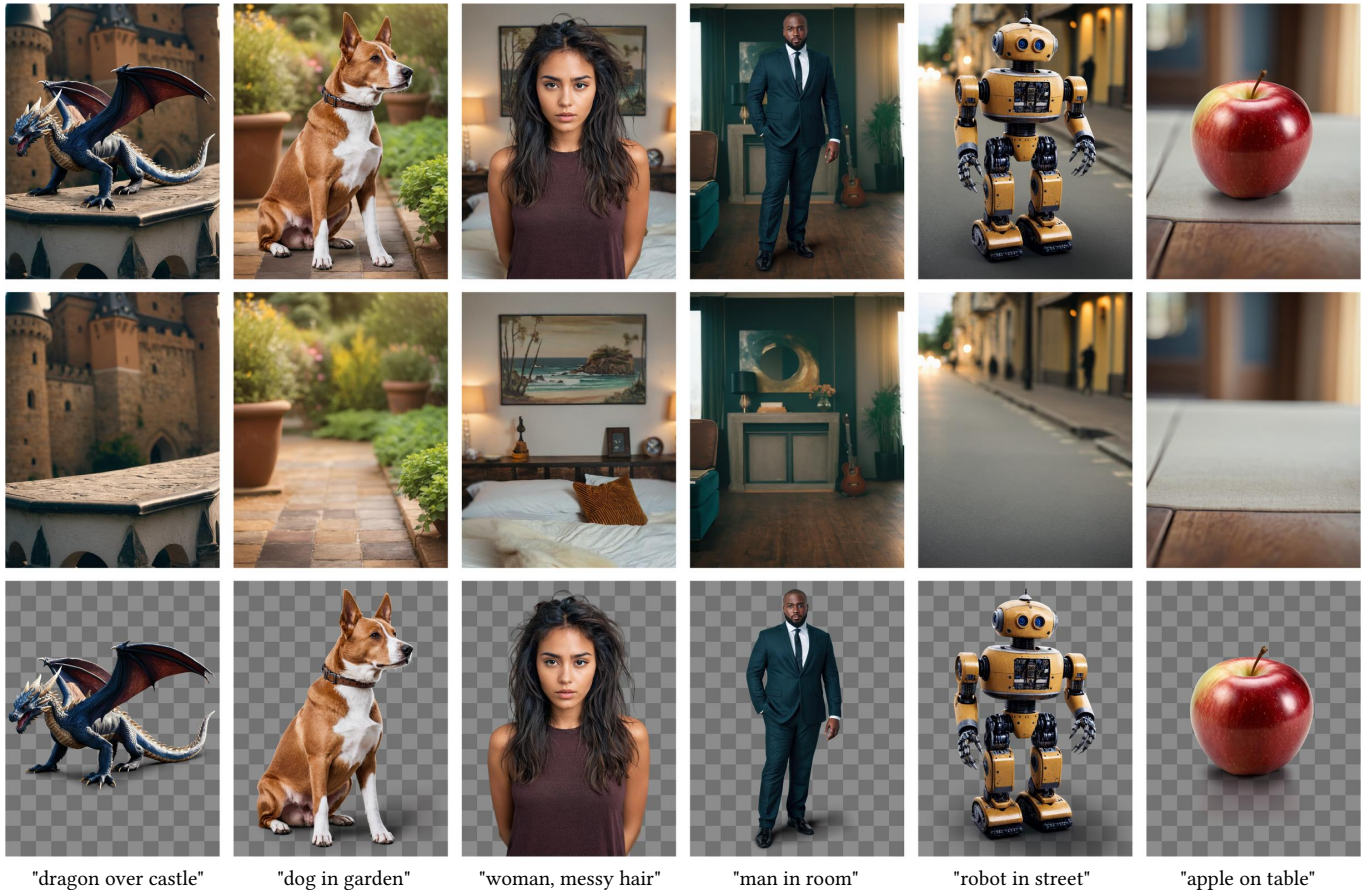


Fig. 7. **Multi-Layer Qualitative Results.** We presents qualitative results generated by our model using prompts with diverse topics. For each example, we show the blended image, and two output layers. More results are available in supplementary materials.

(Section 4.5), and then discuss the difference and connection between our approach and image matting (Section 4.6). Finally, we conduct perceptual user study (Section 4.7) and present a range of discussions to further study the behaviors of our framework (Section 4.8,4.10,4.12).

4.1 Qualitative Results

We present qualitative results in Fig. 6 with a diverse set of transparent images generated using our single-image base model. These results showcase the model’s capability to generate *natively* transparent images that yield high-quality glass transparency, hair, fur, and semi-transparent effects like glowing light, fire, magic effect, *etc.* These results also demonstrate the model’s capability to generalize to diverse content topics.

We further present multi-layer results in Fig. 7 with transparent layers generated by our multi-layer model and the blended images. These results showcase the model’s capability to generate harmonious compositions of objects that can be blended together seamlessly. The layers are not only consistent with respect to illumination and geometric relationships, but also demonstrate the aesthetic quality of Stable Diffusion (*e.g.*, the color choice of the

background and foreground follows a learned distribution that looks harmonious and aesthetic).

4.2 Conditional Layer Generation

We present conditional layer generation results (*i.e.*, foreground-conditioned background and background-conditioned foreground generation) in Fig. 8. We can see that the model is able to generate consistent composition with coherent geometry and illumination. In the “bulb in the church” example, the model tries to generate an aesthetic symmetric design to match the foreground. The “sitting on bench”/“sitting on sofa” examples demonstrate that the model is able to infer the interaction between foreground and background and generate corresponding geometry.

4.3 Iterative Generation

Fig. 9 shows that we can iteratively use the background-conditioned foreground generation model to achieve composition or arbitrary number of layers. For each new layer, we blend all previously generated layers into one RGB image and feed it to the background-conditioned foreground model. We also observe that the model is able to interpret natural language in the context of the background

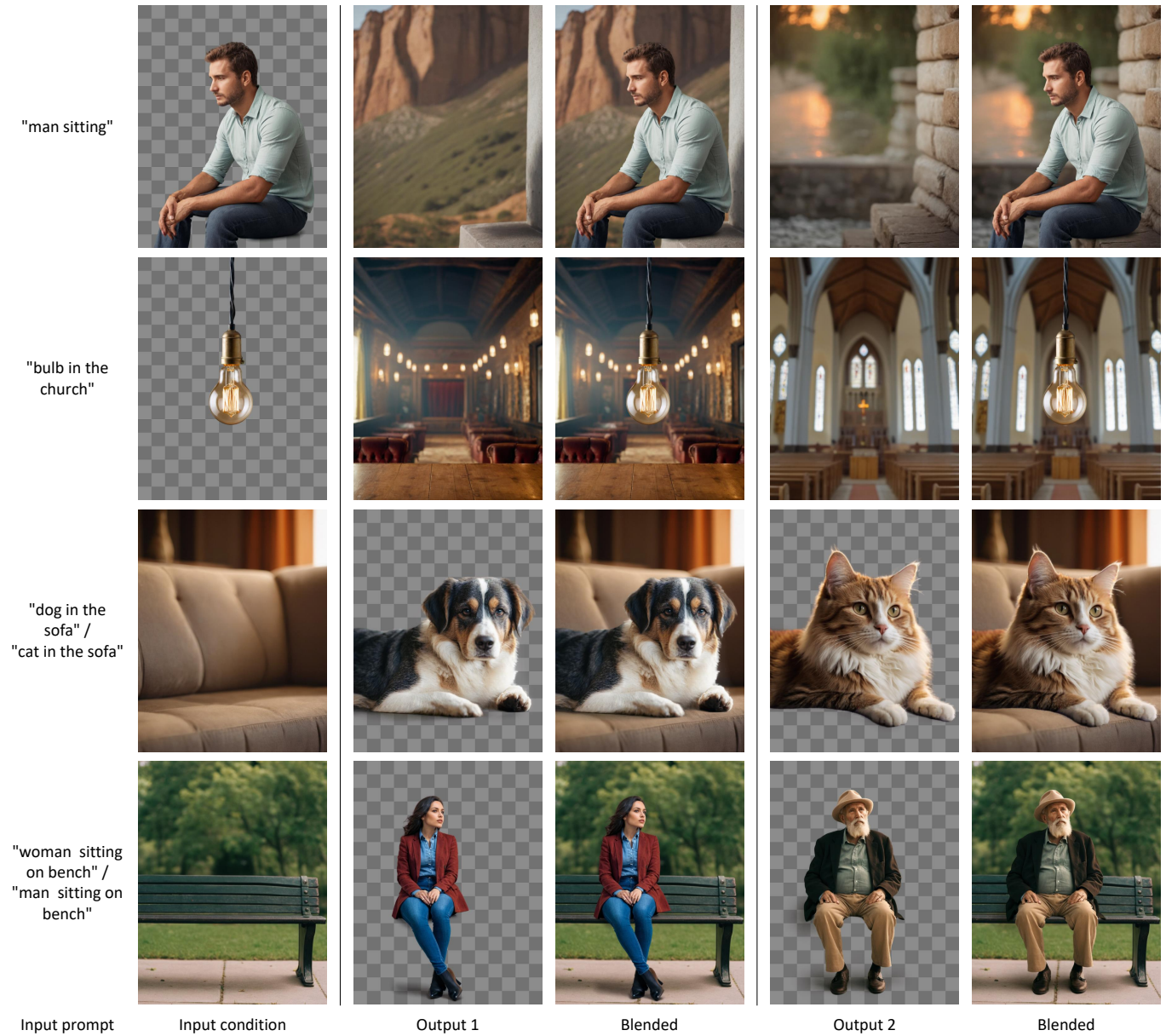


Fig. 8. **Conditional Layer Generating.** We presents results with foreground-conditioned background (the first two rows) and background-conditioned foreground (the last two rows). For each example, we generate two foregrounds/backgrounds.

image, *e.g.*, generating a book in front of the cat. The model displays strong geometric composition capabilities, *e.g.*, composing a human sitting on a box.

4.4 Controllable Generation

As shown in Fig. 10, we demonstrate that existing control models like ControlNet [Zhang and Agrawala 2023] can be applied to our model for enriched functionality. We can see that the model is able to preserve the global structure according to the ControlNet signal to generate harmonious compositions with consistent illumination

effects. We also use a “reflective ball” example to show that the model is able to interact with the content of the foreground and background to generate consistent illumination like the reflections.

4.5 Ablative Study

We conduct an ablative study to evaluate the contribution of each component in our framework. We are interested in a possible architecture that does not modify Stable Diffusion’s latent VAE encoder/decoder, but only adds channels to the UNet. In the original Stable Diffusion, a $512 \times 512 \times 3$ image is encoded to a latent image

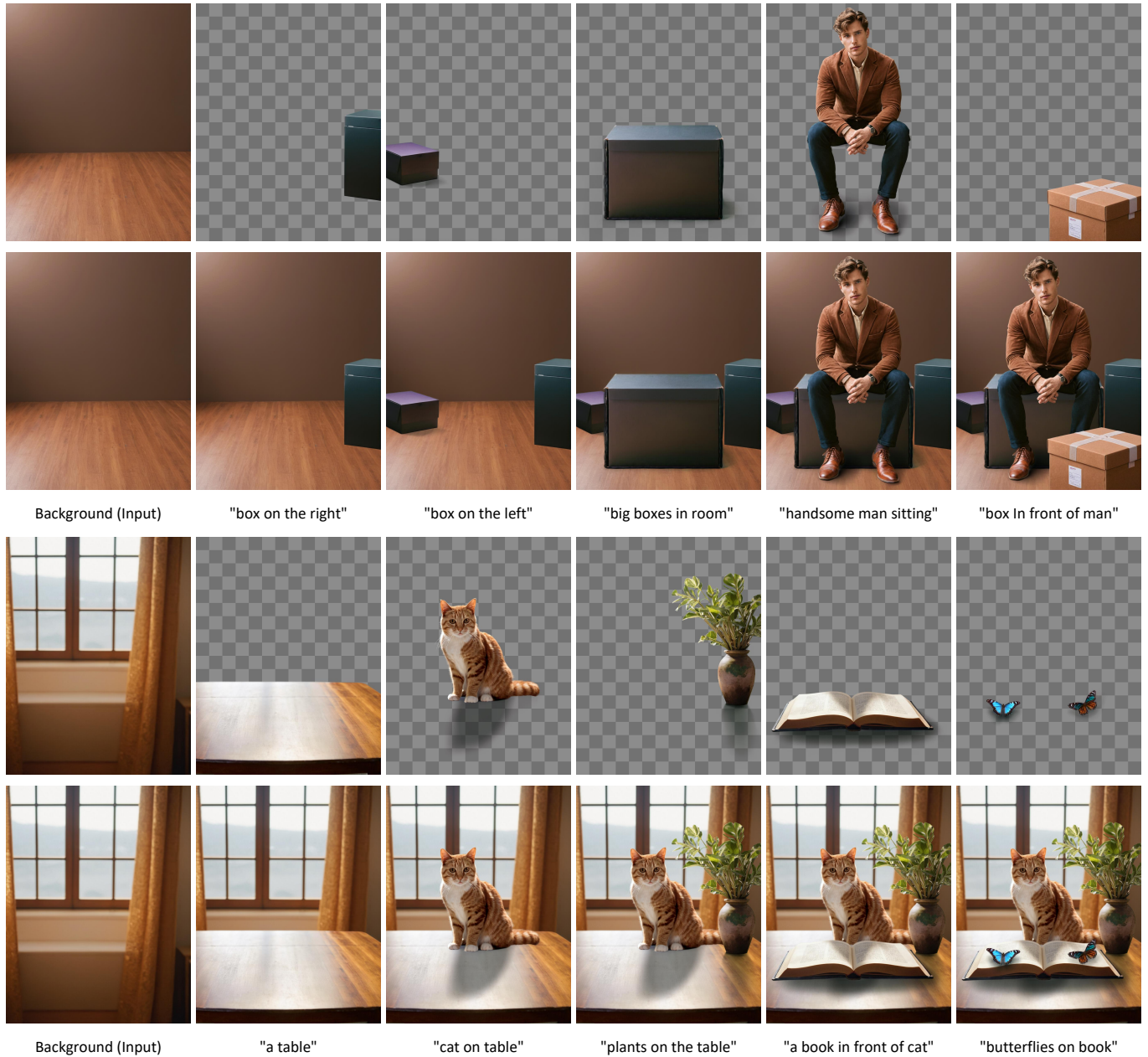


Fig. 9. **Generating Multiple Layers.** We show that our framework can compose multiple layers iteratively, by repeating the background-conditioned foreground model. At each step, we blend all existing layers and feed the blended result to the background-conditioned generator. The prompts at each step is at the bottom of outputs.

of size $64 \times 64 \times 4$. This indicates that if we duplicate the $512 \times 512 \times 1$ alpha channel 3 times into a $512 \times 512 \times 3$ matrix, the alpha could be directly encoded into a $64 \times 64 \times 4$ latent image. By concatenating this with the original latent image, the final latent image would form a $64 \times 64 \times 8$ matrix. This means we could add 4 channels to Stable Diffusion UNet to force it support an alpha channel. We present the results of this approach in Fig. 11-(a). We can see that this method

severely degrades the generation quality of the pretrained large model, because its latent distribution is changed; although the VAE is unchanged (it is frozen), the additional 4 channels significantly change the feature distribution after the first convolution layer in the VAE UNet. Note that this is different from adding a control signal to the UNet — the UNet must generate and recognize the added channels all at the same time because diffusion is an iterative process,



Fig. 10. **Combining with Control Models.** We show that our approach can directly be combined with control models like ControlNet [Zhang and Agrawala 2023] to enhance the functionality. The prompts are “human in street”, “human in forest”, “big reflective ball in street”, and “big reflective ball in forest”.

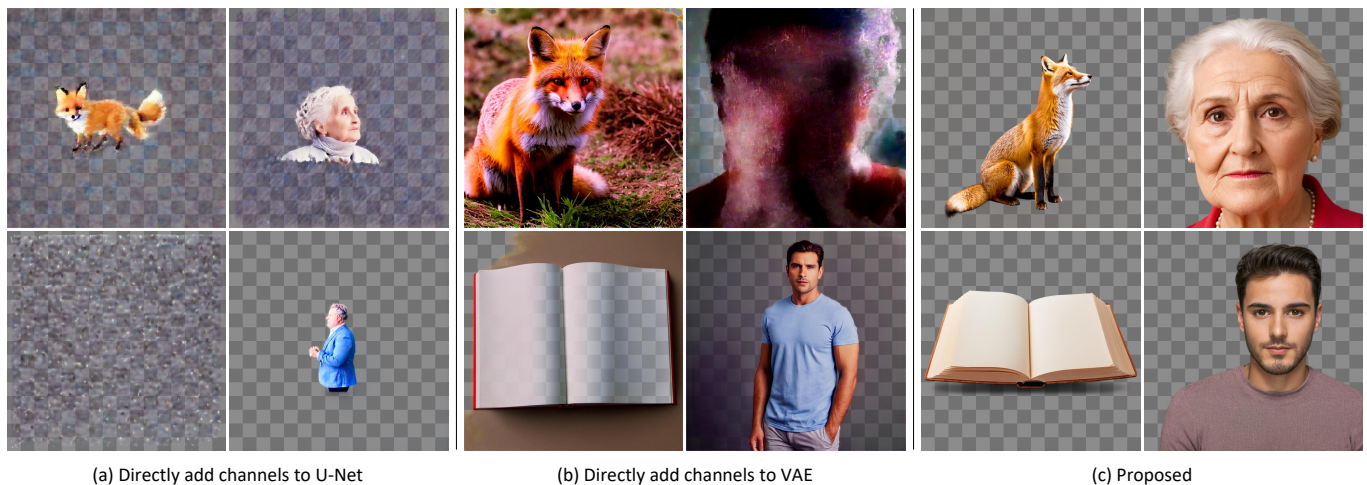


Fig. 11. **Ablative Study.** We compare our approach to two alternative architecture: directly adding channels to UNet and directly adding channels to VAE. When adding channel to UNet, we directly encode alpha channel as an external image and add 4 channels to UNet. When adding channels to VAE, the UNet is finetuned on the latent images encoded by the newer VAE. The test prompts are “fox”, “elder woman”, “a book”, “man”.

and the outputs of any diffusion step become the input of the next diffusion step.

In Fig. 11-(b), we test another architecture that directly adds a channel to the VAE encoder and decoder. We train the VAE to include an alpha channel, and then further train the UNet. We observe that such training is very unstable, and the results suffer from different types of collapse from time to time. The essential reason leading to this phenomenon is that the latent distribution is changed too much during in the VAE fine-tuning.

We also introduce several alternative architectures in Fig. 12 for more complicated workflows. We can add zero-initialized channels to the UNet and use the VAE (with or without latent transparency)

to encode the foreground, or background, or layer combinations into conditions, and train the model to generate foreground or background or directly generate blended images (e.g., Fig. 12-(a, b, c)). We visualize examples of this two-stage pipeline in Fig. 12-(d, e).

4.6 Relationship to Image Matting

We discuss the difference and connection between native transparent image generation and image matting. To be specific, we test the following matting methods: (1) *PPMatting* [Chen et al. 2022] is a state-of-the-art neural network image matting model. This model reports to achieve the highest precision among all “classic” neural network based matting methods, *i.e.*, neural models trained from

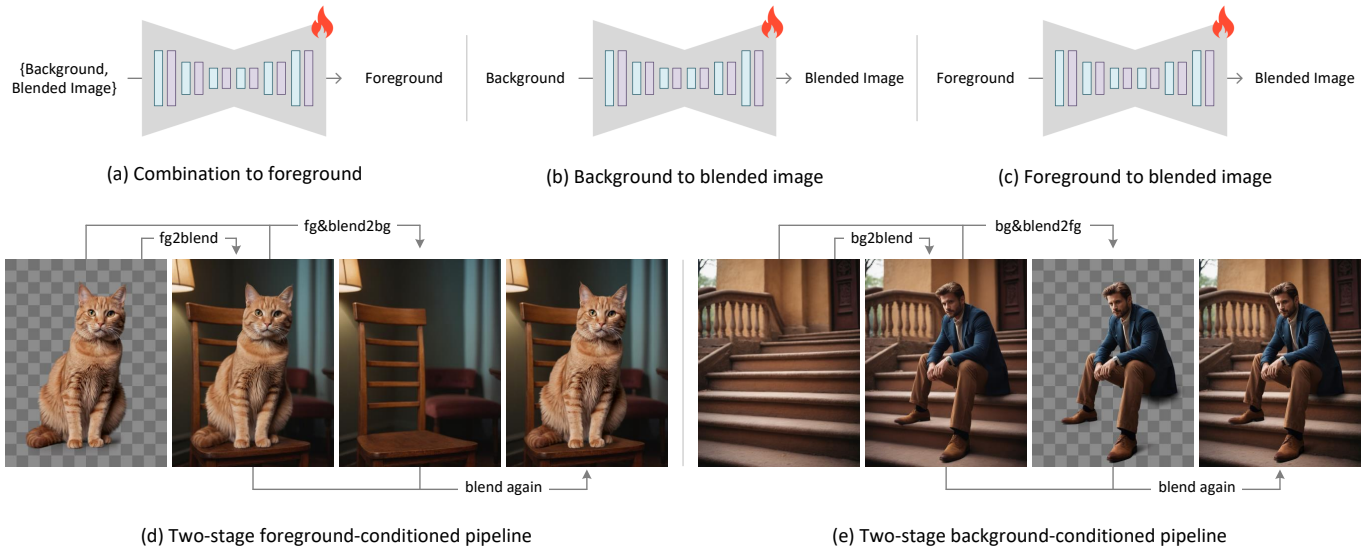


Fig. 12. **Additional Ablative Architectures** We also include several alternative models for more complicated workflows. These include generating a blended image from background or foreground, as well as generating background/foreground from other combined layers. We also demonstrate the use of two-step pipelines to generate/decompose independent layers.

Table 2. **User Study.** We present the results from user study. We conduct user study in two groups: the first group compares outputs between different methods, while the second group directly compare our generated results to the search result of a commercial transparent image assets (Adobe Stock). Higher is better and best in bold.

Candidate	Group 1	Group 2
SD + PPMatting [Chen et al. 2022]	2.1±1.2%	/
SD + Matting Anything [Li et al. 2023b]	0.8±0.5%	/
Ours (base model)	97.1±1.9%	45.3±9.1%
Commercial Transparent Asset Stock	/	54.7±8.3%

scratch on a collected dataset of transparent images. This model is fully automatic and does not need a user-specified tri-map. (2) *Matting Anything* [Li et al. 2023b] is a new type of image matting model based on the recently released Segment Anything Model (SAM) [Kirillov et al. 2023]. This model uses pretrained SAM as a base and finetunes it to perform matting. This model also does not need a user-specified tri-map. We also include a tri-map-based method to study the potential for user-guided matte extraction. (3) *VitMatte* [Yao et al. 2024] is a state-of-the-art matting model that uses tri-maps. The architecture is a Vision Transformer (ViT) and represents the highest quality of current user-guided matting models.

As shown in Fig. 13, we can see that several types of patterns are difficult for matting approaches, e.g., semi-transparent effects like fire, pure white fur against a pure white background, shadow separation, etc. For semi-transparent contents like fire and shadows, once these patterns are blended with complicated background, separating them becomes a nearly impossible task. To obtain perfectly

clean elements, probably the only method is to synthesize elements from scratch, using a native transparent layer generator. We further notice the potential to use outputs of our framework to train matting models.

4.7 Perceptual User Study

In order to perceptually evaluate and compare our approach with existing methods, we perform a perceptual user study focusing on human aspects of our native transparent results and ad-hoc methods like Stable Diffusion + generation-and-matting. We target real-world use cases where users want to get transparent elements given specific demands (prompts). Our study tests multiple types of methods (native transparent generation, generating-and-matting, online commercial stock) to see how they fulfill such demands (by asking users which they prefer).

Specifically, our user study involves 14 individuals, where 11 individuals are online crowd-source workers, 1 is a computer science student, and the other 2 are professional content creators. We sample 100 results using the 3 methods (prompts are randomly sampled from PickaPic [Kirstain et al. 2023]), and this leads to 100 result groups, with each group containing 3 results from 3 methods. The participants are invited to rank the results in each group. When ranking the results in each group, we ask users the question – “Which of the following results do you prefer most? Please rank the following transparent elements according to your preference”. We use the preference rate as the testing metric. This process is repeated 4 times to compute the standard deviation. Afterwards, we calculate the average preference rate of each method. We call this user study “group 1”.

We compare our approach with SD+PPMatting [Chen et al. 2022], SD+Matting Anything [Li et al. 2023b]. Herein, “SD+” means we first use Stable Diffusion XL to generate an RGB image, and then

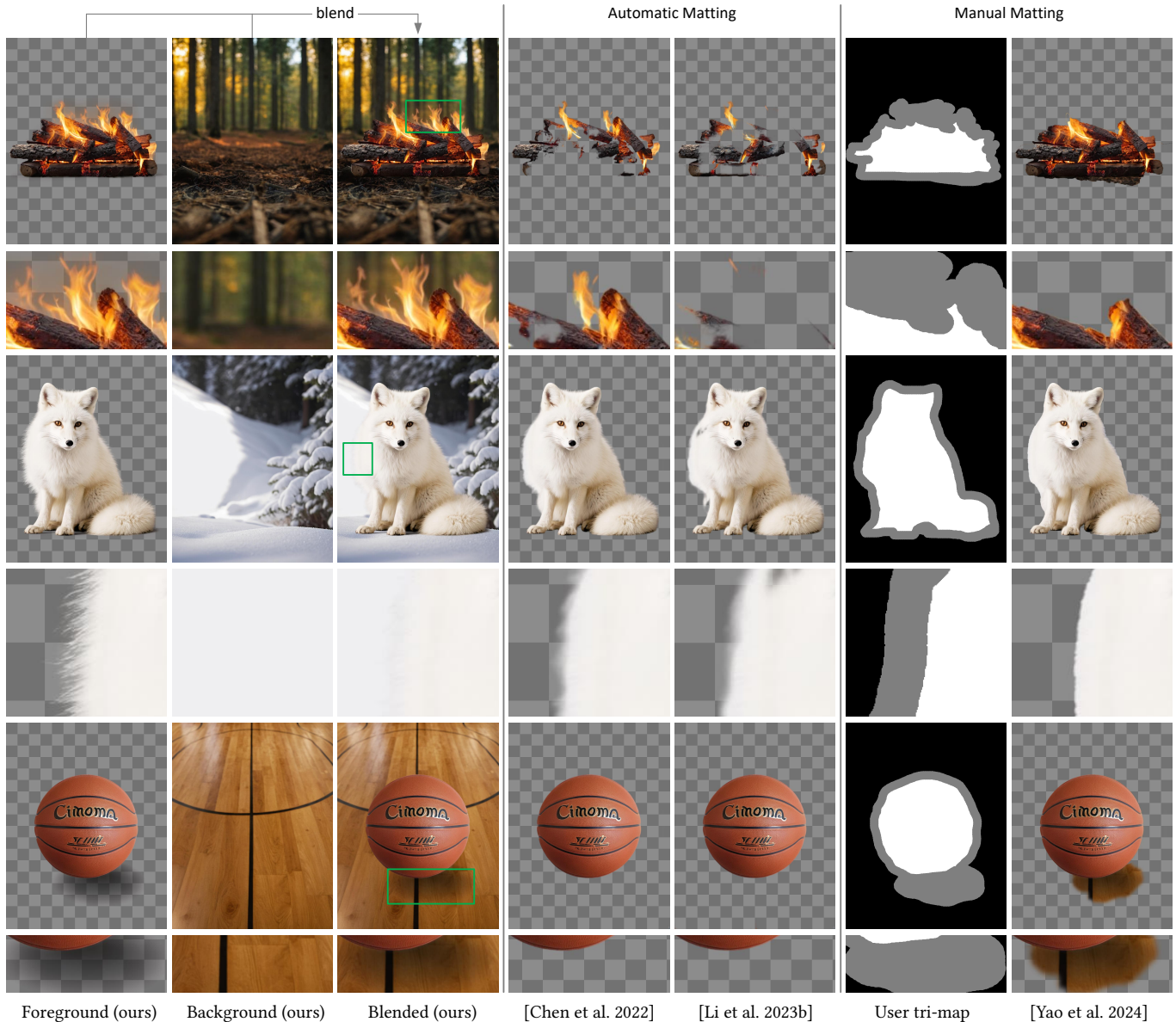


Fig. 13. **Difference between Joint Layer Generating and Generating-then-matting.** This is *not* a result comparison since the left images are outputs layers of our method. The blended images are alpha blending of the generated layers. (Our method does not decompose images.) We try to reproduce similar results using matting approaches. The prompts are “fire on burning wood in forest”, “white fox in white snow ground, all white, very white”, and “basketball”.

perform matting using the corresponding method. Results are shown in Table. 2, group 1. We find that users prefer our approach over all other approaches (in more than 97% cases). This demonstrates the advantage of native transparent image generation over ad-hoc solution like generation-then-matting.

We also perform another user preference experiment in “group 2”, comparing our results against searching for commercial transparent assets from Adobe Stock, using the same aforementioned user preference metric. In Table. 2, group 2, we report that the preference rate of our method is close to commercial stock (45.3% v.s. 54.7%).

Though the high-quality paid content from commercial stock is still preferred marginally. This result suggests that our generated transparent content is competitive to commercial sources that require users to pay for each image.

4.8 Raw RGBA Channels

Fig. 14 shows the raw outputs with each channel in our generated transparent images. We can see that the model avoids aliasing by padding the RGB channel with smooth “bleeding” colors. This



Fig. 14. **Raw outputs of the RGB channels and alpha channel.** We present the raw RGB and alpha channel for evaluation. The prompts are “woman with messy hair”, “boy with messy hair”, and “glass cup”.

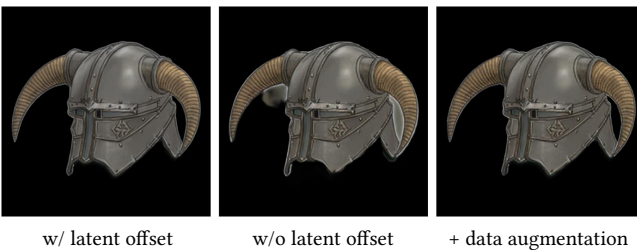


Fig. 15. **Robust decoder with data augmentations.** We show that it is possible to use data augmentation methods to train a robust decoder to handle situations when the UNet cannot diffuse the desired latent offsets. Samples are transparent images on black backgrounds.

approach ensure high-quality foreground color in areas of alpha blending.

4.9 Robust Decoder with Data Augmentations

Fig. 15 shows that we can use data augmentation methods to achieve more robust decoder to handle situations when the latent offset is missing or wrong. This can be useful when certain community

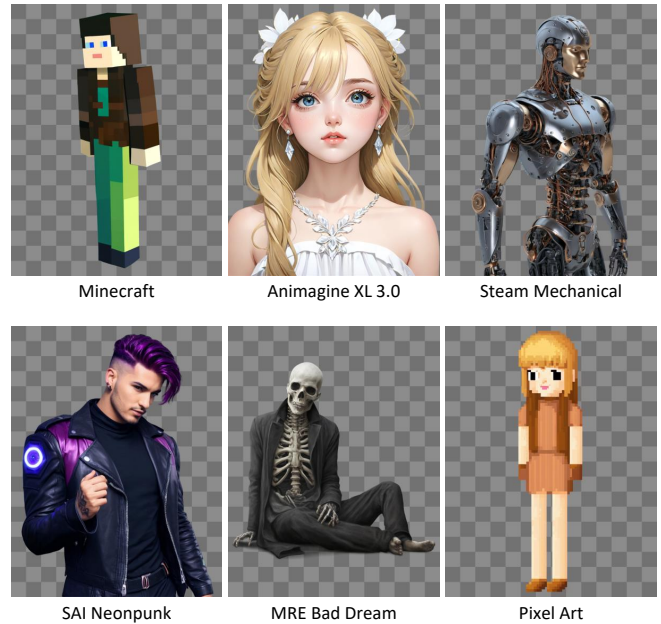


Fig. 16. **Applying to Community Models.** We show that our model can be applied to community LoRAs/Models/PromptStyles to achieve diverse results. All images are achieved using prompt “person”, excepting Animate XL using “1girl, masterpiece, fantastic art”.

models (e.g., anime, cartoon, etc.) fail to produce the desired latent offset during the diffusion process, and we still want to decode useful transparent images from those slightly mismatched latent spaces yielded by those fine-tuned models. To be specific, Fig. 15 simply dropout 30% offsets when training the decoder.

4.10 Community Models

As shown in Fig. 16, our method can be applied to various community models, LoRAs, and prompt styles, without additional training. More specifically, we try a Minecraft LoRA, a pixel art LoRA, an anime model [cagliostrolab 2024], and several community prompt styles. We can see that applying to different models neither degrades the quality of target model/LoRAs nor degrades the quality of image transparency. This integration capability suggests the model potential for wider use in diverse creative and professional domains.

4.11 Inference Speed

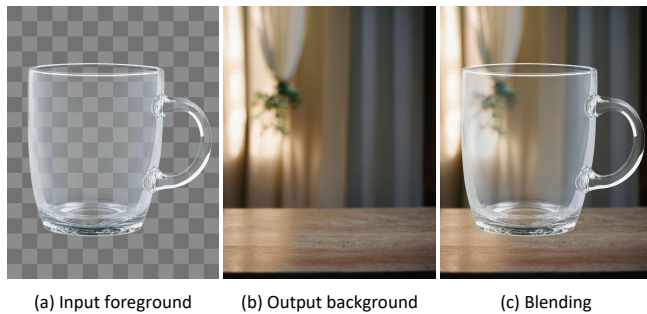
In Table 3, we report the inference speed with different base diffusion models and architectures. All tests are based on personal level computation devices. We tested SD1.5 and SDXL with Nvidia RTX 3070 and RTX 4090. Our tests include generating single transparent images, generating multiple layers jointly, and generating multiple layers using the two-stage pipelines.

4.12 Limitations

As shown in Fig. 17, one trade-off with our framework is between generating “clean transparent elements” and “harmonious blending”. For instance, if the transparent image is a clean and reusable

Table 3. **Inference Speed.** We report on the inference speed of our framework under different base diffusion models and architectures. We test with a Nvidia RTX 3070 device and a RTX 4090 device. All models use 30 diffusion sampling steps. The reported data are averages of 64 runs. All two-stage model needs to be run twice for two layers, doubling the inference time. This speed may be affected by different inference software.

Candidate	SD 1.5	SDXL
Only transparent image (RTX4090)	1.85s	7.3s
Only transparent image (RTX3070)	4.13s	12.5s
Two layers (joint, RTX4090)	5.71s	21.5s
Two layers (joint, RTX3070)	13.01s	41.28s
Two layers (two-stage method, RTX4090)	$1.92s \times 2$	$4.37s \times 2$
Two layers (two-stage method, RTX3070)	$4.77s \times 2$	$14.71s \times 2$



(d) Background-conditioned foreground

Fig. 17. **Limitation.** The prompt in this example is “glass cup on table in a warm room”. If the input foreground is a clean transparent object without any illumination or shadow effects, harmonious blending is very difficult since the alpha blending does not create deformation of light or casting of shadows. This can be resolved to some extent when using the background as a condition to generate the foreground. But in this case, getting a clean and reusable transparent object without the influence of illumination is difficult.

element without any special illumination or shadow effects, generating a background that can be harmoniously blended with the foreground can be very challenging and the model may not succeed in every cases (Fig. 17-(c) is a failure case). This phenomenon can be cured to some extent if we only use backgrounds as conditions to generate foregrounds to force a harmonious blending (Fig. 17-(d)). Nevertheless, this will also lead to illumination influencing the transparent object, making the transparent objects less reusable. One may argue that the image in Fig. 17-(a) is much more reusable for designers and in-the-wild applications than the transparent images in Fig. 17-(d) which contain many specific patterns bound to the background.

5 CONCLUSION

In summary, this paper introduces “latent transparency”, an approach to create either individual transparent images or a series of coherent transparent layers. The method encodes the transparent alpha channel into the latent distribution of Stable Diffusion. This process ensures that the high-quality output of large-scale image diffusion models, by regulating an offset added to the latent space. The training of the models involved 1M pairs of transparent image layers, gathered using a human-in-the-loop collection scheme. We present a range of applications, such as generating layers conditioned on foreground/background, combining layers, structure-controlled layer generating, *etc.* User study results indicate that in a vast majority of cases, users favor the transparent content produced natively by our method over traditional methods like generation-then-matting. The quality of the transparent images generated was found to be comparable to the assets in commercial stocks.

ACKNOWLEDGMENTS

This work was partially supported by Google through their affiliation with Stanford Institute for Human-centered Artificial Intelligence (HAI).

REFERENCES

- Yağız Aksoy, Tunç Ozan Aydın, and Marc Pollefeys. 2017a. Designing Effective Inter-Pixel Information Flow for Natural Image Matting. In *Proc. CVPR*.
- Yağız Aksoy, Tunç Ozan Aydın, Marc Pollefeys, and Aljoša Smolić. 2016. Interactive High-Quality Green-Screen Keying via Color Unmixing. *ACM Trans. Graph.* 35, 5 (2016), 152:1–152:12.
- Yağız Aksoy, Tunç Ozan Aydın, Aljoša Smolić, and Marc Pollefeys. 2017b. Unmixing-Based Soft Color Segmentation for Image Manipulation. *ACM Trans. Graph.* 36, 2 (2017), 19:1–19:19.
- Yağız Aksoy, Tae-Hyun Oh, Sylvain Paris, Marc Pollefeys, and Wojciech Matusik. 2018. Semantic Soft Segmentation. *ACM Trans. Graph. (Proc. SIGGRAPH)* 37, 4 (2018), 72:1–72:13.
- Omri Avrahami, Dani Lischinski, and Ohad Fried. 2022. Blended diffusion for text-driven editing of natural images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 18208–18218.
- Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, et al. 2022. ediffi: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324* (2022).
- cagliostrolab. 2024. animagine-xl-3.0. *huggingface* (2024).
- Mingdeng Cao, Xintao Wang, Zhongang Qi, Ying Shan, Xiaoju Qie, and Yinqiang Zheng. 2023. MasaCtrl: Tuning-Free Mutual Self-Attention Control for Consistent Image Synthesis and Editing. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*. 22560–22570.
- Guowei Chen, Yi Liu, Jian Wang, Juncai Peng, Yuying Hao, Lutao Chu, Shiyu Tang, Zewu Wu, Zeyu Chen, Zhiliang Yu, Yuning Du, Qingqing Dang, Xiaoguang Hu, and Dianhai Yu. 2022. PP-Matting: High-Accuracy Natural Image Matting.

- Jianqi Chen, Yilan Zhang, Zhengxia Zou, Keyan Chen, and Zhenwei Shi. 2023. Dense Pixel-to-Pixel Harmonization via Continuous Image Representation. *IEEE Transactions on Circuits and Systems for Video Technology* (2023), 1–1. <https://doi.org/10.1109/TCSVT.2023.3324591>
- Wenyan Cong, Jianfu Zhang, Li Niu, Liu Liu, Zhixin Ling, Weiyuan Li, and Liqing Zhang. 2020. Dovenet: Deep image harmonization via domain verification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 8394–8403.
- Guillaume Couairon, Jakob Verbeek, Holger Schwenk, and Matthieu Cord. 2023. DiffEdit: Diffusion-based semantic image editing with mask guidance. In *International Conference on Learning Representations (ICLR)*.
- diffusers. 2024. stable-diffusion-xl-1.0-inpainting-0.1. *diffusers* (2024).
- Zheng-Jun Du, Liang-Fu Kang, Jianchao Tan, Yotam Gingold, and Kun Xu. 2023. Image vectorization and editing via linear gradient layer decomposition. *ACM Transactions on Graphics (TOG)* 42, 4 (Aug. 2023).
- Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. 2022. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618* (2022).
- Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. 2015. Explaining and Harnessing Adversarial Examples. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7–9, 2015, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.). arXiv:1412.6572
- Julian Jorge Andrade Guerreiro, Mitsuru Nakazawa, and Björn Stenger. 2023. PCT-Net: Full Resolution Image Harmonization Using Pixel-Wise Color Transformations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 5917–5926.
- Zonghui Guo, Haiyong Zheng, Yufeng Jiang, Zhaorui Gu, and Bing Zheng. 2021. Intrinsic image harmonization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16367–16376.
- Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. 2023. Prompt-to-Prompt Image Editing with Cross-Attention Control. In *International Conference on Learning Representations (ICLR)*.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. *Denoising diffusion probabilistic models*. NeurIPS.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. LoRA: Low-Rank Adaptation of Large Language Models. *arXiv preprint arXiv:2106.09685* (2021).
- Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. 2021. OpenCLIP.
- Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. 2023. Imagic: Text-Based Real Image Editing with Diffusion Models. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Gwanghyun Kim, Taesung Kwon, and Jong Chul Ye. 2022a. DiffusionCLIP: Text-Guided Diffusion Models for Robust Image Manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2426–2435.
- Gwanghyun Kim, Taesung Kwon, and Jong Chul Ye. 2022b. DiffusionCLIP: Text-Guided Diffusion Models for Robust Image Manipulation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2416–2425.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. 2023. Segment Anything. *arXiv:2304.02643* (2023).
- Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. 2023. Pick-a-Pic: An Open Dataset of User Preferences for Text-to-Image Generation.
- Zhifeng Kong and Wei Ping. 2021. On fast sampling of diffusion probabilistic models. *CoRR* 2106 (2021).
- Yuki Koyama and Masataka Goto. 2018. Decomposing Images into Layers with Advanced Color Blending. *Computer Graphics Forum* 37, 7 (Oct. 2018), 397–407. <https://doi.org/10.1111/cgf.13577>
- Jiachen Li, Jitesh Jain, and Humphrey Shi. 2023b. Matting Anything. *arXiv: 2306.05399* (2023).
- Pengzhi Li, Qinxuan Huang, Yikang Ding, and Zhiheng Li. 2023a. LayerDiffusion: Layered Controlled Image Editing with Diffusion Models. arXiv:2305.18676 [cs.CV]
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual Instruction Tuning. In *NeurIPS*.
- Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. 2023. Null-text inversion for editing real images using guided diffusion models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 6038–6047.
- Chong Mou, Xintao Wang, Liangbin Xie, Jian Zhang, Zhongang Qi, Ying Shan, and Xiaoju Qie. 2023. T2I-Adapter: Learning Adapters to Dig out More Controllable Ability for Text-to-Image Diffusion Models. *arXiv preprint arXiv:2302.08453* (2023).
- Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. 2021. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741* (2021).
- Li Niu, Junyan Cao, Wenyan Cong, and Liqing Zhang. 2023. Deep Image Harmonization with Learnable Augmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 7482–7491.
- Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. 2023. SDXL: Improving Latent Diffusion Models for High-Resolution Image Synthesis. (July 2023). arXiv:2307.01952 [cs.CV]
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. (Oct. 2019). <https://doi.org/10.48550/ARXIV.1910.10683> arXiv:1910.10683 [cs.LG]
- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125* (2022).
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10684–10695.
- Natanuel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. 2022. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. *arXiv preprint arXiv:2208.12242* (2022).
- Robin San-Roman, Eliya Nachmani, and Lior Wolf. 2021. Noise estimation for generative diffusion models. *CoRR* 2104 (2021).
- Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade W Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa R Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. 2022. LAION-5B: An open large-scale dataset for training next generation image-text models. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Christoph Schuhmann and Peter Bevan. 2023. LAION POP: 600,000 High-Resolution Images With Detailed Descriptions. <https://huggingface.co/datasets/laion/laion-pop>.
- Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. 2015. Deep unsupervised learning using nonequilibrium thermodynamics. *CoRR* 1503 (2015).
- Jiaming Song, Chenlin Meng, and Stefano Ermon. 2021. *Denoising diffusion implicit models*. In *ICLR*. OpenReview.net.
- Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. 2020. Score-based generative modeling through stochastic differential equations. *CoRR* 2011 (2020), 13456.
- Stability. 2022a. Stable Diffusion v1.5 Model Card, <https://huggingface.co/runwayml/stable-diffusion-v1-5>.
- Stability. 2022b. Stable Diffusion v2 Model Card, Stable-Diffusion-2-Depth, <https://huggingface.co/stabilityai/stable-diffusion-2-depth>.
- Jianchao Tan, Stephen DiVerdi, Jingwan Lu, and Yotam Gingold. 2019. Pigmento: Pigment-Based Image Analysis and Editing. *Transactions on Visualization and Computer Graphics (TVCG)* 25, 9 (2019). <https://doi.org/10.1109/TVCG.2018.2858238>
- Jianchao Tan, Marek Dvorožňák, Daniel Šýkora, and Yotam Gingold. 2015. Decomposing Time-Lapse Paintings into Layers. *ACM Transactions on Graphics (TOG)* 34, 4, Article 61 (July 2015), 10 pages. <https://doi.org/10.1145/2766960>
- Jianchao Tan, Jose Echevarria, and Yotam Gingold. 2018. Efficient palette-based decomposition and recoloring of images via RGBXY-space geometry. *ACM Transactions on Graphics (TOG)* 37, 6, Article 262 (Dec. 2018), 10 pages. <https://doi.org/10.1145/3272127.3275054>
- Jianchao Tan, Jyh-Ming Lien, and Yotam Gingold. 2016. Decomposing Images into Layers via RGB-space Geometry. *ACM Transactions on Graphics (TOG)* 36, 1, Article 7 (Nov. 2016), 14 pages. <https://doi.org/10.1145/2988229>
- Lin Feng Tan, Jiangtong Li, Li Niu, and Liqing Zhang. 2023. Deep image harmonization in dual color spaces. In *Proceedings of the 31st ACM International Conference on Multimedia*. 2159–2167.
- Jingwei Tang, Yağız Aksoy, Cengiz Öztireli, Markus Gross, and Tunç Ozan Aydın. 2019. Learning-based Sampling for Natural Image Matting. In *Proc. CVPR*.
- Yi-Hsuan Tsai, Xiaohui Shen, Zhe Lin, Kalyan Sunkavalli, Xin Lu, and Ming-Hsuan Yang. 2017. Deep image harmonization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 3789–3797.
- Narek Tumanyan, Michal Geyer, Shai Bagon, and Tali Dekel. 2023. Plug-and-Play Diffusion Features for Text-Driven Image-to-Image Translation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 1921–1930.
- Menghan Xia, Xueting Liu, and Tien-Tsin Wong. 2018. Invertible Grayscale. *ACM Transactions on Graphics (SIGGRAPH Asia 2018 issue)* 37, 6 (Nov. 2018), 246:1–246:10.
- Mingqing Xiao, Shuxin Zheng, Chang Liu, Yaolong Wang, Di He, Guolin Ke, Jiang Bian, Zhouchen Lin, and Tie-Yan Liu. 2020. *Invertible Image Rescaling*. Springer International Publishing, 126–144.
- Ning Xu, Brian Price, Scott Cohen, and Thomas Huang. 2017. Deep Image Matting. (March 2017). <https://doi.org/10.48550/ARXIV.1703.03872> arXiv:1703.03872 [cs.CV]
- Xingqian Xu, Zhangyang Wang, Eric Zhang, Kai Wang, and Humphrey Shi. 2022. Versatile diffusion: Text, images and variations all in one diffusion model. *arXiv preprint arXiv:2211.08332* (2022).

- Jingfeng Yao, Xinggang Wang, Shusheng Yang, and Baoyuan Wang. 2024. ViTMatte: Boosting image matting with pre-trained plain vision transformers. *Information Fusion* 103 (2024), 102091.
- Hu Ye, Jun Zhang, Sibio Liu, Xiao Han, and Wei Yang. 2023. IP-Adapter: Text Compatible Image Prompt Adapter for Text-to-Image Diffusion Models. (2023).
- Lvmin Zhang and Maneesh Agrawala. 2023. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543* (2023).
- Xinyang Zhang, Wentian Zhao, Xin Lu, and Jeff Chien. 2023. Text2Layer: Layered Image Generation using Latent Diffusion Model. arXiv:2307.09781 [cs.CV]
- Jun-Yan Zhu, Philipp Krahenbuhl, Eli Shechtman, and Alexei A Efros. 2015. Learning a discriminative model for the perception of realism in composite images. In *Proceedings of the IEEE International Conference on Computer Vision*. 3943–3951.
- Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. 2017. Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks. In *Computer Vision (ICCV), 2017 IEEE International Conference on*.

A PADDED RGB CHANNELS

In the RGB channel of a transparent RGBA image, we refer to pixels that are completely invisible as “undefined” pixels, *i.e.*, pixels with alpha value strictly equal to zero. Since these pixels are strictly invisible, processing them with arbitrary color does not influence the appearance of images after alpha blending. Nevertheless, since neural networks tends to produce high-frequency patterns surrounding image edges, we avoid unnecessary edges in the RGB channels to avoid potential artifacts. We define a local Gaussian filter

$$G(I_c)_p = \begin{cases} \phi(I_c)_p & , \text{if } (I_a)_p = 0 \\ (I_c)_p & , \text{otherwise} \end{cases} \quad (10)$$

where $\phi(\cdot)$ is a standard Gaussian filter with $13 * 13$ kernel, and p is pixel position. We perform this filter 64 times to completely propagate colors to all “undefined” pixels.

B NEURAL NETWORK ARCHITECTURE

The latent transparency encoder has exactly same neural network architecture with Stable Diffusion latent VAE encoder [Podell et al. 2023] (but the input contains 4 channels for RGBA). This model is trained from scratch. The output convolution layer is zero-initialized to avoid initial harmful noise.

The latent transparency decoder is a UNet. The encoding part of this UNet has same architecture as Stable Diffusion’s latent VAE encoder, while the decoding part has same architecture as Stable Diffusion’s VAE decoder. The input latent is added to the middle block, and all the encoder’s feature maps are added to the input of each decoder block with skip connection. To be specific, assuming the input image is $512 \times 512 \times 3$ and the input latent is $64 \times 64 \times 4$, the feature map goes through $512 \times 512 \times 3 \rightarrow 512 \times 512 \times 128 \rightarrow 256 \times 256 \times 256 \rightarrow 128 \times 128 \times 512 \rightarrow 64 \times 64 \times 512$ where each \rightarrow is two resnet blocks. Then input latent is projected by a convolution layer to match channel and then added to the middle feature. Then the decoder goes through $64 \times 64 \times 512 \rightarrow 128 \times 128 \times 512 \rightarrow 256 \times 256 \times 256 \rightarrow 512 \times 512 \times 128 \rightarrow 512 \times 512 \times 3$ and here each \rightarrow also adds the skip features from the encoder’s corresponding layers.

C PATCHGAN DISCRIMINATOR

We use exactly same PatchGAN Discriminator architecture, learning objective, and training scheduling with Latent Diffusion VAE [Rombach et al. 2022]. We directly use the python class *LPIPSWithDiscriminator* from their official code base (the input channel is set to 4). The generator side objective (from [Rombach et al. 2022]) can

be written as

$$\mathbb{L}_{\text{disc}}(\mathbf{z}) = \text{relu}(1 - D_{\text{disc}}(\mathbf{z})), \quad (11)$$

where \mathbf{z} is a matrix with shape $h \times w \times 4$ and $\text{relu}(\cdot)$ is rectified linear unit. The $D_{\text{disc}}(\cdot)$ is a neural network with 5 convolution-normalization-silu layers $512 \times 512 \times 3 \rightarrow 512 \times 512 \times 64 \rightarrow 256 \times 256 \times 128 \rightarrow 128 \times 128 \times 256 \rightarrow 64 \times 64 \times 512 \rightarrow 64 \times 64 \times 1$ and the last layer is a patch-wise real/fake classification layer. The last layer does not use normalization and activation.

D SINGLE TRANSPARENT IMAGES

We present additional results for single transparent images, from Figure 18 to Figure 33.

E MULTIPLE TRANSPARENT LAYERS

We present additional results for multiple transparent layers, from Figure 34 to Figure 36.

F FOREGROUND-CONDITIONED BACKGROUNDS

We present additional results for foreground-conditioned backgrounds, from Figure 37 to Figure 38.

G BACKGROUND-CONDITIONED FOREGROUNDS

We present additional results for background-conditioned foregrounds in Figure 39.

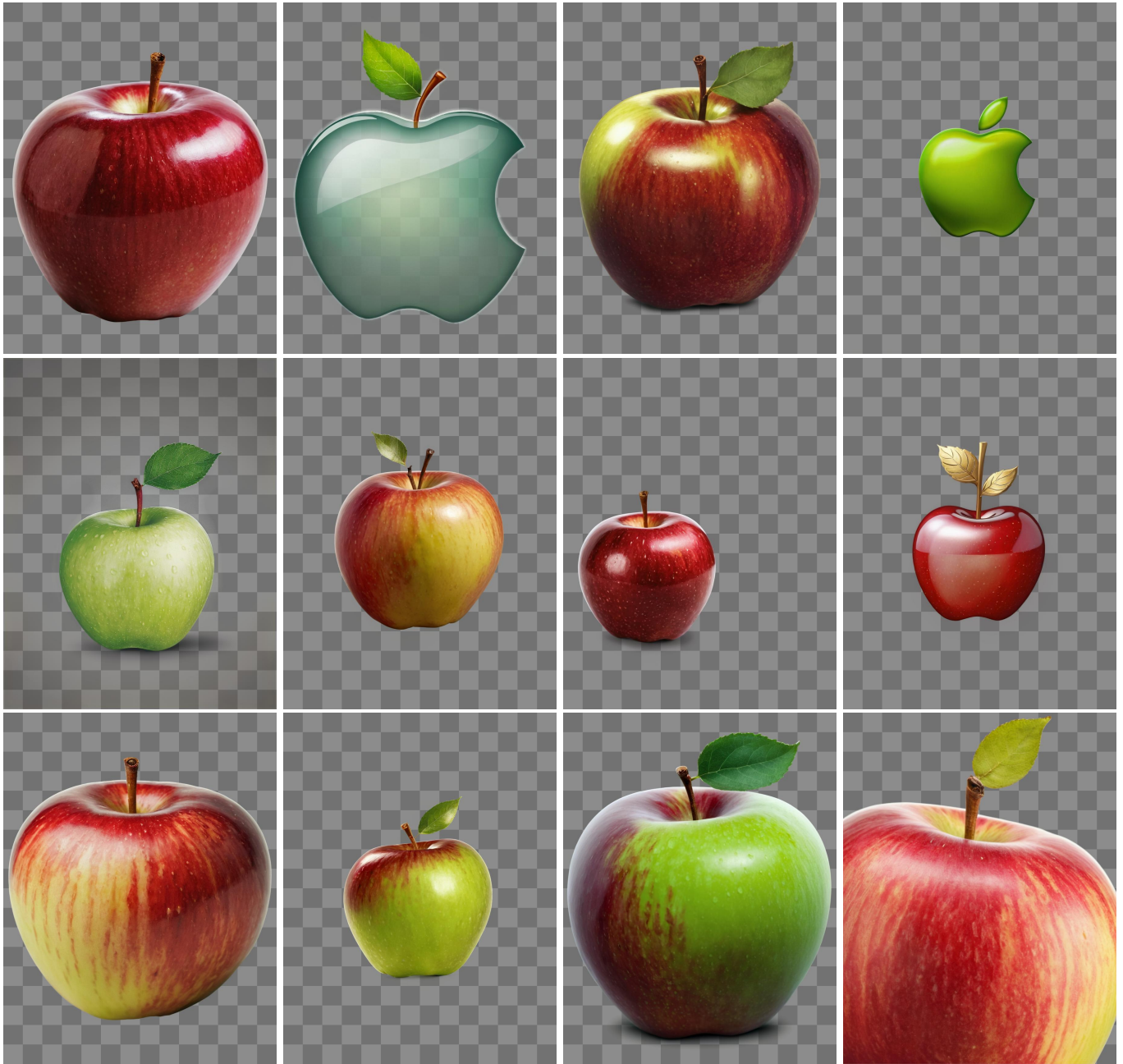


Fig. 18. Single Transparent Image Results #1. The prompt is “apple”. Resolution is 896 × 1152.

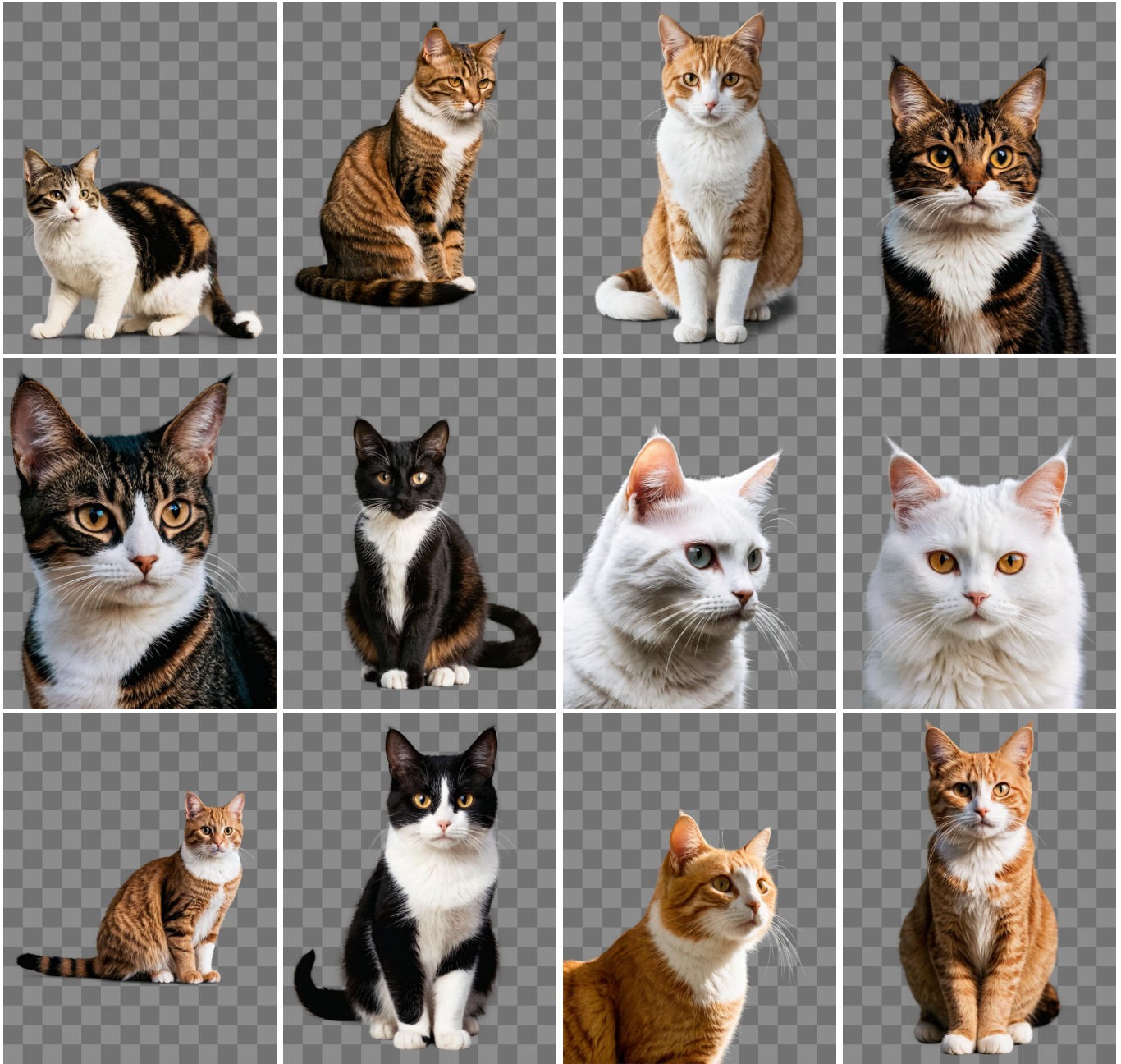


Fig. 19. Single Transparent Image Results #2. The prompt is “a cat”. Resolution is 896×1152 .



Fig. 20. Single Transparent Image Results #3. The prompt is “a man”. Resolution is 896×1152 .



Fig. 21. Single Transparent Image Results #4. The prompt is “a man with messy hair”. Resolution is 896×1152 .



Fig. 22. Single Transparent Image Results #5. The prompt is "woman". Resolution is 896×1152 .



Fig. 23. Single Transparent Image Results #6. The prompt is “woman with messy hair”. Resolution is 896×1152 .



Fig. 24. Single Transparent Image Results #7. The prompt is “dog”. Resolution is 1024×1024 .



Fig. 25. Single Transparent Image Results #8. The prompt is "glass cup". Resolution is 1024×1024 .



Fig. 26. Single Transparent Image Results #9. The prompt is “dragon”. Resolution is 1152 × 896.



Fig. 27. Single Transparent Image Results #10. The prompt is “car”. Resolution is 1152 × 896.



Fig. 28. Single Transparent Image Results #11. The prompt is “magic book”. Resolution is 1152 × 896.

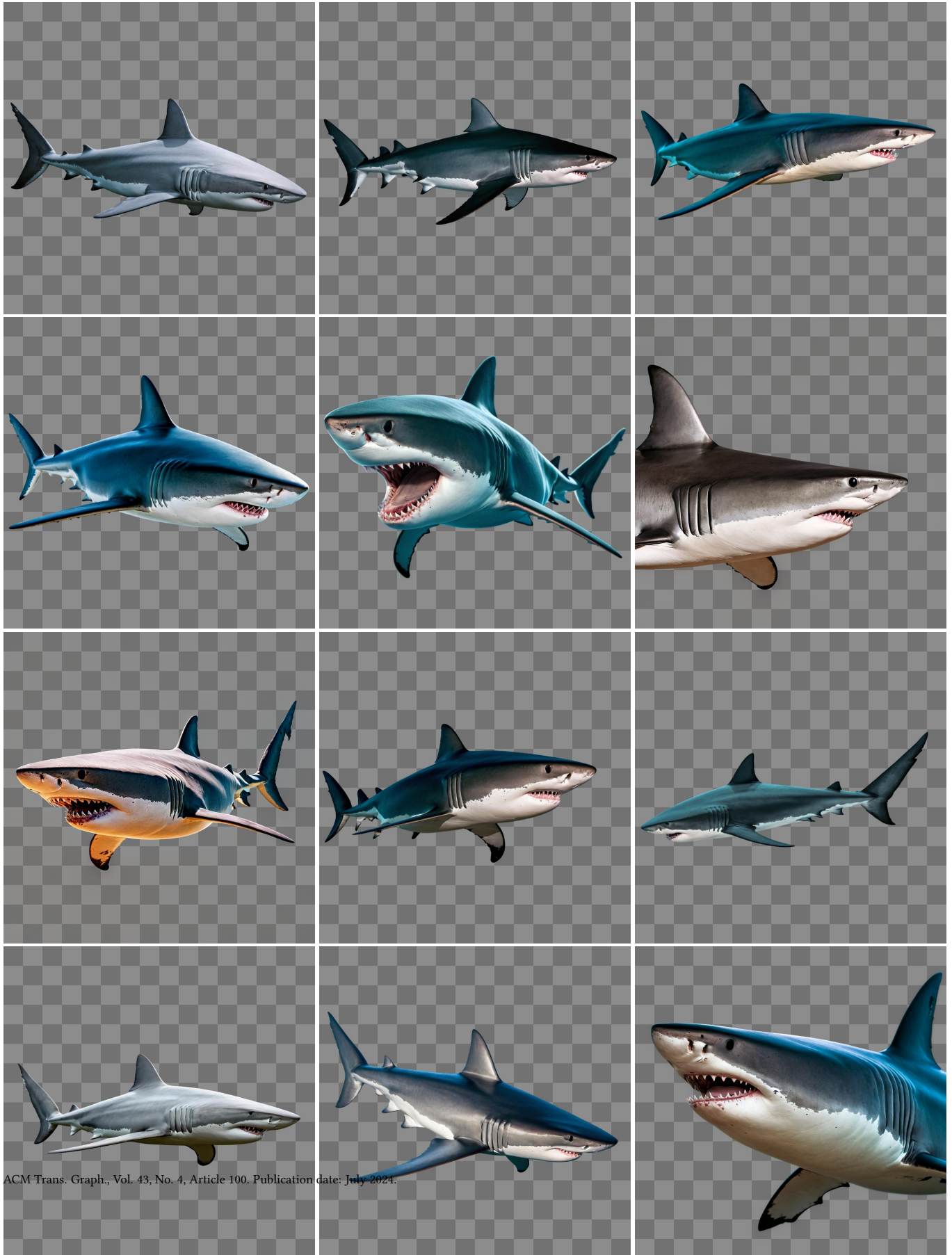


Fig. 29. Single Transparent Image Results #12. The prompt is “shark”. Resolution is 1024×1024 .

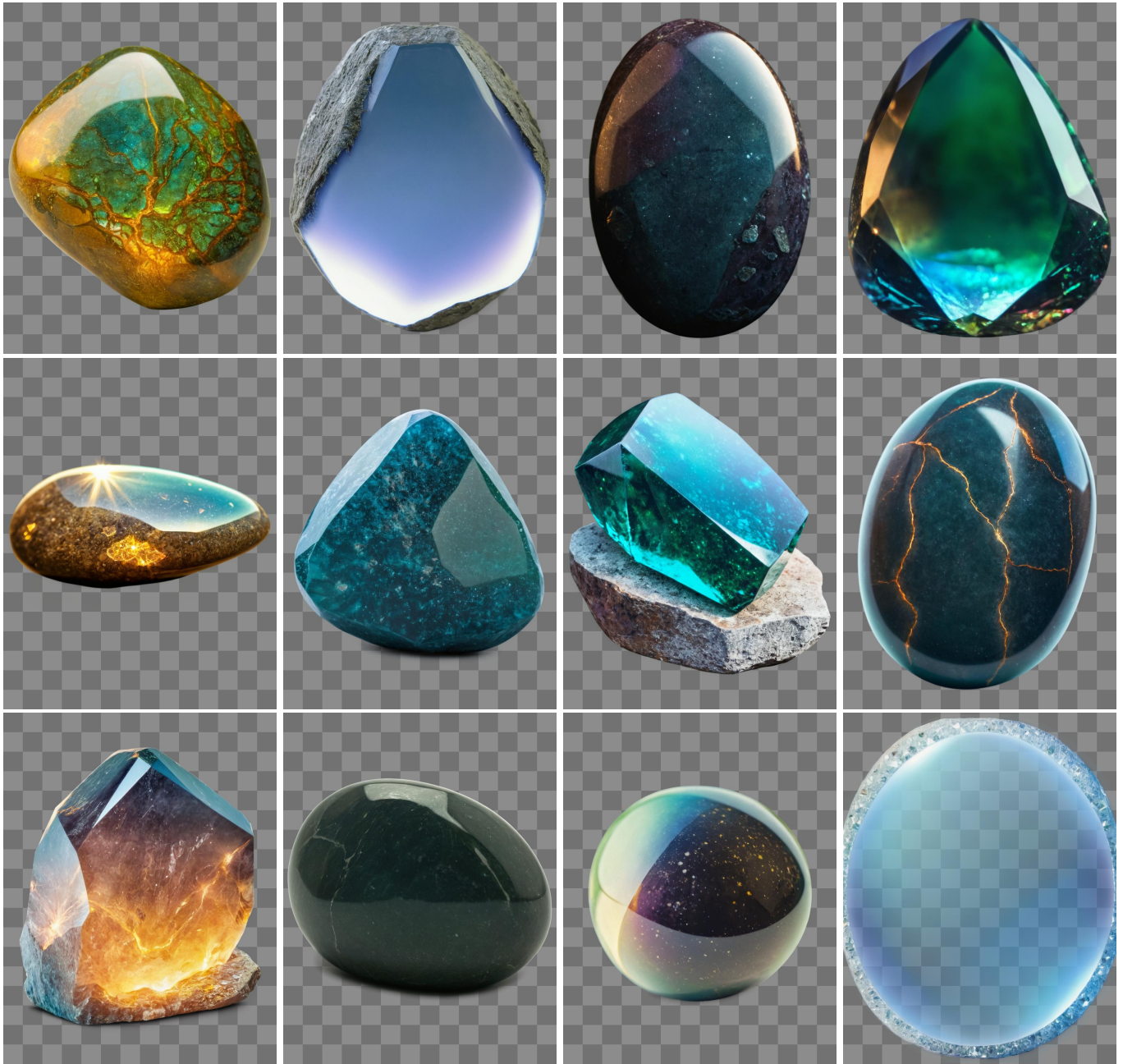


Fig. 30. Single Transparent Image Results #13. The prompt is "magic stone". Resolution is 896×1152 .



Fig. 31. Single Transparent Image Results #14. The prompt is "parrot, green fur". Resolution is 896×1152 .

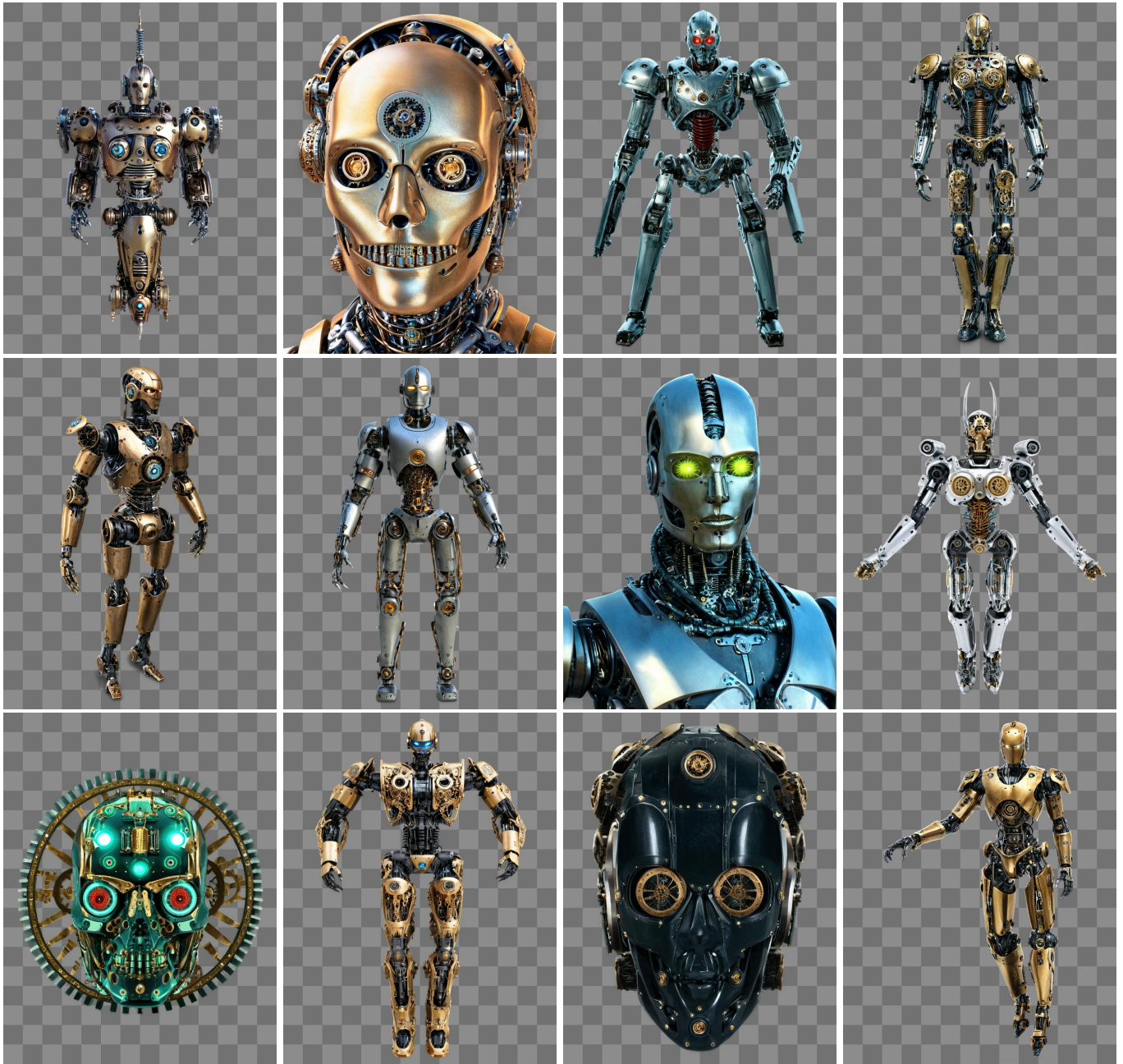


Fig. 32. Single Transparent Image Results #15. The prompt is “cyber steampunk robot”. Resolution is 896×1152 .



Fig. 33. Single Transparent Image Results #16. The prompt is “necromancer”. Resolution is 896×1152 .



Fig. 34. Multi-layer Results #1. The prompts are “plant on table”, “woman in room”, “dog on floor”, “man walking on street”. Resolution is 896×1152 .

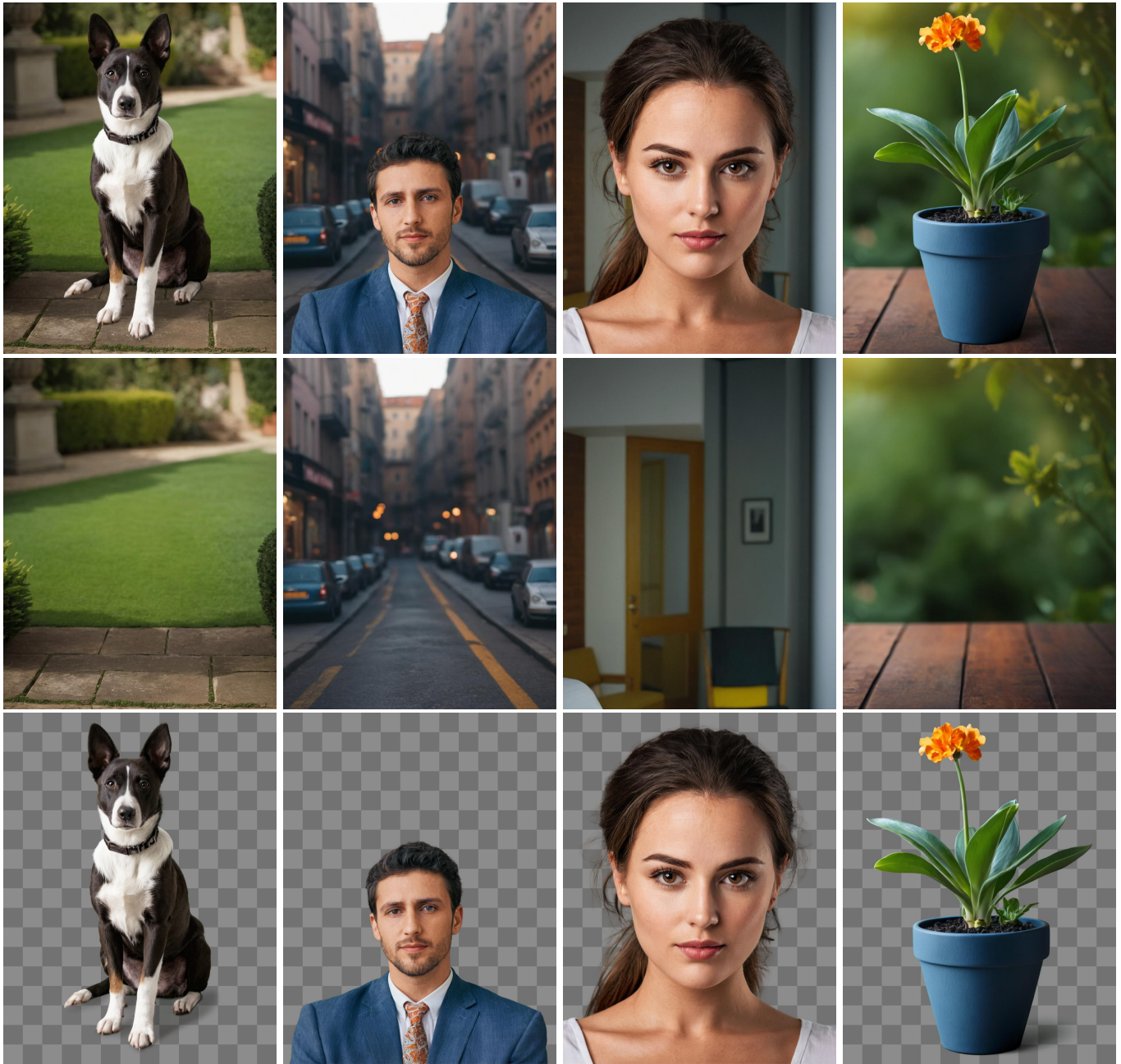


Fig. 35. Multi-layer Results #2. The prompts are “dog in garden”, “man in street”, “woman, closeup”, “plants on table”. Resolution is 896×1152 .



Fig. 36. Multi-layer Results #3. The prompts are “cat on floor”, “woman in room”, “man in room”, “golden cup”. Resolution is 896×1152 .

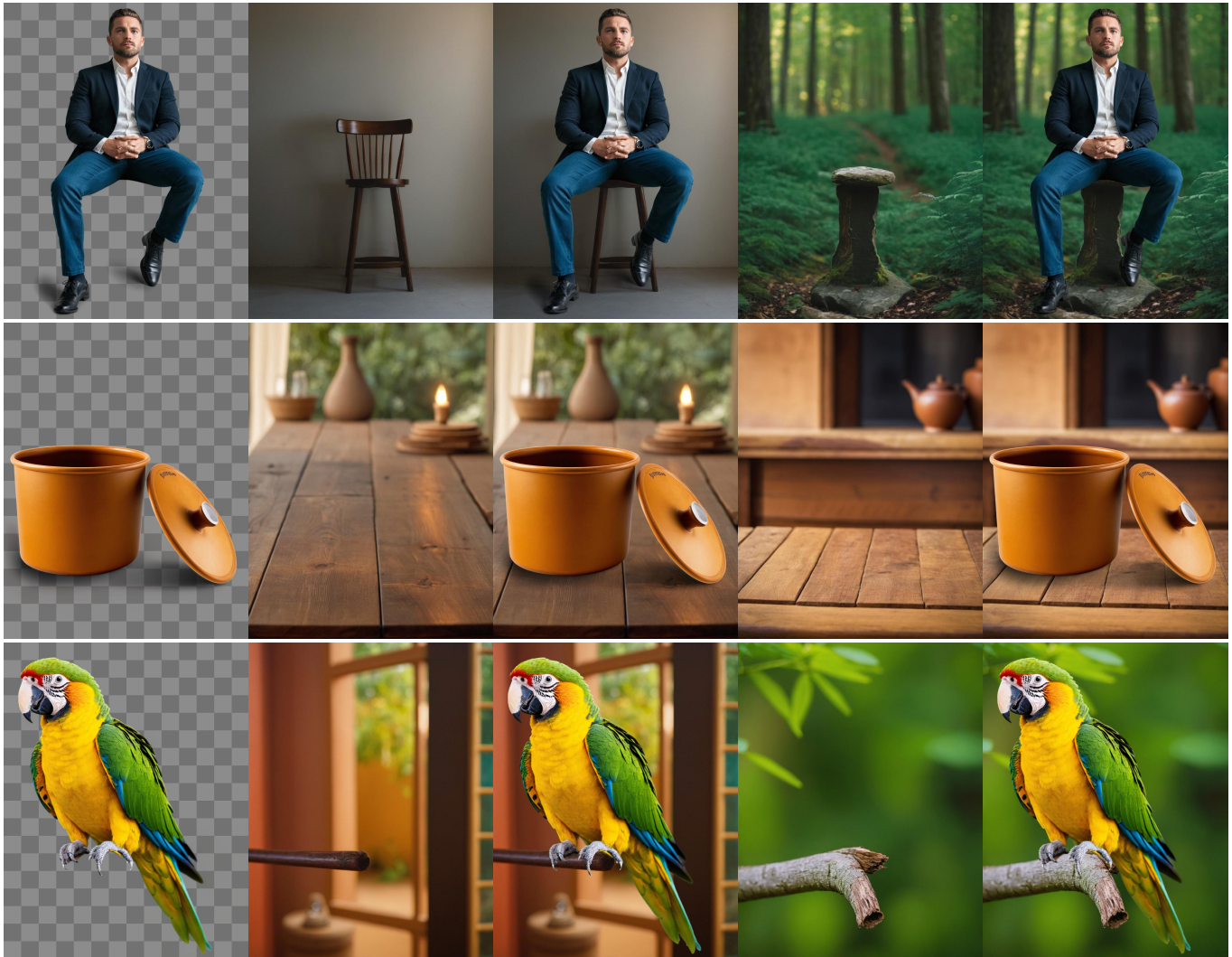


Fig. 37. Foreground-conditioned Background Results #1. The left-most images are inputs. The prompts are “man sitting on chair”, “man sitting in forest”, “pots on wood table”, “parrot in room”, “parrot in forest”. Resolution is 896×1152 .

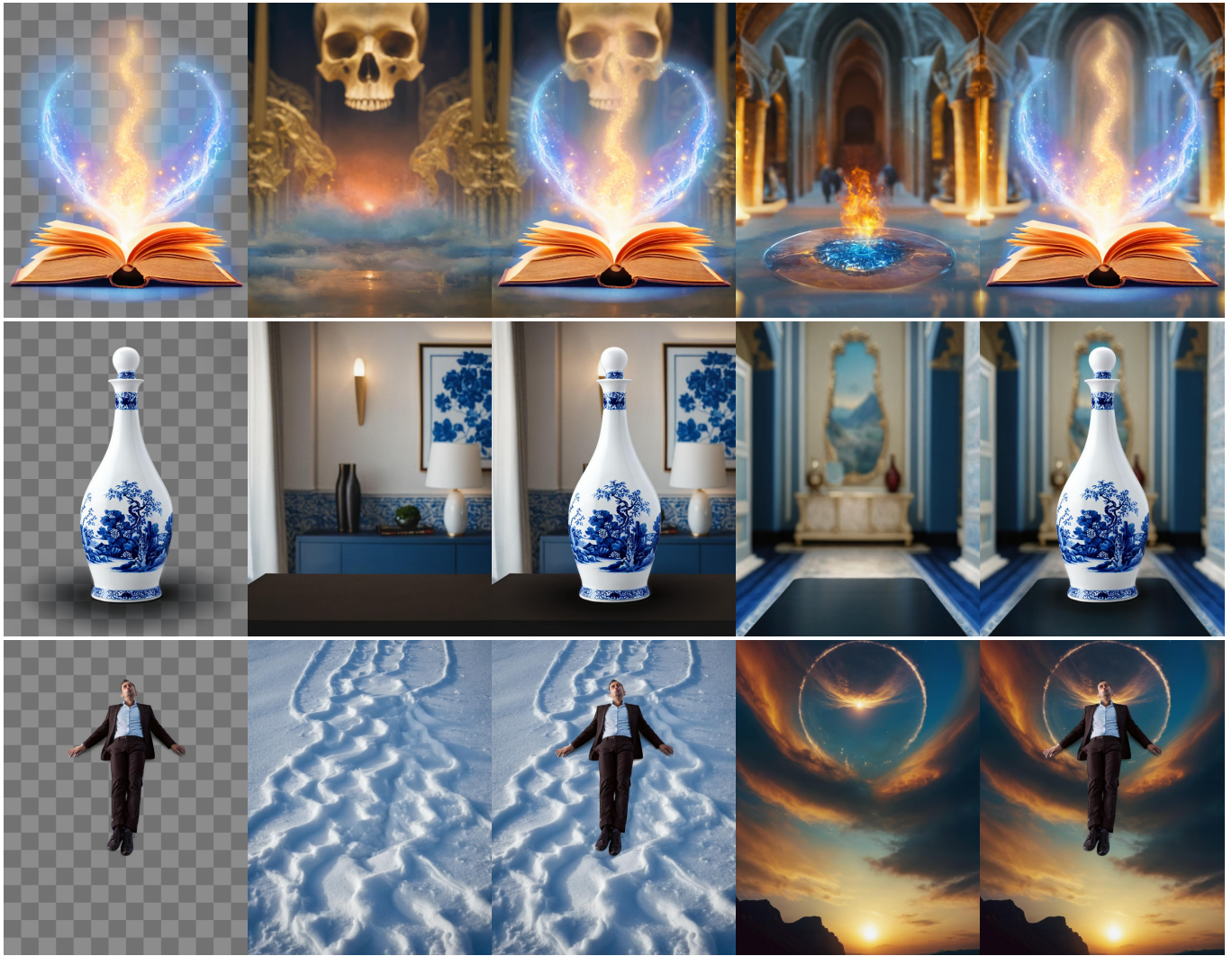


Fig. 38. Foreground-conditioned Background Results #2. The left-most images are inputs. The prompts are “magic book of death”, “magic book of life”, “blue and white porcelain vase in my home”, “blue and white porcelain vase in the museum”, “the man in the snow”, “god of infinity”. Resolution is 896×1152 .

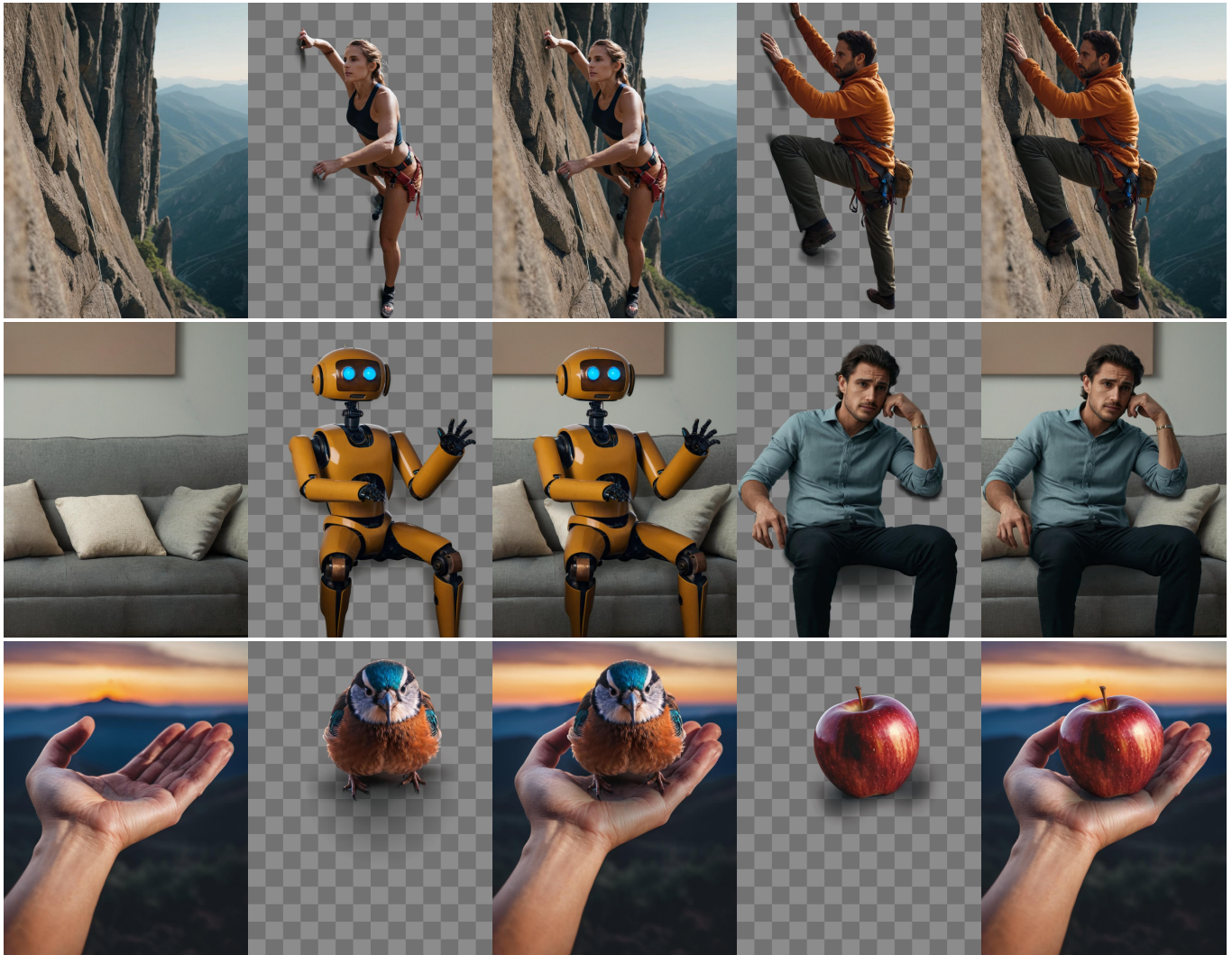


Fig. 39. Background-conditioned Foreground Results #1. The left-most images are inputs. The prompts are “woman climbing mountain”, “man climbing mountain”, “robot in sofa waving hand”, “man in sofa”, “bird on hand”, “apple on hand”. Resolution is 896×1152 .