

MirrorVerse: Pushing Diffusion Models to Realistically Reflect the World

Anonymous CVPR submission

Paper ID 35

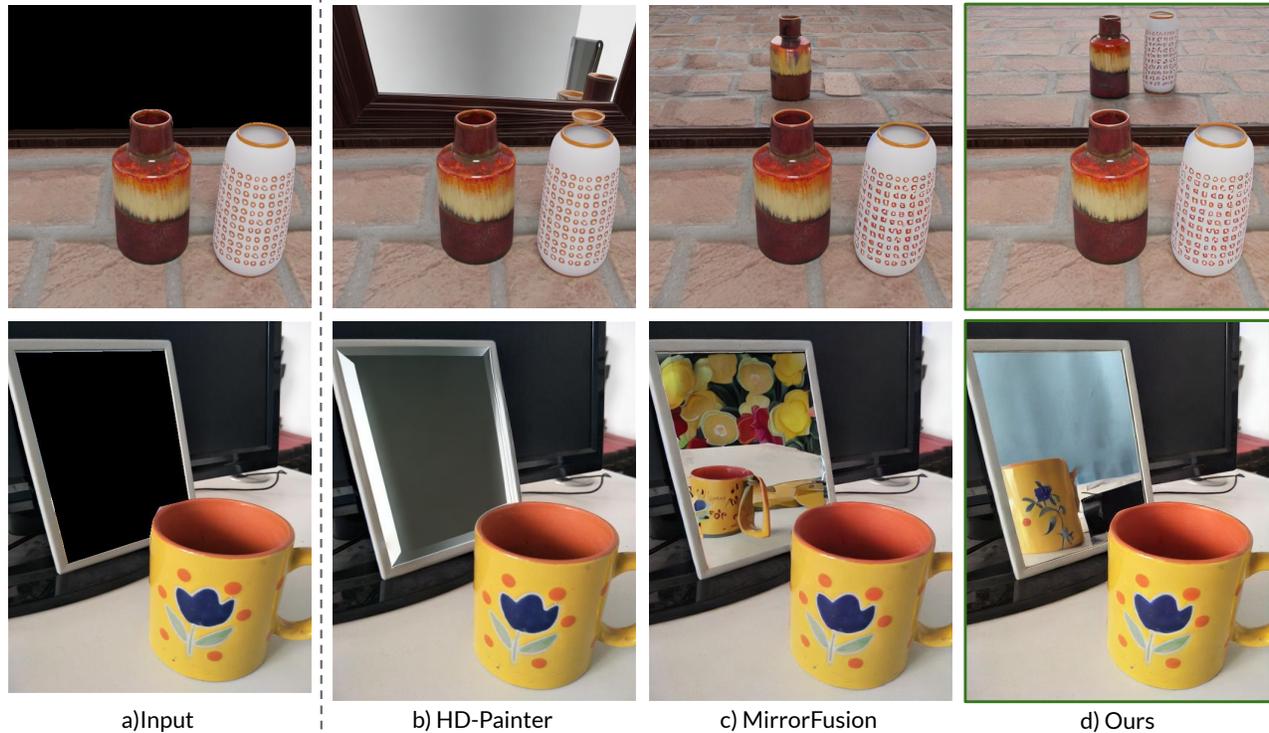


Figure 1. Our model MirrorFusion 2.0, trained on our enhanced dataset SynMirrorV2 surpasses previous state-of-the-art diffusion-based inpainting models at the task of generating mirror reflections. All images were created by appending the prompt: “A perfect plane mirror reflection of ” to the object description. All text prompts can be found in the supplementary.

Abstract

001 Diffusion models have become central to various image
 002 editing tasks, yet they often fail to fully adhere to physical
 003 laws, particularly with effects like shadows, reflections,
 004 and occlusions. In this work, we address the challenge of
 005 generating photorealistic mirror reflections using diffusion-
 006 based generative models. Despite extensive training data,
 007 existing diffusion models frequently overlook the nuanced
 008 details crucial to authentic mirror reflections. Recent ap-
 009 proaches have attempted to resolve this by creating syn-
 010 thetic datasets and framing reflection generation as an in-
 011 painting task; however, they struggle to generalize across
 012 different object orientations and positions relative to the
 013 mirror. Our method overcomes these limitations by intro-
 014 ducing key augmentations into the synthetic data pipeline:

(1) random object positioning, (2) randomized rotations,
 and (3) grounding of objects, significantly enhancing gener-
 alization across poses and placements. To further address
 spatial relationships and occlusions in scenes with multi-
 ple objects, we implement a strategy to pair objects during
 dataset generation, resulting in a dataset robust enough to
 handle these complex scenarios. Achieving generalization
 to real-world scenes remains a challenge, so we introduce
 a three-stage training curriculum to develop the MirrorFu-
 sion 2.0 model to improve real-world performance. We pro-
 vide extensive qualitative and quantitative evaluations to
 support our approach.

015
 016
 017
 018
 019
 020
 021
 022
 023
 024
 025
 026



Figure 2. We observe that current state-of-the-art T2I models, SD3.5 [2] (top row) and Flux [22] (bottom row), face significant challenges in producing consistent and geometrically accurate reflections when prompted to generate reflections in the scene.

027 1. Introduction

028 In recent years, diffusion-based generative models have re-
 029 defined what is possible in fields spanning from image
 030 generation to video synthesis, producing impressive re-
 031 sults across various applications [14, 17, 18, 22, 35, 38].
 032 The evolution of these models has been accompanied by
 033 a range of methods designed to fine-tune the generation
 034 process through conditional inputs, such as edge maps,
 035 sketches, depth maps, and segmentation maps [30, 54, 56,
 036 58]. However, there remains a significant gap in their ca-
 037 pacity to replicate intricate physical effects—particularly
 038 those rooted in the subtlety of real-world physics, includ-
 039 ing shadows [41], specular reflections [49], and perspective
 040 cues [46]. More challenging still, these techniques struggle
 041 to authentically generate mirror reflections, a task requir-
 042 ing a nuanced understanding of light, geometry, and real-
 043 ism that current methods do not adequately address. In this
 044 work, we address the question: “*Can current methods be*
 045 *fine-tuned to generate plausible mirror reflections?*”

046 We motivate the problem further by providing genera-
 047 tions from current text-to-image (T2I) generation models.
 048 We prompt Stable Diffusion 3.5 [2] and FLUX [22] with
 049 prompts to generate a scene with a mirror reflection. Fig. 2
 050 shows that these methods fail to generate plausible mir-
 051 ror reflections. Specifically, check the reflection of “teddy-
 052 bear” in the generated outputs from both the methods. Fur-
 053 ther, inpainting methods like HD-Painter [27] also fail for
 054 this task, as shown in Fig. 1. A contemporary method called
 055 MirrorFusion [12], claiming to generate mirror reflections,
 056 falls short on real-world and challenging scenes as appar-
 057 ent in Fig. 1.

058 Despite their impressive capabilities, powerful diffusion
 059 models struggle to generate mirror reflections accurately.

Table 1. Our proposed dataset, **SynMirrorV2**, surpasses existing mirror datasets in terms of attribute diversity and variability. While recent work [12] introduced the synthetic SynMirror dataset, it lacks key augmentations and scenario, limiting its performance in complex and real-world settings (See Fig. 1).

Dataset	Type	Size (#Images)	Attributes
MSD [52]	Real	4,018	RGB, Masks
Mirror-NeRF [55]	Real & Synthetic	9 scenes	RGB, Masks, Multi-View
DLSU-OMRS [15]	Real	454	RGB, Mask
TROSD [44]	Real	11,060	RGB, Mask
PMD [24]	Real	6,461	RGB, Masks
RGBD-Mirror [28]	Real	3,049	RGB, Depth
Mirror3D [45]	Real	7,011	RGB, Masks, Depth
SynMirror [12]	Synthetic	198,204	Single Fixed Objects: RGB, Depth, Masks, Normals, Multi-View
SynMirrorV2 (Ours)	Synthetic with Single & Multiple Objects	207,610	Single + Multiple Objects : RGB, Depth, Masks, Normals, Multi-View, Augmentations

This limitation stems from the models’ reliance on poorly
 learned priors, a consequence of the quality and quantity
 of their training data. The scarcity of high-quality, real-
 world images featuring mirrors and their reflections, as ev-
 idenced in Tab. 1, poses a significant challenge. While re-
 cent work [12] has attempted to address this issue by train-
 ing on a synthetic dataset, the results, as illustrated in Fig. 1,
 suggest that the method’s performance suffers in complex
 scenes and real-world settings. We hypothesize that this is
 due to inherent limitations in the synthetic data generation
 process and the training dynamics of the model.

To address the shortcomings in the synthetic data gener-
 ation pipeline, we create an enhanced pipeline incorporat-
 ing useful augmentations such as randomizing object po-
 sition and rotation. We also ensure that the objects are
 anchored to the ground level in the 3D world. We ob-
 serve that this diverse data improves the generalization of
 a trained model across the pose and position of objects in
 the scene. However, it does not generalize to more complex
 scenes with multiple objects. To address this, we propose
 a novel pipeline that places multiple objects in the scene
 based on their semantic categories, further enhancing the
 quality and utility of the proposed synthetic dataset. Draw-
 ing inspiration from previous works, such as those that have
 improved the generation quality on various tasks, notably
 image-editing [29], multilingual T2I generation [23, 50, 53]
 and several others, we aim to leverage the stage-wise train-
 ing approach that enhanced the results in these methods.

We briefly sum up our contributions as follows:

- We propose SynMirrorV2, a large-scale synthetic dataset with diversity in objects and their relative position and orientation in the scene.
- Further, we create a pipeline to add multiple objects to a scene in SynMirrorV2.
- We show that with a curriculum strategy of training on SynMirrorV2, a generative method can also generalize to real-world scenes. We show this generalization capability on the challenging real-world MSD [52] dataset.

098

2. Related Work

099

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

Image Generative models. Diffusion models [43] have become quite popular for image generation tasks. Diffusion models work by gradually adding noise to data and then learning to reverse this process to generate data from a variety of distributions [8, 11, 17]. Subsequent works have expanded the scope of image generation by incorporating text guidance [37, 40] into the diffusion process, simplifying the reverse process [48], and reformulating diffusion to occur in a latent space [38] for improved speed. [31] explore advancements in diffusion models by addressing bias through distribution-guided debiasing techniques. Further, methods [32, 33] are developed to provide more fine-grained generation control to these models. Building on the success of vision transformers [47], recent approaches [7, 34, 60] have replaced the U-Net architecture in diffusion models with transformer-based designs, leading to high-quality image generation results. Further, there are popular methods [2, 22] for high-quality image generation. However, these methods also fail for the task of generating reflections on the mirror as shown in Fig. 2

Image Inpainting. Building on the advancements in image diffusion models, methods like Palette [39] and Repaint [26] leverage known regions through the denoising process to reconstruct missing parts. Blended Diffusion [3, 4] refines this approach by replacing noise in unmasked areas with known content but struggles with complex scenes and shapes. Stable-Diffusion Inpainting [38] (SDI) enhances results by fine-tuning the denoiser with noisy latents, masks, and masked images. Recent methods, such as HD-Painter [27], PowerPoint [61], SmartBrush [51] build on SDI with additional training. Recently, BrushNet [20] introduces a plug-and-play architecture that preserves unmasked content while improving coherence with textual prompts. However, Fig. 1 highlights the limitations of these methods in generating reflections on the mirror.

Diffusion Models and 3D concepts. Recently, LRM [19] based methods predict 3D model from a single image. Some methods [21] utilize diffusion-based methods to enable editing of these 3D presentations. Other diffusion-based methods [29] use synthetic image pairs for 3D-aware image editing. However, the synthetic-to-real domain gap can limit their applicability. Further, ObjectDrop [49] trains a diffusion model for object insertion/removal using a counterfactual dataset that can handle shadows and specular reflections. Sarkar et al. [41] shows that generated images have different geometric features such as shadows and reflections from the real images. Upadhyay et al. [46] proposed a geometric constraint in the training process to improve the perspective cues in the generated images. Alchemist [42] provides control over the material properties of an object by proposing an object-centric synthetic dataset with physically-based materials.

3. Dataset

3.1. Data Generation Pipeline

Fig. 2 highlights the failure of state-of-the-art models in handling the reflection generation task. MirrorFusion [12] addresses this challenge by proposing a synthetic dataset but struggles in complex scenarios involving multiple objects and real-world scenes (Fig. 1). We attribute this limitation to the lack of diversity in their dataset. To mitigate these shortcomings, we introduce SynMirrorV2, a large-scale dataset which significantly expands diversity with varied backgrounds, floor textures, objects, camera poses, mirror orientation, object positions, and rotations. Tab. 1 compares existing mirror datasets, while Fig. 3 showcases samples from SynMirrorV2.

Object Sources. We source objects from Objaverse [9] and Amazon Berkeley Objects (ABO) [6] datasets. Objaverse, a large-scale dataset, contains 800K diverse 3D assets, while ABO contributes 7,953 common household objects. To ensure quality, we refine our selection using a curated list of 64K objects from OBJECT 3DIT [29] and the filtering procedure discussed in [12], eliminating low-quality textures and sub-par renderings. After filtering, we get 58,109 objects from Objaverse. In total, we utilize 66,062 objects.

Scene Resources. To create a realistic scene, we require assets such as a mirror, floor and background. We create a plane for the floor and apply diverse textures sourced from CC-textures [10]. We use HDRI samples provided by PolyHaven [16] to represent the background. In our experiments, we use different kinds of mirrors: full-wall mirrors and tall rectangular mirrors. For lighting, we position an area-light slightly above and behind the object at a 45° angle, directing it towards both the object and the mirror.

Object Placement in the scene. To begin, we fix the mirror’s position within the scene as a fixed reference point. The sampled object is then scaled to fit within a unit cube, ensuring uniformity in size across all objects. We proceed by sampling the object’s x-y position from a pre-computed region that guarantees both visibility of the object in the mirror and camera. This pre-computed region is determined by identifying the intersection between the mirror’s viewing frustum and the camera’s viewing frustum. Once the position is set, we randomly sample an angle for the object’s rotation around the y-axis to introduce variability. However, even with these steps, there may be instances where the object appears to float in the air, which can undermine the dataset’s utility. To address this, we apply a straightforward grounding technique, detailed in the supplementary material. Together, these strategies contribute to the diversity and overall quality of the proposed dataset.

Multiple Objects. A typical scene includes multiple objects arranged in varied layouts, producing a range of depth and occlusion scenarios that enhance scene realism. To

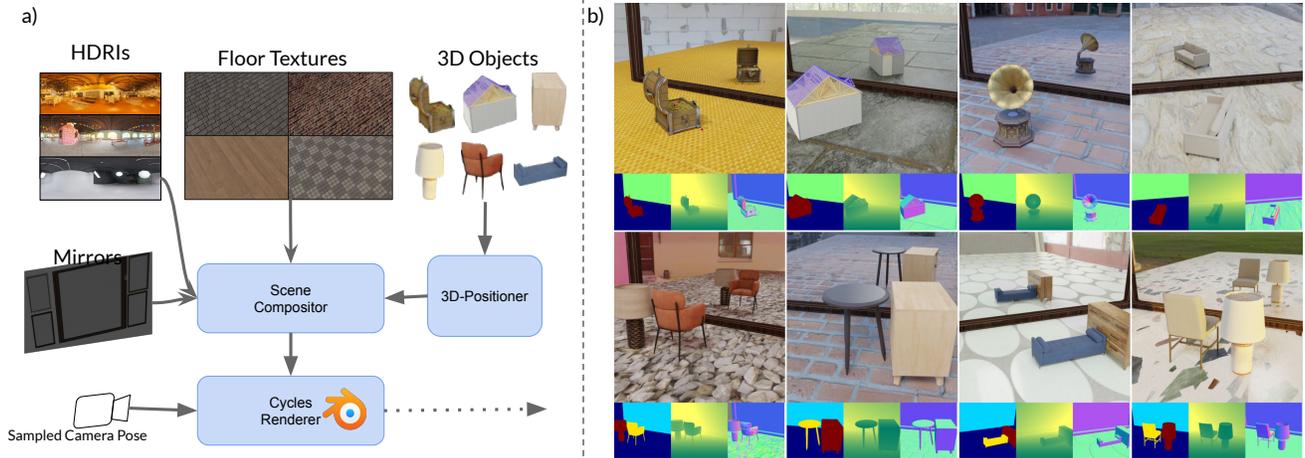


Figure 3. **Dataset Generation Pipeline.** Our dataset generation pipeline introduces key augmentations such as **random positioning, rotation, and grounding** of objects within the scene using the 3D-Positioner. Additionally, we pair objects in semantically consistent combinations to simulate complex spatial relationships and occlusions, capturing realistic interactions for **multi-object** scenes.

Algorithm 1 Procedure to Render Multiple Objects

Require: Input 3D model \mathcal{M}_1

- 1: **Function** GETPAIREDOBJECTCATEGORY(M)
- 2: $c \leftarrow$ GetSemanticCategory(M)
- 3: $L \leftarrow$ GetPairedCategoriesList(c)
- 4: $c_{paired} \leftarrow$ SampleCategory(L)
- 5: **return** c_{paired}
- 6: **Function** SAMPLEOBJECT(c)
- 7: $L_{obj} \leftarrow$ GetListObjects(c)
- 8: $\mathcal{M} \leftarrow$ Sample3DObject(L_{obj})
- 9: **return** \mathcal{M}

Main Algorithm

- 11: $c_{paired} \leftarrow$ GETPAIREDOBJECTCATEGORY(\mathcal{M}_1)
- 12: $\mathcal{M}_2 \leftarrow$ SAMPLEOBJECT(c_{paired})
- 13: Initialize position of \mathcal{M}_2 at X
- 14: **while** \mathcal{M}_2 collides with \mathcal{M}_1 **do**
- 15: $T_r \leftarrow$ SampleRandomPosition()
- 16: $X \leftarrow T_r$
- 17: **end while**

203 capture this complexity, our dataset incorporates scenes
 204 with multiple objects, as described in Algorithm 1. We
 205 start by sampling K objects from the original ABO dataset
 206 and identifying each object’s class from [6]. Categories
 207 are manually paired to ensure semantic coherence—for in-
 208 stance, pairing a chair with a table. During rendering, af-
 209 ter positioning and rotating the primary object K_1 , an addi-
 210 tional object K_2 from the paired category is sampled and ar-
 211 ranged to prevent overlap, ensuring distinct spatial regions
 212 within the scene. This process yields 3,140 scenes featuring
 213 diverse object configurations and spatial relationships, pro-
 214 viding a robust foundation for realistic scene representation.

Rendering. Following scene composition, we randomly
 215 sample three camera poses from a predefined list of 19 camera
 216 positions and render each scene using BlenderProc [10]
 217 to obtain RGB, depth, normal, and semantic label outputs.
 218 All renderings are produced at a resolution of 512×512 pix-
 219 els. We set the “cycles rendering” parameter to 1024, which
 220 is necessary for accurately capturing reflections. Representa-
 221 tive samples are provided in Fig. 3 and additional exam-
 222 ples are available in the supplementary material. 223

4. Method 224

Preliminaries Diffusion models are generative models
 225 that can construct data samples by progressively removing
 226 noise. In the forward diffusion process, Gaussian noise
 227 $\epsilon \sim \mathcal{N}(0, 1)$ is incrementally added to an initial clean sam-
 228 ple x_0 over T timesteps to create a noisy sample x_T . In
 229 the reverse process, a clean image x_0 is reconstructed by
 230 iteratively denoising x_T . This denoising process is carried
 231 out by a denoising network ϵ_θ which is conditioned on the
 232 timestep $t \in \{1, T\}$ and optional additional conditioning c
 233 (e.g. text prompts, inpainting masks). Training loss of the
 234 denoiser is as follows: 235

$$L_{DM} = E_{x_0, \epsilon \sim \mathcal{N}(0, I), t} \|\epsilon - \epsilon_\theta(z_t, t, c)\|^2 \quad (1) \quad 236$$

Model Architecture. Building upon MirrorFusion [12],
 237 we also formulate this task as an inpainting task. Our model
 238 employs a base dual branch network similar to Brush-
 239 Net [20] and additionally uses depth map conditioning for
 240 the condition branch of BrushNet. In particular, we con-
 241 catenate the noisy latent z_t , masked image z_m , inpainting
 242 mask x_m and depth map x_d , and provide this as an input to
 243

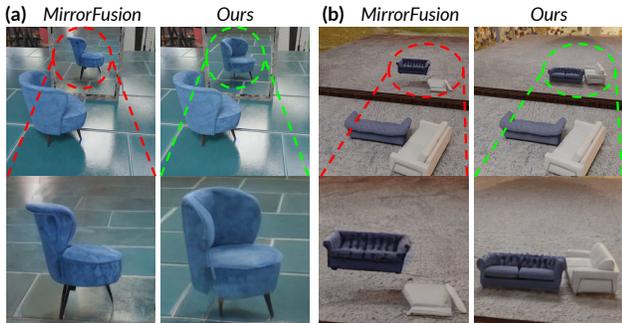


Figure 4. **Comparison on MirrorBenchV2.** The baseline fails to maintain accurate reflections and spatial consistency, showing (a) incorrect chair orientation and (b) distorted reflections of multiple objects. In contrast, our method correctly renders (a) the chair and (b) the sofas with accurate position, orientation, and structure, demonstrating superior performance.

244 the conditioning U-Net branch. Each layer of the generation
245 U-Net ϵ_i is conditioned with the corresponding layer of the
246 conditioning U-Net ϵ' with the help of zero-convolutions
247 (\mathcal{Z}) as follows:

$$\epsilon_{\theta}(z_t, t, c)_i = \epsilon_{\theta}(z_t, t, c)_i + w \cdot \mathcal{Z} \left(\epsilon'_{\theta}([z_t, z_m, x_m, x_d], t)_i \right) \quad (2)$$

248 w is the preservation scale to adjust the influence of condi-
249 tioning. We set w to be 1.0 for all our experiments. We,
250 train the model with the loss in Eq. (1).
251

252 **Training details.** We follow a 3 stage training curriculum
253 to improve the generalization of the model on real-world
254 scenes. We utilize the AdamW [25] optimizer with a learn-
255 ing rate of $1e^{-5}$ and a batch size of 4 per GPU. We train on
256 4 NVIDIA A100 GPUs in all stages.

- 257 • **Stage 1.** In the first stage, we initialize the weights of
258 both the conditioning and generation branch with the Sta-
259 ble Diffusion v1.5 checkpoint and finetune the model on
260 the single object train split of our proposed SynMirrorV2.
261 In contrast to [12], we do not keep the generation branch
262 frozen and train the model till 40,000 iterations. The
263 variation in the position and rotation in the SynMirrorV2
264 compared to SynMirror allows us to train the model for
265 longer iterations without any degradation in the genera-
266 tion quality compared to [12].
- 267 • **Stage 2.** In the second stage, we finetune the model
268 for 10,000 iterations on the multiple objects train split
269 of SynMirrorV2 to incorporate the concepts of occlusions
270 as present in realistic scenes.
- 271 • **Stage 3.** We propose a third stage training on real-world
272 data from the MSD [52] dataset for another 10,000 it-
273 erations to bridge the domain gap between synthetic and
274 real-world image inpainting.

275 In the first two stages, we use ground truth depth maps
276 and for the third stage, we generate depth maps using a

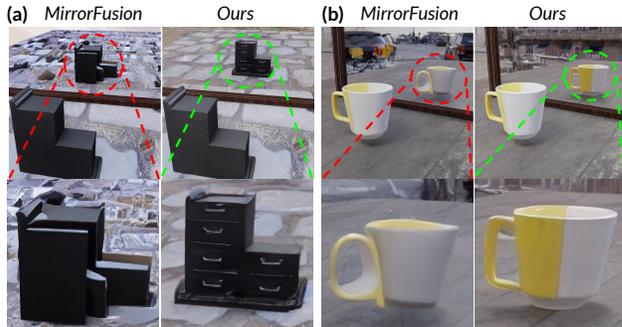


Figure 5. **Comparison on GSO [13] dataset.** In (a), the baseline method misrepresents object structure, while our method preserves spatial integrity and produces realistic reflections. In (b), the baseline yields incomplete and distorted reflections of the mug, whereas our approach generates accurate geometry, color, and detail, showing superior performance on out-of-distribution objects.

monocular depth estimator [5]. To enhance learning and
reduce reliance on text prompts, we randomly drop them
20% of the time during training, enabling the model to uti-
lize depth information better.

Inference. During inference, we use a CFG value of 7.5 and
utilize the UniPC scheduler [59] for 50 time steps. During
inference, we allow the user to provide the mask depicting
the mirror and estimate the input depth map using Depth-
Pro [5] by passing the masked image as input.

5. Experiments & Results

We discuss the evaluation strategy and compare our current
method with the previous state-of-the-art method, Mirror-
Fusion [12], referring to this as the baseline. Additionally,
we also provide ablation studies on different design choices
in Sec. 5.1.

Dataset. Compared to MirrorBench, MirrorBenchV2 con-
sists of renderings of single and multiple objects in a scene.
Additionally, we qualitatively test our method on several
images from the MSD dataset and renderings from the
Google Scanned Objects(GSO) [13] dataset. For single ob-
ject renderings, we have a total of 2,991 images, which
come from categories that are both seen and unseen dur-
ing training. We create 300 images that contain two objects
from the ABO dataset in the same scene to test the model
on generating reflections for multiple objects.

Metrics. We benchmark various methods on the quality of
the generated reflection and textual alignment of the gener-
ated image with the input prompt.

- **Reflection Generation Quality.** We evaluate reflection
quality using Peak-Signal-to-Noise ratio (PSNR), Struc-
tural Similarity (SSIM) and Learned Perceptual Image
Patch Similarity (LPIPS) [57] on the masked mirror re-
gion.

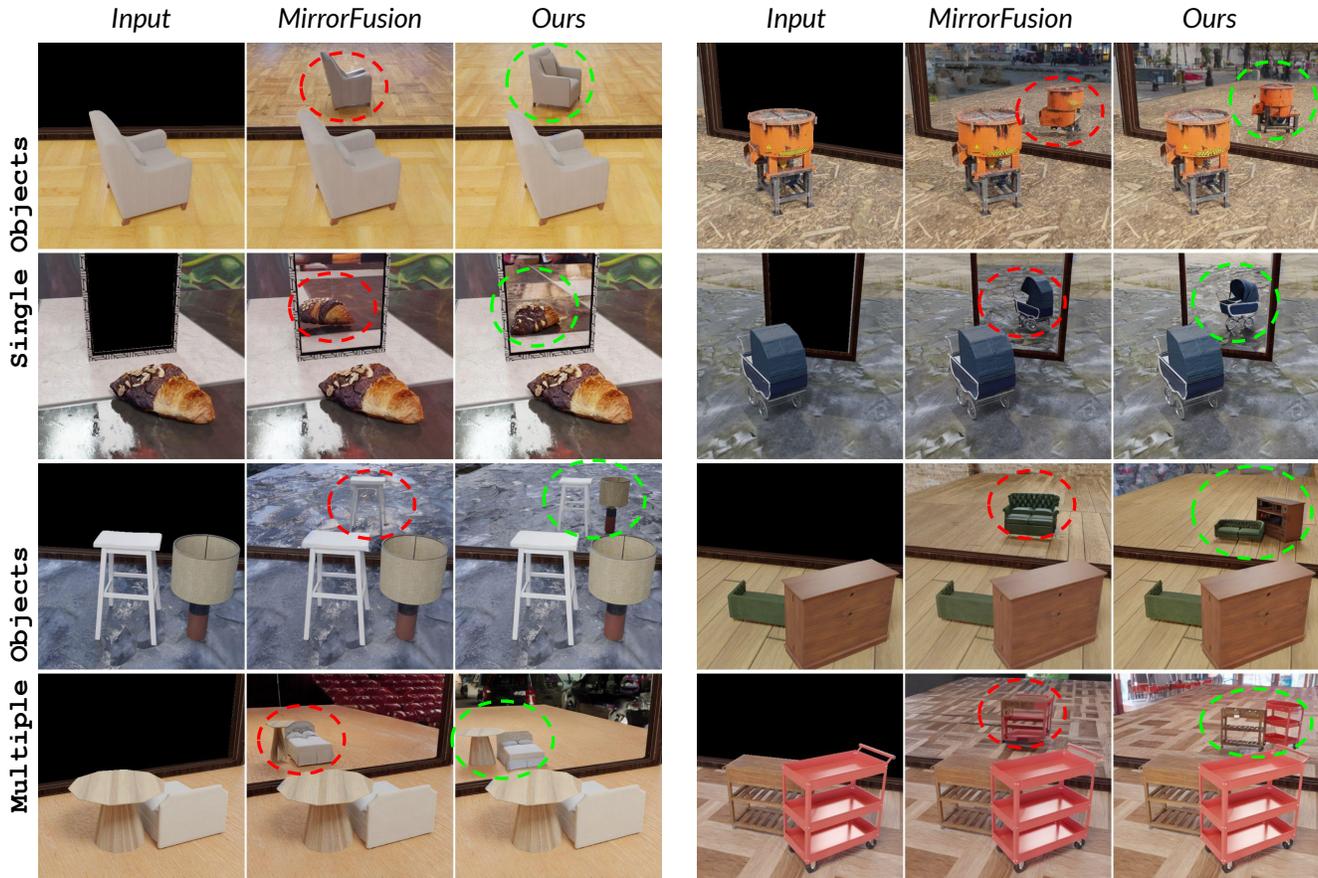


Figure 6. **Results on MirrorBenchV2.** We compare our method with the baseline MirrorFusion [12] on MirrorBenchV2. The baseline method shown struggles with pose variations, even in single-object scenes, and fails to produce accurate reflections for multiple objects. In contrast, our method handles variations in the object orientation effectively and generates geometrically accurate reflections, even in complex, multi-object scenarios.

310 • **Text Alignment.** We use CLIP [36] Similarity for assess-
311 ing textual alignment.

312 **Qualitative results on MirrorBenchV2.** In Fig. 4 (a), a
313 single chair that is slightly rotated is placed in front of a
314 mirror. We observe that the baseline method completely
315 misrepresents the chair’s orientation in the generated reflection
316 as seen in the mirror. Notice the zoomed-in region
317 where the reflection appears as if the object was cut and
318 pasted onto the mirror. In contrast, MirrorFusion 2.0 trained
319 on SynMirrorV2 accurately captures the chair’s orientation
320 in the reflection, as shown in the zoomed-in region high-
321 lighted by the green circle.

322 Fig. 4 (b), shows a scene with a white sofa rotated and
323 placed to the right of a gray sofa. The baseline method pro-
324 duces two artifacts in the reflection: 1) the gray sofa ap-
325 pears to be floating in the air, and 2) the generated reflec-
326 tion of the white sofa is completely incorrect. In contrast,
327 our method accurately generates the scene in the reflection.
328 These results demonstrate the effectiveness of our augmen-

329 tion strategies, as described in Sec. 3. We show more ex-
330 amples with both single and multiple objects in Fig. 6.

331 **Qualitative results on GSO [13].** We further evaluate
332 the generalization ability of MirrorFusion 2.0 on real-world
333 scanned objects from GSO, shown in Fig. 5. MirrorFusion
334 2.0 generates significantly more accurate and realistic reflec-
335 tions. For instance, in Fig. 5 (a), MirrorFusion 2.0 cor-
336 rectly reflects the drawer handles (highlighted in green),
337 while the baseline model produces an implausible reflection
338 (highlighted in red). Likewise, for the “White-Yellow mug”
339 in Fig. 5 (b), MirrorFusion 2.0 delivers a convincing geom-
340 etry with minimal artifacts, unlike the baseline, which fails
341 to accurately capture the object’s geometry and appearance.

342 **Qualitative results on the Real-World MSD dataset.**
343 MirrorFusion 2.0 performs well on MirrorBenchV2 and
344 real-world objects from GSO but struggles with complex
345 scenes, such as cluttered cables on a table and reflections
346 across multiple mirrors (see Fig. 7). To improve coher-
347 ence, we fine-tune it on a subset of the MSD dataset and

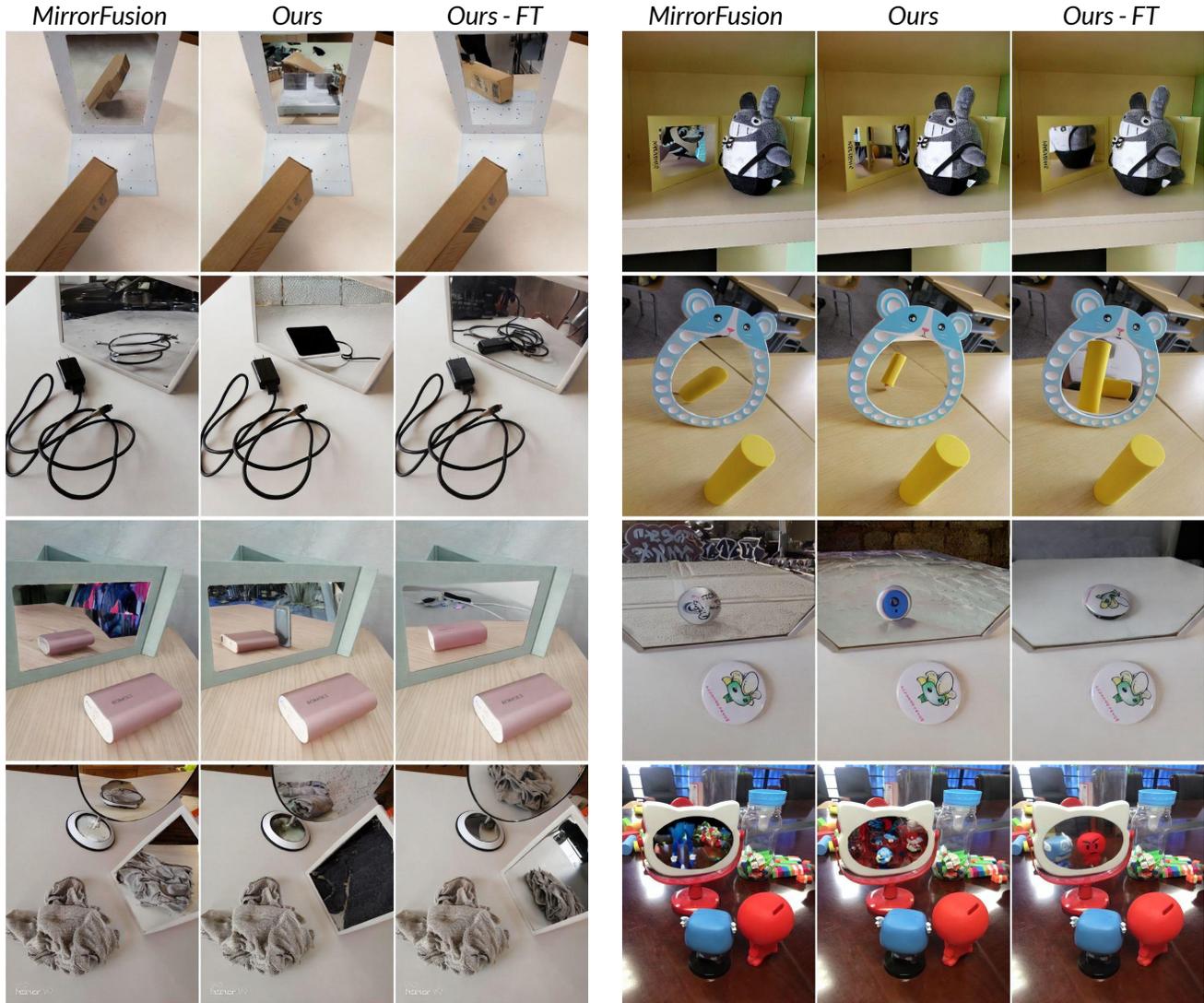


Figure 7. **Real World Scenes.** We show results for MirrorFusion [12], our method and our method fine-tuned on the MSD [52] dataset. We observe that our method can generate reflections capturing the intricacies of complex scenes, such as a cluttered cable on the table and the presence of two mirrors in a 3D scene.

348 test it on a held-out split, enhancing its ability to handle
 349 real-world scenarios. As shown in Fig. 7, this fine-tuning
 350 enables high-fidelity reflections, accurately capturing de-
 351 tails like the “black cable” on the table and the “towel” in
 352 both mirrors. These results demonstrate how our dataset im-
 353 proves diffusion models, enabling more realistic reflections
 354 in challenging settings. Fig. 7 illustrates further examples
 355 on the real-world MSD dataset.

356 **Quantitative results with baselines.** For evaluating the
 357 metrics, we generate images using four seeds for a particu-
 358 lar prompt and select the image that has the best SSIM score
 359 on the unmasked region. For a particular metric, we report
 360 the average value across MirrorBenchV2 by averaging the
 361 metric for all the selected images. Tabs. 2 and 3 show that

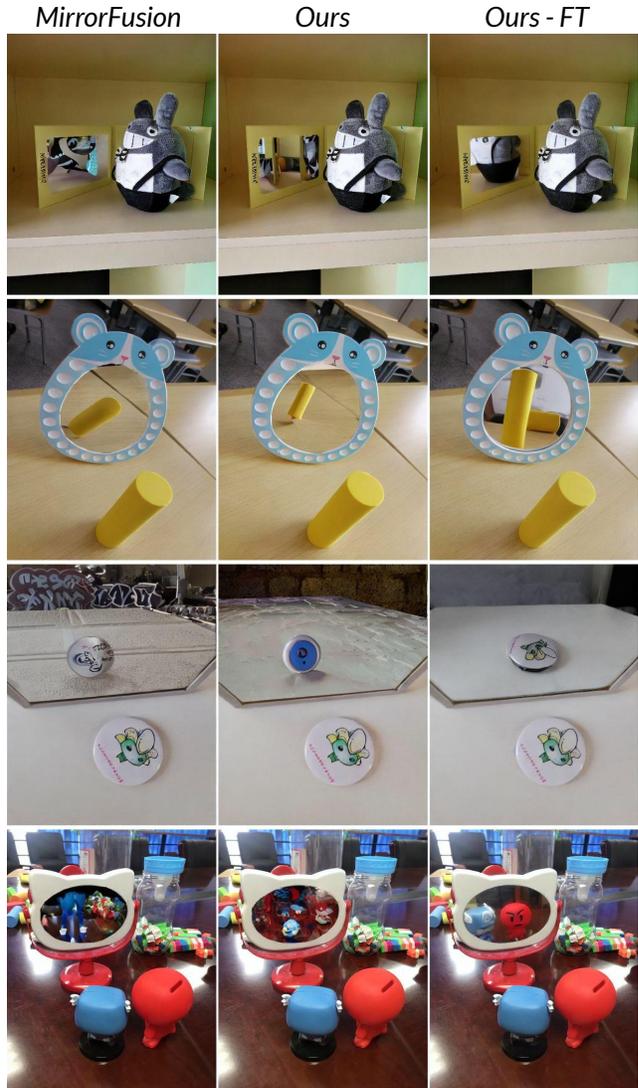


Table 2. **Single Object Reflection Generation Quality.** We compare the quantitative results between the baseline and MirrorFusion 2.0 on the **single object** split of MirrorBenchV2. The best results are shown in **bold**. This shows the effectiveness of the dataset by achieving improved scores.

Metrics	Reflection Generation Quality			Text Alignment
Models	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	CLIP Sim \uparrow
baseline [12]	18.31	0.76	0.122	26.00
Ours 40k	18.79	0.77	0.108	25.96

our method outperforms the baseline method and finetuning
 on multiple objects improves the results on complex scenes.

362

363

Table 3. **Multiple Object Reflection Generation Quality.** We compare the quantitative results between MirrorFusion 2.0 trained without multiple objects and MirrorFusion 2.0 trained with multiple objects on the **multiple object** split of MirrorBenchV2. The best results are shown in **bold**. This shows the effectiveness of finetuning further on multiple objects.

Metrics	Reflection Generation Quality			Text Alignment
Models	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	CLIP Sim \uparrow
Ours 40k	17.77	0.743	0.126	26.17
Ours 50k	18.00	0.744	0.119	26.09

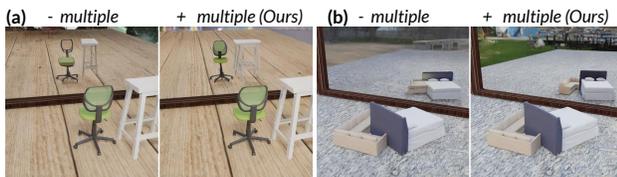


Figure 8. **Impact of adding multiple objects.** We observe that training without multiple objects leads to (a) poor reflection generation and (b) artifacts like object blending, supporting the need for finetuning the model on such scenarios.

364 **User study.** To evaluate the effectiveness of our proposed
365 strategy, we also conducted a user study where we provided
366 users with 40 different samples containing single, multiple,
367 GSO objects, and real-world generations from the baseline
368 and MirrorFusion 2.0. **84% of users preferred generations
369 from MirrorFusion 2.0 over the baseline method.** We
370 provide more details in Appendix D.5.

371 **Limitations.** Fig. 10 illustrates examples where our method
372 accurately captures overall geometry but introduces minor
373 artifacts which can be easily addressed by synthesizing ad-
374 ditional training data and fine-tuning the model.

375 5.1. Ablation Studies

376 **Impact of multiple objects dataset.** To evaluate the im-
377 pact of adding multiple objects to our dataset, we com-
378 pare MirrorFusion 2.0 with (“+ multiple”) and without (“-
379 multiple”) object training in Fig. 8. “MirrorFusion 2.0-w/o
380 multiple” struggles to generate plausible mirror reflections,
381 as evident in Fig. 8 (b), where the bed and sofa appear to
382 blend together. In contrast, “MirrorFusion 2.0-with multi-
383 ple” accurately captures the spatial relationships between
384 objects within the mirror reflection. These results highlight
385 the importance of including multiple objects in the dataset,
386 enabling the model to learn spatial relationships and effec-
387 tively handle occlusions.

388 **Ablation on architecture.** To further validate our architec-
389 tural choice, we adapt Stable Diffusion Inpainting to accept
390 depth maps as input similar to the changes made for Mir-
391 rorFusion 2.0 and train this modified model on our pro-

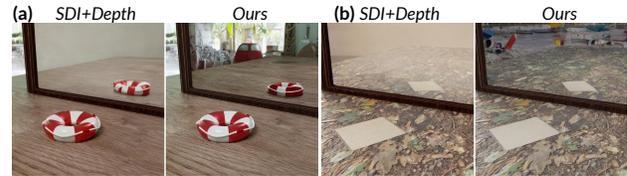


Figure 9. **Comparison with SDI+Depth baseline.** We observe color leakage issues in “SDI+Depth” generations. A dual-branch architecture proves to be a better choice, yielding superior outcomes.



Figure 10. **Limitations.** Our method performs well in multi-object scenes (more than two objects) but retains some artifacts, which can be reduced by synthesizing the dataset through the proposed data-generation pipeline and further increasing the diversity and scale.

posed dataset referring to it as “SDI+Depth”. We compare
“SDI+Depth” with MirrorFusion 2.0 in Fig. 9. While
“SDI+Depth” accurately positions objects in the mirror, it
suffers from significant artifacts, including color leakage
in contrast to MirrorFusion 2.0. We suspect that this happens
due to the early combination of the noisy latent features,
mask, and conditioning information in the initial convolu-
tion layer, restricting later layers from accessing clean fea-
tures. These findings suggest that a dual branch architecture
to provide the conditioning information separately as done
in MirrorFusion 2.0 is a better choice.

6. Conclusion

We introduce SynMirrorV2, a novel large-scale synthetic
dataset designed to advance mirror reflection generation
significantly. By employing targeted data augmentations,
we achieved robust variability in object pose, position,
and occlusion, alongside the ability to handle multi-object
scenes. Our qualitative and quantitative evaluations demon-
strate SynMirrorV2’s efficacy in reflection generation, with
promising generalization to real-world scenes using cur-
riculum training. This dataset holds substantial potential
for driving progress in various mirror-related tasks. Future
research will explore advanced data augmentation techni-
ques to enhance real-world performance further.

416

References

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

- [1] Adobe. Adobe firefly. <https://www.adobe.com/products/firefly.html>, 2025. Accessed: 2025-03-23. 8, 9
- [2] Stability AI. Stable diffusion 3.5. <https://huggingface.co/stabilityai/stable-diffusion-3.5-large>, 2025. Accessed: 2025-03-23. 2, 3
- [3] Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18208–18218, 2022. 3
- [4] Omri Avrahami, Ohad Fried, and Dani Lischinski. Blended latent diffusion. *ACM transactions on graphics (TOG)*, 42(4):1–11, 2023. 3
- [5] Aleksei Bochkovskii, Amaël Delaunoy, Hugo Germain, Marcel Santos, Yichao Zhou, Stephan R. Richter, and Vladlen Koltun. Depth pro: Sharp monocular metric depth in less than a second. *arXiv preprint arXiv: 2410.02073*, 2024. 5
- [6] Jasmine Collins, Shubham Goel, Kenan Deng, Achleshwar Luthra, Leon Xu, Erhan Gundogdu, Xi Zhang, Tomas F Yago Vicente, Thomas Dideriksen, Himanshu Arora, et al. Abo: Dataset and benchmarks for real-world 3d object understanding. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 21126–21136, 2022. 3, 4
- [7] Katherine Crowson, Stefan Andreas Baumann, Alex Birch, Tanishq Mathew Abraham, Daniel Z Kaplan, and Enrico Shippole. Scalable high-resolution pixel-space image synthesis with hourglass diffusion transformers. In *Forty-first International Conference on Machine Learning*, 2024. 3
- [8] Giannis Daras, Mauricio Delbracio, Hossein Talebi, Alexandros G. Dimakis, and Peyman Milanfar. Soft diffusion: Score matching with general corruptions. *Trans. Mach. Learn. Res.*, 2023, 2023. 3
- [9] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13142–13153, 2023. 3
- [10] Maximilian Denninger, Dominik Winkelbauer, Martin Sundermeyer, Wout Boerdijk, Markus Wendelin Knauer, Klaus H Strobl, Matthias Humt, and Rudolph Triebel. Blenderproc2: A procedural pipeline for photorealistic rendering. *Journal of Open Source Software*, 8(82):4901, 2023. 3, 4
- [11] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 3
- [12] Ankit Dhiman, Manan Shah, Rishubh Parihar, Yash Bhalgat, Lokesh R Boregowda, and R Venkatesh Babu. Reflecting reality: Enabling diffusion models to produce faithful mirror reflections, 2024. 2, 3, 4, 5, 6, 7
- [13] Laura Downs, Anthony Francis, Nate Koenig, Brandon Kinman, Ryan Hickman, Krista Reymann, Thomas B McHugh, and Vincent Vanhoucke. Google scanned objects: A high-quality dataset of 3d scanned household items. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 2553–2560. IEEE, 2022. 5, 6, 8
- [14] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12873–12883, 2021. 2
- [15] Mark Edward M. Gonzales, Lorene C. Uy, and Joel P. Ila. Designing a lightweight edge-guided convolutional neural network for segmenting mirrors and reflective surfaces. *Computer Science Research Notes*, 3301:107–116, 2023. 2
- [16] Poly Haven. Poly haven : The public 3d asset library, 2025. Accessed: 2025-03-23. 3
- [17] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 2, 3
- [18] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633–8646, 2022. 2
- [19] Yicong Hong, Kai Zhang, Jiuxiang Gu, Sai Bi, Yang Zhou, Difan Liu, Feng Liu, Kalyan Sunkavalli, Trung Bui, and Hao Tan. Lrm: Large reconstruction model for single image to 3d. *ArXiv*, abs/2311.04400, 2023. 3
- [20] Xu Ju, Xian Liu, Xintao Wang, Yuxuan Bian, Ying Shan, and Qiang Xu. Brushnet: A plug-and-play image inpainting model with decomposed dual-branch diffusion. In *European Conference on Computer Vision*, 2024. 3, 4
- [21] Kunal Kathare, Ankit Dhiman, K Vikas Gowda, Siddharth Aravindan, Shubham Monga, Basavaraja Shanthappa Vandrotti, and Lokesh R Boregowda. Instructive3d: Editing large reconstruction models with text instructions. In *Proceedings of the Winter Conference on Applications of Computer Vision (WACV)*, pages 3246–3256, 2025. 3
- [22] Black Forest Labs. Flux. <https://github.com/black-forest-labs/flux>, 2025. Accessed: 2025-03-23. 2, 3
- [23] Zhimin Li, Jianwei Zhang, Qin Lin, Jiangfeng Xiong, Yanxin Long, Xincheng Deng, Yingfang Zhang, Xingchao Liu, Minbin Huang, Zedong Xiao, Dayou Chen, Jiajun He, Jiahao Li, Wenyue Li, Chen Zhang, Rongwei Quan, Jianxiang Lu, Jiabin Huang, Xiaoyan Yuan, Xiaoxiao Zheng, Yixuan Li, Jihong Zhang, Chao Zhang, Meng Chen, Jie Liu, Zheng Fang, Weiyan Wang, Jinbao Xue, Yangyu Tao, Jianchen Zhu, Kai Liu, Sihuan Lin, Yifu Sun, Yun Li, Dongdong Wang, Mingtao Chen, Zhichao Hu, Xiao Xiao, Yan Chen, Yuhong Liu, Wei Liu, Di Wang, Yong Yang, Jie Jiang, and Qinglin Lu. Hunyuan-dit: A powerful multi-resolution diffusion transformer with fine-grained chinese understanding, 2024. 2
- [24] Jiaying Lin, Guodong Wang, and Rynson W.H. Lau. Progressive mirror detection. In *Proc. CVPR*, 2020. 2
- [25] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2017. 5
- [26] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting

- 531 using denoising diffusion probabilistic models. In *Proceed-*
532 *ings of the IEEE/CVF conference on computer vision and*
533 *pattern recognition*, pages 11461–11471, 2022. 3
- [27] Hayk Manukyan, Andranik Sargsyan, Barsegh Atanyan,
534 Zhangyang Wang, Shant Navasardyan, and Humphrey Shi.
535 Hd-painter: high-resolution and prompt-faithful text-guided
536 image inpainting with diffusion models. *arXiv preprint*
537 *arXiv:2312.14091*, 2023. 2, 3
- [28] Haiyang Mei, Bo Dong, Wen Dong, Pieter Peers, Xin Yang,
539 Qiang Zhang, and Xiaopeng Wei. Depth-aware mirror seg-
540 mentation. In *Proceedings of the IEEE/CVF conference on*
541 *computer vision and pattern recognition*, pages 3044–3053,
542 2021. 2
- [29] Oscar Michel, Anand Bhattad, Eli VanderBilt, Ranjay Kr-
544 ishna, Aniruddha Kembhavi, and Tanmay Gupta. Object
545 3dit: Language-guided 3d-aware image editing. *Advances*
546 *in Neural Information Processing Systems*, 36, 2024. 2, 3
- [30] Sicheng Mo, Fangzhou Mu, Kuan Heng Lin, Yanli Liu,
548 Bochen Guan, Yin Li, and Bolei Zhou. Freecontrol:
549 Training-free spatial control of any text-to-image diffusion
550 model with any condition. In *Proceedings of the IEEE/CVF*
551 *Conference on Computer Vision and Pattern Recognition*,
552 pages 7465–7475, 2024. 2
- [31] Rishubh Parihar, Abhijnya Bhat, Abhipsa Basu, Saswat
554 Mallick, Jogendra Nath Kundu, and R Venkatesh Babu.
555 Balancing act: distribution-guided debiasing in diffusion
556 models. In *Proceedings of the IEEE/CVF conference on*
557 *computer vision and pattern recognition*, pages 6668–6678,
558 2024. 3
- [32] Rishubh Parihar, Harsh Gupta, Sachidanand VS, and
560 R Venkatesh Babu. Text2place: Affordance-aware text
561 guided human placement. In *European Conference on Com-*
562 *puter Vision*, pages 57–77. Springer, 2024. 3
- [33] Rishubh Parihar, VS Sachidanand, Sabariswaran Mani, Te-
564 jan Karmali, and R Venkatesh Babu. Precisecontrol: En-
565 hancing text-to-image diffusion models with fine-grained at-
566 tribute control. In *European Conference on Computer Vision*,
567 pages 469–487. Springer, 2024. 3
- [34] William Peebles and Saining Xie. Scalable diffusion models
569 with transformers. In *Proceedings of the IEEE/CVF Inter-*
570 *national Conference on Computer Vision*, pages 4195–4205,
571 2023. 3
- [35] Dustin Podell, Zion English, Kyle Lacey, Andreas
573 Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and
574 Robin Rombach. Sdxl: Improving latent diffusion mod-
575 els for high-resolution image synthesis. *arXiv preprint*
576 *arXiv:2307.01952*, 2023. 2
- [36] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya
578 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,
579 Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning
580 transferable visual models from natural language supervi-
581 sion. In *International conference on machine learning*, pages
582 8748–8763. PMLR, 2021. 6
- [37] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu,
584 and Mark Chen. Hierarchical text-conditional image gener-
585 ation with clip latents. *arXiv preprint arXiv:2204.06125*, 1
586 (2):3, 2022. 3
- [38] Robin Rombach, Andreas Blattmann, Dominik Lorenz,
588 Patrick Esser, and Björn Ommer. High-resolution image
589 synthesis with latent diffusion models. In *Proceedings of*
590 *the IEEE/CVF conference on computer vision and pattern*
591 *recognition*, pages 10684–10695, 2022. 2, 3
- [39] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee,
593 Jonathan Ho, Tim Salimans, David Fleet, and Mohammad
594 Norouzi. Palette: Image-to-image diffusion models. In
595 *ACM SIGGRAPH 2022 conference proceedings*, pages 1–10,
596 2022. 3
- [40] Chitwan Saharia, William Chan, Saurabh Saxena, Lala
598 Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour,
599 Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans,
600 et al. Photorealistic text-to-image diffusion models with deep
601 language understanding. *Advances in neural information*
602 *processing systems*, 35:36479–36494, 2022. 3
- [41] Ayush Sarkar, Hanlin Mai, Amitabh Mahapatra, Svetlana
604 Lazebnik, David A Forsyth, and Anand Bhattad. Shad-
605 ows don’t lie and lines can’t bend! generative models don’t
606 know projective geometry... for now. In *Proceedings of*
607 *the IEEE/CVF Conference on Computer Vision and Pattern*
608 *Recognition*, pages 28140–28149, 2024. 2, 3
- [42] Prafull Sharma, Varun Jampani, Yuanzhen Li, Xuhui Jia,
610 Dmitry Lagun, Fredo Durand, Bill Freeman, and Mark
611 Matthews. Alchemist: Parametric control of material proper-
612 ties with diffusion models. In *Proceedings of the IEEE/CVF*
613 *Conference on Computer Vision and Pattern Recognition*,
614 pages 24130–24141, 2024. 3
- [43] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan,
616 and Surya Ganguli. Deep unsupervised learning using
617 nonequilibrium thermodynamics. In *International confer-*
618 *ence on machine learning*, pages 2256–2265. PMLR, 2015.
619 3
- [44] Tianyu Sun, Guodong Zhang, Wenming Yang, Jing-Hao
621 Xue, and Guijin Wang. Trosd: A new rgb-d dataset for trans-
622 parent and reflective object segmentation in practice. *IEEE*
623 *Transactions on Circuits and Systems for Video Technology*,
624 33(10):5721–5733, 2023. 2
- [45] Jiaqi Tan, Weijie Lin, Angel X Chang, and Manolis Savva.
626 Mirror3D: Depth refinement for mirror surfaces. In *Proceed-*
627 *ings of the IEEE Conference on Computer Vision and Pattern*
628 *Recognition (CVPR)*, 2021. 2
- [46] Rishi Upadhyay, Howard Zhang, Yunhao Ba, Ethan Yang,
630 Blake Gella, Sicheng Jiang, Alex Wong, and Achuta
631 Kadambi. Enhancing diffusion models with 3d perspec-
632 tive geometry constraints. *ACM Transactions on Graphics*
633 *(TOG)*, 42(6):1–15, 2023. 2, 3
- [47] A Vaswani. Attention is all you need. *Advances in Neural*
635 *Information Processing Systems*, 2017. 3
- [48] Bram Wallace, Akash Gokul, and Nikhil Naik. Edict: Exact
637 diffusion inversion via coupled transformations. In *Proceed-*
638 *ings of the IEEE/CVF Conference on Computer Vision and*
639 *Pattern Recognition*, pages 22532–22541, 2023. 3
- [49] Daniel Winter, Matan Cohen, Shlomi Fruchter, Yael Pritch,
641 Alex Rav-Acha, and Yedid Hoshen. Objectdrop: Bootstrap-
642 ping counterfactuals for photorealistic object removal and in-
643 sertion, 2024. 2, 3

- 645 [50] Xiaojun Wu, Dixiang Zhang, Ruyi Gan, Junyu Lu, Ziwei
646 Wu, Renliang Sun, Jiaying Zhang, Pingjian Zhang, and Yan
647 Song. Taiyi-diffusion-xl: Advancing bilingual text-to-image
648 generation with large vision-language model support, 2024.
649 2
- 650 [51] Shaoan Xie, Zhifei Zhang, Zhe Lin, Tobias Hinz, and Kun
651 Zhang. Smartbrush: Text and shape guided object inpainting
652 with diffusion model. In *Proceedings of the IEEE/CVF Con-
653 ference on Computer Vision and Pattern Recognition*, pages
654 22428–22437, 2023. 3
- 655 [52] Xin Yang, Haiyang Mei, Ke Xu, Xiaopeng Wei, Baocai Yin,
656 and Rynson W.H. Lau. Where is my mirror? In *The IEEE
657 International Conference on Computer Vision (ICCV)*, 2019.
658 2, 5, 7, 8, 9
- 659 [53] Fulong Ye, Guangyi Liu, Xinya Wu, and Ledell Yu Wu. Alt-
660 diffusion: A multilingual text-to-image diffusion model. In
661 *AAAI Conference on Artificial Intelligence*, 2023. 2
- 662 [54] Hu Ye, Jun Zhang, Sibao Liu, Xiao Han, and Wei Yang. Ip-
663 adapter: Text compatible image prompt adapter for text-to-
664 image diffusion models. *arXiv preprint arXiv:2308.06721*,
665 2023. 2
- 666 [55] Junyi Zeng, Chong Bao, Rui Chen, Zilong Dong, Guofeng
667 Zhang, Hujun Bao, and Zhaopeng Cui. Mirror-nerf: Learn-
668 ing neural radiance fields for mirrors with whitted-style ray
669 tracing. In *Proceedings of the 31st ACM International Con-
670 ference on Multimedia*, pages 4606–4615, 2023. 2
- 671 [56] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding
672 conditional control to text-to-image diffusion models. In
673 *Proceedings of the IEEE/CVF International Conference on
674 Computer Vision*, pages 3836–3847, 2023. 2
- 675 [57] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shecht-
676 man, and Oliver Wang. The unreasonable effectiveness of
677 deep features as a perceptual metric. In *Proceedings of the
678 IEEE conference on computer vision and pattern recogni-
679 tion*, pages 586–595, 2018. 5
- 680 [58] Shihao Zhao, Dongdong Chen, Yen-Chun Chen, Jianmin
681 Bao, Shaozhe Hao, Lu Yuan, and Kwan-Yee K Wong. Uni-
682 controlnet: All-in-one control to text-to-image diffusion
683 models. *Advances in Neural Information Processing Sys-
684 tems*, 36, 2024. 2
- 685 [59] Wenliang Zhao, Lujia Bai, Yongming Rao, Jie Zhou, and
686 Jiwen Lu. Unipc: A unified predictor-corrector framework
687 for fast sampling of diffusion models. *Advances in Neural
688 Information Processing Systems*, 36, 2024. 5
- 689 [60] Hongkai Zheng, Weili Nie, Arash Vahdat, and Anima
690 Anandkumar. Fast training of diffusion models with masked
691 transformers. *Trans. Mach. Learn. Res.*, 2024, 2023. 3
- 692 [61] Junhao Zhuang, Yanhong Zeng, Wenran Liu, Chun Yuan,
693 and Kai Chen. A task is worth one word: Learning with
694 task prompts for high-quality versatile image inpainting. In
695 *European Conference on Computer Vision*, 2023. 3