
HarmonyDreamer: Rectifying and Harmonizing Dynamic Objects in Backgrounds for 4D Video Generation

Anonymous Author(s)

Affiliation

Address

email

Abstract

Recent advancements in score distillation have markedly enhanced the generation of both static and dynamic 3D tasks by effectively harnessing the power of 2D pretrained diffusion models for artistic applications. However, traditional distillation methods primarily concentrate on object-centric representations without adequately addressing the integration of objects with their backgrounds, which includes matching rendering styles and motions. This paper presents a new framework that incorporates a background-enhanced distillation score along with an iterative editing-reconstruction process tailored for 4D video generation from basic dynamic 4D Gaussian splatting. Our experiments demonstrate that our method not only excels in producing high-fidelity 4D videos that harmoniously blend objects with varied backgrounds, but also facilitates the creation of multi-view videos. Additionally, our approach outperforms existing baselines in both spatio-temporal and semantic consistency, demonstrating its effectiveness in 4D video generation. More video results can be viewed in our project page: <https://harmonydreamer.github.io/web/>.

1 Introduction

Recent advancements in generative models, particularly text-to-image and text-to-video diffusion models, have paved the way for a new era of content creation. Studies show that these models have achieved remarkable success in generating various digital formats, encompassing 2D images [34, 31, 30], videos [3, 15, 5], and 3D scenes [28, 43, 40]. Beyond producing contents from above-mentioned modalities, considerable effort has been directed towards developing more sophisticated content such as videos and 4D scenes. The majority of research focuses on text-to-4D [36, 2, 46, 24], video-to-4D [8, 33, 49], and image-to-4D approaches [33, 51, 47]. However, most of these methods are far from practical usage as they can only do object-centric generation without integrating these objects under arbitrary background with matching rendering styles, shadows, and motion.

In this paper, we introduce the inversion score for both image and video harmonization. We also propose a workflow as illustrated in Fig. 1 that processes rough 3D or 4D inputs for multi-view video generation. Our pipeline begins by stacking rough Gaussian splatting 3D inputs onto a layered image background. We then distill the background semantics to the front 3D input, enabling both rectification and harmonization using our proposed image layered inversion score. Following this, the edited 3D result can be attached to a motion network to produce a rough 4D output, which can be further rectified using our proposed video inversion score.

In the text-to-3D and editing tasks, score-based distillation [28] methods have been extensively studied to distill object-centric 3D objects that align with monocular views or text prompts. Another

35 approach, starting with Instruct-Nerf2Nerf [13], proposes a 3D scene editing method that continually
36 trains a 3D representation while editing partial multi-view rendering images using off-the-shelf image
37 editing method InstructPix2Pix [4]. This variation of the score distillation method can edit more
38 general 3D scenes with high resolution but is time-consuming and limited to text-guided editing.

39 To harmonize visual artifacts and background semantics, a distillation method must be flexible,
40 not constrained to object-centric information, and robust against unsatisfactory inputs with layered
41 artifacts and camera intrinsic inconsistencies. Our method demonstrates that the text-to-image off-the-
42 shelf model can perform harmonization using layered inversion. Additionally, we show that the same
43 scheme works for video inversion, enabling 4D harmonization, including rectification and interaction
44 dynamics discovery. Our contributions can be listed as follows:

- 45 • We propose a versatile 3D editing method for harmonization and rectification that is not
46 constrained to object-centric scenes.
- 47 • We explore the semantics derived from layered backgrounds, in addition to those from text
48 prompts.
- 49 • We develop an inversion score that is applicable to both images and videos, extending 3D
50 harmonization results to 4D with arbitrary input motions. This allows for the decoupling of
51 editing appearance and motion style in multi-view video creation.

52 2 Related Work

53 **Text-to-3D Generation** The limited availability of extensive text-annotated 3D training datasets
54 hampers progress in 3D content generation. To overcome such limitations, DreamField [17] combined
55 neural rendering [27] and multi-modal representation [29] to directly generate 3D assets without
56 3D supervision. However, most of the recent advancement utilized the pretrained Text-to-Image
57 models such as Stable Diffusion [34] to serve as a strong generative prior for text-to-3D generation.
58 Dreamfusion [28] pioneered such pipeline by translating these image-based priors into coherent
59 3D models with an iterative distillation approach, named as Score Distillation Sampling (SDS). It
60 leverages and back-propagates random view samples from a 2D diffusion model to optimize 3D
61 representations such as NeRF [27], through a mode seeking manner. Sparked by the novel SDS loss,
62 research in this area has concentrated on several key objectives. Several of them aim to introducing
63 better update rule for score distillation to overcome the artifacts in the generated 3D assets such as
64 the low-diversity, over-saturation, and the *Janus* problem [43, 42, 41, 23, 52], while the other choose
65 to optimize over more efficient 3D representation such as Gaussian Splatting [22] to accelerate the
66 content generation process [40, 48, 7, 11].

67 **4D Generation** 4D content Generation targets to generate dynamic and spatial-temporal consistent
68 3D scenes viewed at arbitrary angles. While previous methods focus on generating dynamic object
69 of specific classes such as animals [18] and humans [32, 26, 50], Make-A-Video3D (MAV3D) [36]
70 explores this area by extending the SDS-based approach with 4D Dynamic NeRF and HexPlane repre-
71 sentation [6] to generate arbitrary objects and flexibly model large scene motions. Consistent4D [20]
72 proposes the interpolation-driven consistency loss with per-frame SDS to generate dynamic objects
73 under uncalibrated monocular videos. 4Dfy [2] blends the SDS supervision from both Text-to-image
74 and Text-to-video diffusion models to alleviates the tradeoff between appearance, 3D structure,
75 and motion in 4D assets generation. 4DGen [49] combined DDIM-sampled pseudo labels with
76 SDS to ensure spatial-temporal consistent and high-resolution generation. TC4D [1] proposed a
77 trajectory-conditioned method to supervise the global and local scene animation individually via
78 video score distillation sampling. Most of the existing works focus on object-centric generation
79 without background information. Our work, however, aims to fill the gap by incorporating a novel
80 and enhanced background information distillation score.

81 3 Method

82 3.1 Preliminaries

83 **Sampling Method** Diffusion-based generative models, applicable to both images and videos, utilize
84 a score matching function for training. These models can be sampled using expensive stochastic
85 methods, such as the Denoising Diffusion Probabilistic Models (DDPM) [16]. An advancement in
86 this field, the Denoising Diffusion Implicit Models (DDIM) [37], accelerates the sampling process by

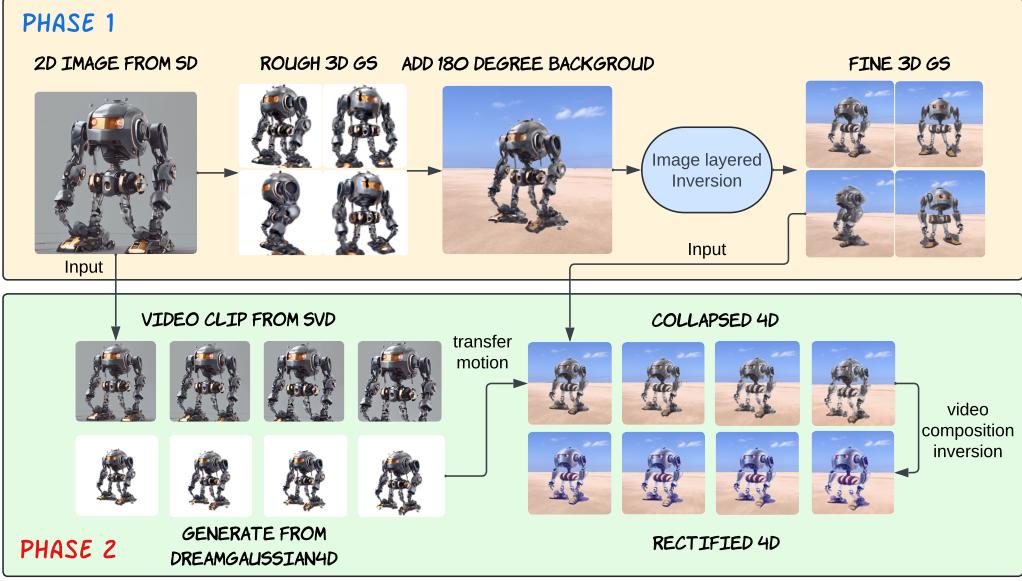


Figure 1: Our Method contains two phases. Phase 1 takes a 2D generated model from Stable diffusion, generate the low-quality 3D Gaussian splatting model and refine it with by adding background and our image layered inversion. Phase 2 takes the generated 2D image to get the 3D reference video from Stable Video diffusion. Then we use DreamGaussian4D to learn the motion from the video and transfer to the refined 3D object and further improve the harmonization with our video composition inversion to get the final output.

87 conceptualizing it as an Ordinary Differential Equation (ODE) probability flow. This method has
 88 gained popularity as the predominant sampling technique for pre-trained image generation models.
 89 Further developments include the Enhanced Diffusion Model (EDM), which explains both Stochastic-
 90 Multiscale-Langevin Dynamics (SMLD) [38] and DDIM within a single ODE sampling framework,
 91 allowing for flexible selection of sampling time steps. Owing to its continuous nature, the EDM
 92 schedule proves to be particularly well-suited for Stable Video Diffusion and is thus commonly
 93 employed as the default method.

94 During the inference phase of text-to-image diffusion model, the deterministic sampling of DDIM [37]
 95 follows an ODE-like deterministic trajectory:

$$\frac{z_{t-1}}{\sqrt{\alpha_{t-1}}} = \frac{z_t}{\sqrt{\alpha_t}} + \left(\frac{\sqrt{1-\alpha_{t-1}}}{\sqrt{\alpha_{t-1}}} - \frac{\sqrt{1-\alpha_t}}{\sqrt{\alpha_t}} \right) \epsilon_\theta(z_t, t) \quad (1)$$

96 And during the inference phase of image-to-video diffusion model, the deterministic sampling of
 97 EDM is as follows:

$$z_{t-1} = z_t + (\sigma(t-1) - \sigma(t)) \frac{D(z; \sigma(t))}{\sigma(t)}, \quad (2)$$

98 where z_t represents the latent sampled at time step t , α_t and $\sigma(t)$ are the time-dependent diffusion
 99 coefficients, and ϵ_θ represents the denoiser module parameterized by a neural network. We follow the
 100 same convention in SVD default sampling, setting $\sigma(t) = t$ and $s(t) = 1$.

101 **Gaussian Splatting and Dynamic Gaussian splatting** 3D Gaussian Splatting [22] aims to recon-
 102 struct a 3D scene by initializing a set of 3D Gaussian primitives and minimizing the photometric
 103 discrepancy between a set of input images $\{I\}_m$ paired with known camera views $\{V\}_m$ and their
 104 corresponding rendered images $\{I_r\}_m$. In practice, we denote one Gaussian splatting representation
 105 as G , with G^t being the Gaussian splatting at time frame t . Here, G_c^t represents colors, G_s^t denotes
 106 the scales, G_r^t indicates the rotations, and G_x^t signifies the Gaussian splats' position centers. We
 107 denote $G(v)$ as the rendered image for a given camera view v and $G_a(v)$ as the corresponding alpha
 108 mask. For simplicity, in the following text, we will abbreviate Gaussian Splatting as GS.

109 **3.2 Inversion Score**

110 **3.2.1 ODE Sampling Inversion**

111 z_0 represents the encoded form of the provided real image or video. Rather than the typical z_T
 112 to z_0 sampling, the inversion process operates inversely, transitioning from z_0 to z_T . z_T can be
 113 used to extract rough semantic information for further editing purposes. To leverage the inverse
 114 characteristics of the ODE path, we can apply a similar inversion scheme to the DDIM[37] sampler
 115 within the text-to-image model:

$$\frac{z_{t+1}}{\sqrt{\alpha_{t+1}}} = \frac{z_t}{\sqrt{\alpha_t}} + \left(\frac{\sqrt{1 - \alpha_{t+1}}}{\sqrt{\alpha_{t+1}}} - \frac{\sqrt{1 - \alpha_t}}{\sqrt{\alpha_t}} \right) \epsilon_\theta(z_t, t, C), \quad (3)$$

116 and EDM[21] sampler with image-to-video model:

$$z_{t+1} = z_t + (\sigma(t+1) - \sigma(t)) \frac{D(z; \sigma(t), C)}{\sigma(t)}. \quad (4)$$

117 **3.2.2 Explicit Modified Inversion**

118 To streamline the notation in our discussion, we use $\text{Inv}(I, C_{text}, t(r))$ to denote the inversion
 119 process applied to an image I , with the textual condition C_{text} , and reversing to the fraction $t = r \cdot T$,
 120 r being the ratio of inversion steps completed over T being the scheduled total sampling steps. In a
 121 similar vein, $\text{Inv}(V, C_{image}, t(r))$ refers to the inversion process for a video sample V , this time
 122 with an image-based condition C_{image} , and using EDM inversion.

123 Inversion process offers additional advantages [14]. [9] applies the clean data estimation $z_{0|t} =$
 124 $\frac{1}{\sqrt{\alpha_t}} (z_t - \sqrt{1 - \alpha_t} \epsilon_\theta(z_t, t))$ from Tweedie's formula [10] to obtain a point estimate of the noisy
 125 data. In our case, we use $\text{Inv}(I, C_{text}, t)$ to denote the inversed latents with t inversion steps,
 126 and further $\text{Inv}(I, C_{text}, t)_{0|t}$ to denote the clean data estimation from this inversed latents
 127 by $\text{Inv}(I, C_{text}, t)_{0|t} = \frac{1}{\sqrt{\alpha_t}} (\text{Inv}(I, C_{text}, t) - \sqrt{1 - \alpha_t} \epsilon_\theta(\text{Inv}(I, C_{text}, t), t))$. In Fig. 2, we
 128 present the visualizations of clean data estimation from the inversed latents, demonstrating the capture
 129 of basic geometric information from the input sample. This serves as guidance for both the reversion
 130 path and the subsequent 3D reconstruction.

131 Our reversion steps modify the direct sampling method in Eq. 1 and Eq. 2 as follows

$$\hat{z}_{0|t} = \mathcal{K}(t) \cdot z_{0|t} + (1 - \mathcal{K}(t)) \cdot \text{Inv}(I, C_{text}, t)_{0|t} \quad (5)$$

$$\hat{\epsilon}_\theta(z_t, t, C, \emptyset) = w \cdot \epsilon_\theta(z_t, t, C) + (1 - w) \cdot \epsilon_\theta(z_t, t, \emptyset) \quad (6)$$

$$\hat{z}_t = \sqrt{\alpha_t} \cdot \hat{z}_{0|t} + \sqrt{1 - \alpha_t} \hat{\epsilon}_\theta(z_t, t, \emptyset) \quad (7)$$

$$\frac{z_{t+1}}{\sqrt{\alpha_{t+1}}} = \frac{\hat{z}_t}{\sqrt{\alpha_t}} + \left(\frac{\sqrt{1 - \alpha_{t+1}}}{\sqrt{\alpha_{t+1}}} - \frac{\sqrt{1 - \alpha_t}}{\sqrt{\alpha_t}} \right) \hat{\epsilon}_\theta(z_t, t, C), \quad (8)$$

132 Note that, as DDIM is a special case of EDM, we use DDIM notation to illustrate the modified
 133 reversion sampling. w is the default Classifier-free Guidance [16] parameter in sampling process.
 134 $\mathcal{K}(t)$ is a decreasing weight acts as a factor that blends the inversion and reversion paths using
 135 the recorded clean data estimated during inversion. We also present visualizations of the clean
 136 data estimation along the reversion path to demonstrate its effect. For simplicity, we denote the
 137 process of iterating the modified reversion sampling using Eqs. 5, 6, and 8 from timestep t to 0 as
 138 $\text{Rev}(I, C, \omega, \mathcal{K})$. Illustrated in Fig. 2, using the first image on the left of the first row as input, the
 139 visuals of $\text{Rev}(I, C, \omega, \mathcal{K})$ would be the first image on the left in the second row. In practice, we use
 140 15 inversion steps out of 50 total DDIM steps for the text-to-image model, and 10 inversion steps out
 141 of 25 total steps for the image-to-video model.

142 **3.3 Image Layered Inversion**

143 For image inversion, we use a cropped object, layer it over a background image, and apply the
 144 inversion method as described. In our practice, this object could be the rendered Gaussian Splatting
 145 and its mask. As shown in Fig. 1, the rendered robot model layered on the background demonstrates

146 rectification. In Fig. 5, stacking the same rough 3D input onto the background with the same
 147 text prompt results in semantically different outcomes after iterations, with the initial rough model
 148 rectified and harmonized. Fig. 7 illustrates the gradual evolution of the editing results, highlighting
 149 our method’s stability and advantages for downstream applications.

150 **3.4 Video Composition Inversion**

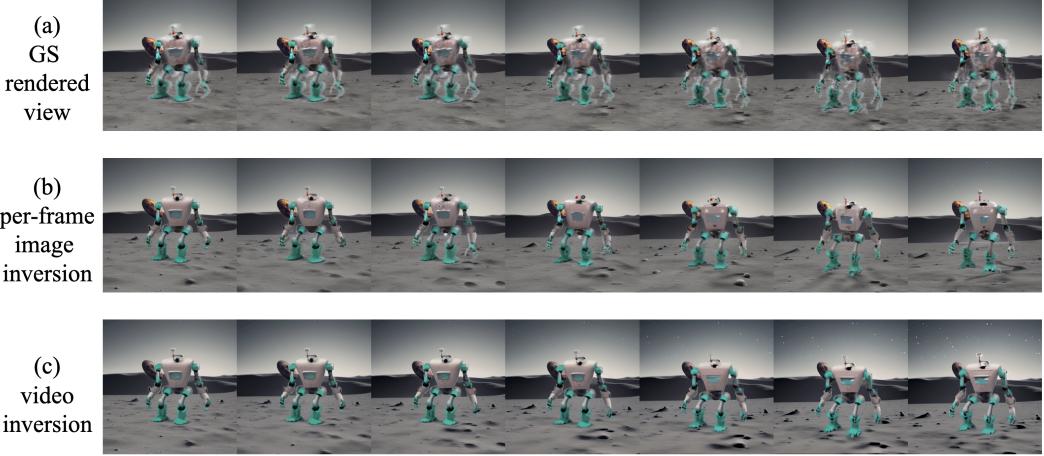


Figure 3: (a) Rendered view of temporal Gaussian Splatting layered with user-specified background, showing 7 frames extracted in sequence from a total of 14 frames. (b) Per-frame image harmonization result using layered image inversion as described in Section 3.3. (c) Temporal harmonization result using video composition inversion as described in Section 3.4.

151 Our inversion method can also be applied to
 152 the sampler in Stable Video Diffusion[3]. For
 153 clarity, we slightly extend our notation, using
 154 $\text{Rev}(I, C_{\text{image}})$ to represent the video inversion
 155 and reversion process. However, directly applying
 156 video inversion to rendered frames results
 157 in deteriorated quality. To address this problem,
 158 we propose **video composition inversion**. First,
 159 we apply image inversion to the rendered frames
 160 as shown in Fig. 3(a), using image layered in-
 161 version to obtain Fig. 3(b). Then, we apply video inversion to the intermediate output to achieve
 162 Fig. 3(c). Our approach effectively balances the temporal consistency and frame-wise quality.

163 **3.5 Scheduled Optimization**

164 Our method initiates with a preliminary Gaussian splatting G and a predefined set of views, delineated
 165 by specific ranges of horizon and azimuth. As discussed in Section 3.3, our inversion score incorpo-
 166 rates semantics from both text prompts and background information. We stack view v on the mapped
 167 background as $I_b(v)$. For multi-view consistency, the mapping from view v to background image I_b
 168 should be continuous, though not necessarily with accurate 3D geometry, as our harmonization steps
 169 will correct any layered artifacts. We propose a simple solution by cropping a small image from a
 170 larger one with changes in horizon and azimuth.

171 **3.5.1 Scheduled Gaussian Splatting Optimization**

172 We denote raw input without any editing denoted as $G(i)$ for $i = 0$, where i represents the i -th
 173 Gaussian splatting after i iterations of editing. A single editing step is described as the following:

174 We first render n views $G(i, v), v \in V$, randomly from the predefined visible range, and then layer
 175 the corresponding mapped background to each view. For example, given view v , the layered rendering

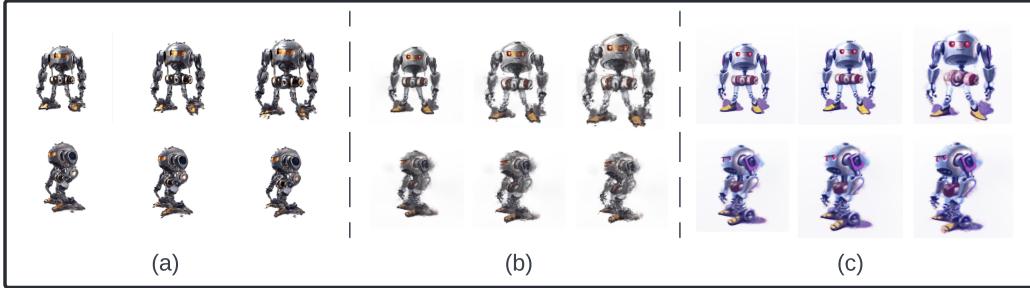


Figure 4: Quality comparison for three frames and two different views. (a) Dreamgaussian4D output. (b) Motion transfer to refined 3D. (c) Rectified 4D with our method.

176 image would be:

$$I_{\text{layer}}(i, v) = G_{\mathbf{a}}(v) \cdot G(i, v) + (1 - G_{\mathbf{a}}(v)) \cdot I_b(v).$$

177 Next, we apply inversion to m randomly selected views out of n . For simplicity, we denote the first
178 m views as the selected ones for inversion $v = 0, 1, \dots, m$. Therefore, the optimization objective for
179 $G_{\mathbf{a}}^t(i, v)$ can be formulated as:

$$\begin{aligned} L = \min_{G(i)} \sum_{v=1}^m & \| (G_{\mathbf{a}}(i, v) \cdot G(i, v) + (1 - G_{\mathbf{a}}(v)) \cdot I_b(v)) - \text{Rev}(I_{\text{layer}}(i, v), C_{\text{Text}}, t(r(v))) \|^2 \\ & + \sum_{v=m+1}^n \| (G_{\mathbf{a}}(i, v) \cdot G(i, v) + (1 - G_{\mathbf{a}}(v)) \cdot I_b(v)) - I_{\text{layer}}(i, v) \|^2. \end{aligned} \quad (9)$$

180 Note that as inversion proceeds, the editing effect applied to the object becomes more divergent
181 from the input. Due to the orbit camera intrinsics, when r is larger, the projected object appears
182 smaller. To adapt and balance the editing effect across different views, we use the empirical formula
183 $r(v) = 0.7 - \text{radius}(v)/10$, where $\text{radius}(v)$ is the radius of view v .

184 However, continuous optimization of a fixed representation can make the Gaussian splatting hard to
185 converge and prone to degradation, preventing recovery in the editing process. Therefore, we propose
186 a rescheduled optimization pipeline for Gaussian splatting reconstruction, which shows robustness in
187 handling multi-view images with subtle inconsistencies, making it suitable for our inversion score.
188 The optimization update step is described as follows:

$$G(i+1)^{(s)} = \begin{cases} G(i)^{(s)} + \mathcal{J}(s)G(i) & \text{if } s = 0 \\ G(i+1)^{(s)} + \mathcal{J}(s)G(i+1) & \text{if } s > 0 \end{cases} \quad (10)$$

$$G(i+1)^{(s+1)} = G(i+1)^{(s)} - \eta \nabla_{G(i+1)} L \quad (11)$$

189 We use $\mathcal{J}(s)$ as a gradually decreasing factor to control the noise added to the optimization. Note that
190 when $s = 0$, the noise is added to the Gaussian splatting $G(i)$ from the previous editing iteration i .

191 Here, we use $G(i+1)$ to represent all learnable attributes in the $i+1$ -th Gaussian splatting. In
192 practice, these attributes can include colors G_c^t only, or all learnable attributes: colors $G_c^t(i+1)$,
193 scales $G_s^t(i+1)$, rotations $G_r^t(i+1)$, and point centers $G_{\mathbf{x}}^t(i+1)$.

194 In the experiment, we set $n = 24$, $m = 4$, and the optimization steps for each editing iteration to 400.
195 After this optimization, we denote the updated Gaussian splatting as $G(i+1)$. A new image set is
196 then rendered from $G(i+1)$, and the process repeats.

197 3.5.2 Scheduled Dynamic Gaussian Splatting Optimization

198 Note that we are not utilizing any additional temporal adjustment structures such as *Gaussian-*
199 *Flow*[12] or *Gaussian4D*[44]. Instead, we adopt per-frame independent Gaussian splatting as our
200 spatio-temporal representation, denoted as $\{G^t\}$. We use $G^t(v)$ to represent the rendered image

201 obtained from a given camera view v at frame t . Temporal consistency is derived from video composition
202 inversion (see Section 3.4). The term Rev^t represents the rectified frame at time t , which is
203 subsequently passed to the Gaussian Splatting rendered loss. The objective for the stacked temporal
204 representation is modified from Equation 9. To capture the dynamics of the layered background, we
205 employ SVD[3] to obtain $I_{\text{layer}}^t(i, v)$.

$$L_{\text{temporal}} = \min_{\{G^t(i)\}} \sum_{v=1}^m \sum_{t=1}^f \| (G_a^t(i, v) \cdot G^t(i, v) + (1 - G_a^t(v)) \cdot I_b(v)) \cdot \text{Rev}^t(\{\text{Rev}(I_{\text{layer}}^t(i, v), C_{\text{Text}}, r(v) \cdot T)\}, \text{Rev}(I_{\text{layer}}^0(i, v), C_{\text{Text}}, r(v) \cdot T)) \|^2 \quad (12)$$

$$+ \sum_{v=m+1}^n \| (G_a^t(i, v) \cdot G^t(i, v) + (1 - G_a^t(v)) \cdot I_b^t(v)) - I_{\text{layer}}^t(i, v) \|^2. \quad (13)$$

206 4 Experiments

207 4.1 3D harmonization

208 First, following the initial stage settings of DreamGaussian4D [33], we use a reference image to
209 generate a rough 3D Gaussian Splatting using reconstruction loss and distillation gradients from a
210 view-conditioned diffusion model. Artifacts, such as rough Gaussian blobs, are visible in these initial
211 models, as shown in Fig. 1 and the first column of Fig. 5. To align with the DreamGaussian output,
212 we use orbit cameras (see Section A.1) and define the visible range.

213 Then, as explained in Section 3.5.1, we apply the scheduled editing process. The iterative editing
214 results are shown in Fig. 7. Additionally, diverse editing results with the same prompt but different
215 backgrounds are displayed in Fig. 5.

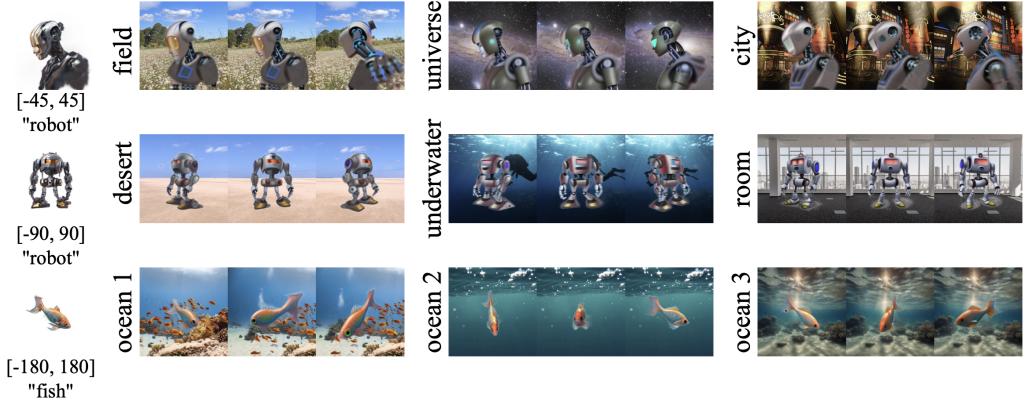


Figure 5: Editing results of input rough GS with three predefined views, harmonized to diverse appearances in different backgrounds.

216 4.2 4D Harmonization

217 Using the same reference image, the reference motion can be prompted from SVD [3]. The video
218 reconstruction loss and view-conditioned distillation score produce a rough 4D model, as shown
219 in Fig. 4(a). This input model is represented as a first frame GS with a deformation network [44],
220 predicting Gaussian blob attribute increments relative to the first frame. This backbone preserves
221 temporal consistency but is challenging to edit and converge. The deformation network accepts a
222 modified first frame GS, including variations in Gaussian blob attributes and the number of Gaussian
223 blobs. However, transferring motion to the edited first frame causes issues, as temporal consistency
224 degrades over time, as shown in Fig. 4(b). In our workflow, we use the edited first frame GS with
225 the deformation network (Fig. 4(a)), and the output (Fig. 4(b)) serves as initialization. Temporal

226 consistency rectification is performed using our method proposed in Section 3.5.2. Results in Fig. 7
 227 show that even after 280 iterations, our method effectively attaches different motions to the edited
 228 character while maintaining appearance consistency. Fig. 6 demonstrates that our method can attach
 229 different first-frame GS to the same motion, serving as a flexible framework for appearance and
 230 motion editing and disentanglement. Emergent temporal semantics are observed after harmonization.
 231 For example, in Fig. 6(b), the flowers wave together, and in Fig. 10(a) in Section A, we see light
 232 reflection, water immersion, and sand resistance.



Figure 6: Input a rough 4D GS model of a flower layered onto two backgrounds (a) and (b), harmonized into two different appearances with the same motion in (c) and (d).

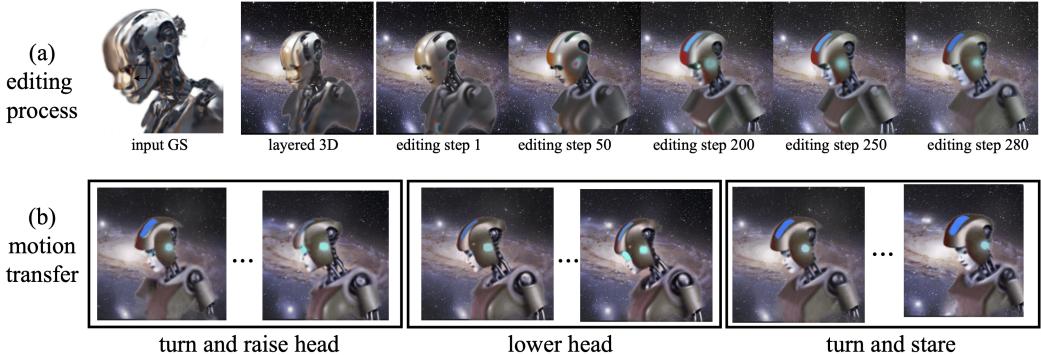


Figure 7: (a) shows input from a rough 3D GS model layered onto a universe background with visual results of the editing process; (b) demonstrates the edited 3D GS attached to different motion styles.

233 4.3 Baseline and Metrics

234 To evaluate the effectiveness of our novel 3D editing method for harmonization, we developed a quan-
 235 titative metric. We used LucidDreamer and DreamGaussian as baselines due to their compatibility
 236 with Gaussian Splatting (GS) for background addition. Baselines involved adding backgrounds to
 237 rendered images using alpha masks from GS, followed by standard distillation scoring. We excluded
 238 methods not specifically designed for GS.

239 As shown in Fig. 8, we present edited images with backgrounds to test harmonization. We calculated
 240 the CLIP similarity [29] between image and text in the format "object + background". We tested

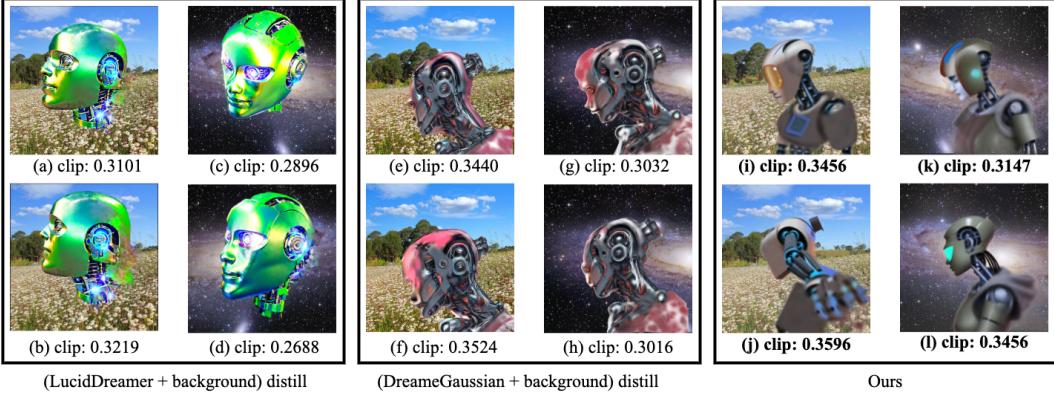


Figure 8: We layer universe and flower field backgrounds during the distillation process of LucidDreamer and DreamGaussian, presenting two randomly sampled views for each method and background. CLIP similarity scores for "robot in universe" and "robot in flower field" are shown below the images.

241 five categories (robot, cat, fish, cat astronaut) with five samples each on five backgrounds (moon,
 242 flower field, universe, sand, ocean), presenting the average CLIP scores in Table 1. And we call this
 243 metric CLIP-Harmonize Our method outperforms LucidDreamer, improving semantic alignment in
 244 backgrounds compared to DreamGaussian. For video quality, we sampled videos from the project
 245 pages of DreamMotion, Control4D, and DreamGaussian4D, which are recent works using distillation
 246 for video editing. Using the metrics CLIP[25], CLIP-Temp[25], and DOVER[45], we present the
 247 average scores of these samples in Table 2. CLIP measures text-frame average alignment, while CLIP
 248 temp assesses temporal semantic consistency by calculating the CLIP similarity between consecutive
 249 frames. Dover is a metric used to evaluate video quality.

Table 1: Average image harmonization CLIP score on 5 categories and 5 backgrounds

Method	LucidDreamer[23]	DreamGaussian[39]	Ours
CLIP-Harmonize↑	0.3270	0.3359	0.3427

Table 2: Average video quality with different metrics.

Method	DreamMotion[19]	Control4D[35]	DreamGaussian4D[33]	Ours
CLIP↑	0.4277	0.4584	0.3574	0.3621
CLIP-Temp↑	0.9900	0.9967	0.9872	0.9899
DOVER↑	0.505	0.5609	0.595	0.6100

250 4.4 Ablations

251 We conduct ablation studies to validate the effectiveness of different components as shown in Fig 9.
 252 Our method allows us to get a better quality and harmonic 4D scene by adding both image and video
 253 inversions.

254 5 Limitations and Conclusions

255 The inversion method in diffusion can perform
 256 harmonization editing for 4D generation, allowing
 257 for 3D harmonization and 4D editing by
 258 decoupling appearances and motions. However,
 259 it is limited to short videos and restricted views
 260 due to base model constraints.



Figure 9: ablations on removing background, without video inversion, image inversion and full model

261 **References**

- 262 [1] S. Bahmani, X. Liu, Y. Wang, I. Skorokhodov, V. Rong, Z. Liu, X. Liu, J. J. Park, S. Tulyakov,
263 G. Wetzstein, et al. Tc4d: Trajectory-conditioned text-to-4d generation. *arXiv preprint arXiv:2403.17920*, 2024.
- 265 [2] S. Bahmani, I. Skorokhodov, V. Rong, G. Wetzstein, L. Guibas, P. Wonka, S. Tulyakov, J. J.
266 Park, A. Tagliasacchi, and D. B. Lindell. 4d-fy: Text-to-4d generation using hybrid score
267 distillation sampling. *arXiv preprint arXiv:2311.17984*, 2023.
- 268 [3] A. Blattmann, T. Dockhorn, S. Kulal, D. Mendelevitch, M. Kilian, D. Lorenz, Y. Levi, Z. English,
269 V. Voletti, A. Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large
270 datasets. *arXiv preprint arXiv:2311.15127*, 2023.
- 271 [4] T. Brooks, A. Holynski, and A. A. Efros. Instructpix2pix: Learning to follow image editing
272 instructions, 2023.
- 273 [5] T. Brooks, B. Peebles, C. Holmes, W. DePue, Y. Guo, L. Jing, D. Schnurr, J. Taylor, T. Luhman,
274 E. Luhman, C. Ng, R. Wang, and A. Ramesh. Video generation models as world simulators.
275 2024.
- 276 [6] A. Cao and J. Johnson. Hexplane: A fast representation for dynamic scenes. In *Proceedings of
277 the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 130–141, 2023.
- 278 [7] R. Chen, Y. Chen, N. Jiao, and K. Jia. Fantasia3d: Disentangling geometry and appearance
279 for high-quality text-to-3d content creation. In *Proceedings of the IEEE/CVF International
280 Conference on Computer Vision*, pages 22246–22256, 2023.
- 281 [8] W.-H. Chu, L. Ke, and K. Fragiadaki. Dreamscene4d: Dynamic multi-object scene generation
282 from monocular videos. *arXiv preprint arXiv:2405.02280*, 2024.
- 283 [9] H. Chung, J. Kim, M. T. Mccann, M. L. Klasky, and J. C. Ye. Diffusion posterior sampling for
284 general noisy inverse problems. *arXiv preprint arXiv:2209.14687*, 2022.
- 285 [10] B. Efron. Tweedie’s formula and selection bias. *Journal of the American Statistical Association*,
286 106(496):1602–1614, 2011.
- 287 [11] D. Epstein, B. Poole, B. Mildenhall, A. A. Efros, and A. Holynski. Disentangled 3d scene
288 generation with layout learning. *arXiv preprint arXiv:2402.16936*, 2024.
- 289 [12] Q. Gao, Q. Xu, Z. Cao, B. Mildenhall, W. Ma, L. Chen, D. Tang, and U. Neumann. Gaussianflow:
290 Splatting gaussian dynamics for 4d content creation. *arXiv preprint arXiv:2403.12365*, 2024.
- 291 [13] A. Haque, M. Tancik, A. Efros, A. Holynski, and A. Kanazawa. Instruct-nerf2nerf: Editing
292 3d scenes with instructions. In *Proceedings of the IEEE/CVF International Conference on
293 Computer Vision*, 2023.
- 294 [14] Y. He, N. Murata, C.-H. Lai, Y. Takida, T. Uesaka, D. Kim, W.-H. Liao, Y. Mitsufuji, J. Z. Kolter,
295 R. Salakhutdinov, et al. Manifold preserving guided diffusion. *arXiv preprint arXiv:2311.16424*,
296 2023.
- 297 [15] R. Henschel, L. Khachatryan, D. Hayrapetyan, H. Poghosyan, V. Tadevosyan, Z. Wang,
298 S. Navasardyan, and H. Shi. Streamingt2v: Consistent, dynamic, and extendable long video
299 generation from text, 2024.
- 300 [16] J. Ho and T. Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*,
301 2022.
- 302 [17] A. Jain, B. Mildenhall, J. T. Barron, P. Abbeel, and B. Poole. Zero-shot text-guided object
303 generation with dream fields. In *Proceedings of the IEEE/CVF conference on computer vision
304 and pattern recognition*, pages 867–876, 2022.
- 305 [18] T. Jakab, R. Li, S. Wu, C. Rupprecht, and A. Vedaldi. Farm3d: Learning articulated 3d animals
306 by distilling 2d diffusion. *arXiv preprint arXiv:2304.10535*, 2023.

- 307 [19] H. Jeong, J. Chang, G. Y. Park, and J. C. Ye. Dreammotion: Space-time self-similarity score
 308 distillation for zero-shot video editing. *arXiv preprint arXiv:2403.12002*, 2024.
- 309 [20] Y. Jiang, L. Zhang, J. Gao, W. Hu, and Y. Yao. Consistent4d: Consistent 360 $\{\backslash \deg\}$ dynamic
 310 object generation from monocular video. *arXiv preprint arXiv:2311.02848*, 2023.
- 311 [21] T. Karras, M. Aittala, T. Aila, and S. Laine. Elucidating the design space of diffusion-based
 312 generative models. *Advances in Neural Information Processing Systems*, 35:26565–26577,
 313 2022.
- 314 [22] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis. 3d gaussian splatting for real-time
 315 radiance field rendering. *ACM Transactions on Graphics*, 42(4):1–14, 2023.
- 316 [23] Y. Liang, X. Yang, J. Lin, H. Li, X. Xu, and Y. Chen. Luciddreamer: Towards high-fidelity
 317 text-to-3d generation via interval score matching. *arXiv preprint arXiv:2311.11284*, 2023.
- 318 [24] H. Ling, S. W. Kim, A. Torralba, S. Fidler, and K. Kreis. Align your gaussians: Text-to-4d with
 319 dynamic 3d gaussians and composed diffusion models. *arXiv preprint arXiv:2312.13763*, 2023.
- 320 [25] Y. Liu, X. Cun, X. Liu, X. Wang, Y. Zhang, H. Chen, Y. Liu, T. Zeng, R. Chan, and Y. Shan.
 321 Evalcrafter: Benchmarking and evaluating large video generation models. *arXiv preprint
 322 arXiv:2310.11440*, 2023.
- 323 [26] M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black. Smpl: A skinned multi-
 324 person linear model. In *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*, pages
 325 851–866. 2023.
- 326 [27] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. Nerf:
 327 Representing scenes as neural radiance fields for view synthesis, 2020.
- 328 [28] B. Poole, A. Jain, J. T. Barron, and B. Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion,
 329 2022.
- 330 [29] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,
 331 P. Mishkin, J. Clark, G. Krueger, and I. Sutskever. Learning transferable visual models from
 332 natural language supervision, 2021.
- 333 [30] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen. Hierarchical text-conditional image
 334 generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.
- 335 [31] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever.
 336 Zero-shot text-to-image generation, 2021.
- 337 [32] D. Rempe, Z. Luo, X. Bin Peng, Y. Yuan, K. Kitani, K. Kreis, S. Fidler, and O. Litany. Trace
 338 and pace: Controllable pedestrian animation via guided trajectory diffusion. In *Proceedings of
 339 the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13756–13766,
 340 2023.
- 341 [33] J. Ren, L. Pan, J. Tang, C. Zhang, A. Cao, G. Zeng, and Z. Liu. Dreamgaussian4d: Generative
 342 4d gaussian splatting. *arXiv preprint arXiv:2312.17142*, 2023.
- 343 [34] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis
 344 with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision
 345 and pattern recognition*, pages 10684–10695, 2022.
- 346 [35] R. Shao, J. Sun, C. Peng, Z. Zheng, B. Zhou, H. Zhang, and Y. Liu. Control4d: Dynamic portrait
 347 editing by learning 4d gan from 2d diffusion-based editor. *arXiv preprint arXiv:2305.20082*,
 348 2023.
- 349 [36] U. Singer, S. Sheynin, A. Polyak, O. Ashual, I. Makarov, F. Kokkinos, N. Goyal, A. Vedaldi,
 350 D. Parikh, J. Johnson, et al. Text-to-4d dynamic scene generation. *arXiv preprint
 351 arXiv:2301.11280*, 2023.
- 352 [37] J. Song, C. Meng, and S. Ermon. Denoising diffusion implicit models. *arXiv preprint
 353 arXiv:2010.02502*, 2020.

- 354 [38] Y. Song and S. Ermon. Generative modeling by estimating gradients of the data distribution.
 355 *Advances in neural information processing systems*, 32, 2019.
- 356 [39] J. Tang, J. Ren, H. Zhou, Z. Liu, and G. Zeng. Dreamgaussian: Generative gaussian splatting
 357 for efficient 3d content creation. *arXiv preprint arXiv:2309.16653*, 2023.
- 358 [40] J. Tang, J. Ren, H. Zhou, Z. Liu, and G. Zeng. Dreamgaussian: Generative gaussian splatting
 359 for efficient 3d content creation, 2024.
- 360 [41] P. Wang, Z. Fan, D. Xu, D. Wang, S. Mohan, F. Iandola, R. Ranjan, Y. Li, Q. Liu, Z. Wang, and
 361 V. Chandra. Steindreamer: Variance reduction for text-to-3d score distillation via stein identity,
 362 2024.
- 363 [42] P. Wang, D. Xu, Z. Fan, D. Wang, S. Mohan, F. Iandola, R. Ranjan, Y. Li, Q. Liu, Z. Wang,
 364 et al. Taming mode collapse in score distillation for text-to-3d generation. *arXiv preprint*
 365 *arXiv:2401.00909*, 2023.
- 366 [43] Z. Wang, C. Lu, Y. Wang, F. Bao, C. Li, H. Su, and J. Zhu. Prolificdreamer: High-fidelity and
 367 diverse text-to-3d generation with variational score distillation, 2023.
- 368 [44] G. Wu, T. Yi, J. Fang, L. Xie, X. Zhang, W. Wei, W. Liu, Q. Tian, and X. Wang. 4d gaussian
 369 splatting for real-time dynamic scene rendering. *arXiv preprint arXiv:2310.08528*, 2023.
- 370 [45] H. Wu, E. Zhang, L. Liao, C. Chen, J. Hou, A. Wang, W. Sun, Q. Yan, and W. Lin. Exploring
 371 video quality assessment on user generated contents from aesthetic and technical perspectives.
 372 In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 20144–
 373 20154, 2023.
- 374 [46] D. Xu, H. Liang, N. P. Bhatt, H. Hu, H. Liang, K. N. Plataniotis, and Z. Wang. Comp4d:
 375 Llm-guided compositional 4d scene generation. *arXiv preprint arXiv:2403.16993*, 2024.
- 376 [47] Q. Yang, M. Feng, Z. Wu, S. Sun, W. Dong, Y. Wang, and A. Mian. Beyond skeletons:
 377 Integrative latent mapping for coherent 4d sequence generation, 2024.
- 378 [48] T. Yi, J. Fang, J. Wang, G. Wu, L. Xie, X. Zhang, W. Liu, Q. Tian, and X. Wang. Gaussian-
 379 dreamer: Fast generation from text to 3d gaussians by bridging 2d and 3d diffusion models,
 380 2024.
- 381 [49] Y. Yin, D. Xu, Z. Wang, Y. Zhao, and Y. Wei. 4dgen: Grounded 4d content generation with
 382 spatial-temporal consistency. *arXiv preprint arXiv:2312.17225*, 2023.
- 383 [50] Y. Yuan, V. Makoviychuk, Y. Guo, S. Fidler, X. Peng, and K. Fatahalian. Learning physically
 384 simulated tennis skills from broadcast videos. *ACM Trans. Graph*, 42(4), 2023.
- 385 [51] Y. Zhao, Z. Yan, E. Xie, L. Hong, Z. Li, and G. H. Lee. Animate124: Animating one image to
 386 4d dynamic scene. *arXiv preprint arXiv:2311.14603*, 2023.
- 387 [52] J. Zhu, P. Zhuang, and S. Koyejo. Hifa: High-fidelity text-to-3d generation with advanced
 388 diffusion guidance. In *The Twelfth International Conference on Learning Representations*,
 389 2023.

390 **A Appendix / supplemental material**

391 **A.1 Experimental Settings**

392 We run the HarmonyDreamer pipeline on a single NVIDIA RTX A6000 GPU with 48G VRAM under
393 Python 3.8 and CUDA 11.8 environment. Each harmonization editing step takes around 1 minute.

394 In our settings, we utilize an orbit camera with predefined and fixed intrinsic parameters. The extrinsic
395 parameters of the camera can be parameterized by the azimuth angle ver, the horizon angle hor, and
396 the radius r , thus $V = V(\text{ver}, \text{hor}, r)$. Additionally, the rendered image inherently includes the alpha
397 mask of the assembly of Gaussian splats, which we denote as $G_{\mathbf{a}}^t(V)$.

398 **A.2 More Qualitative Results**

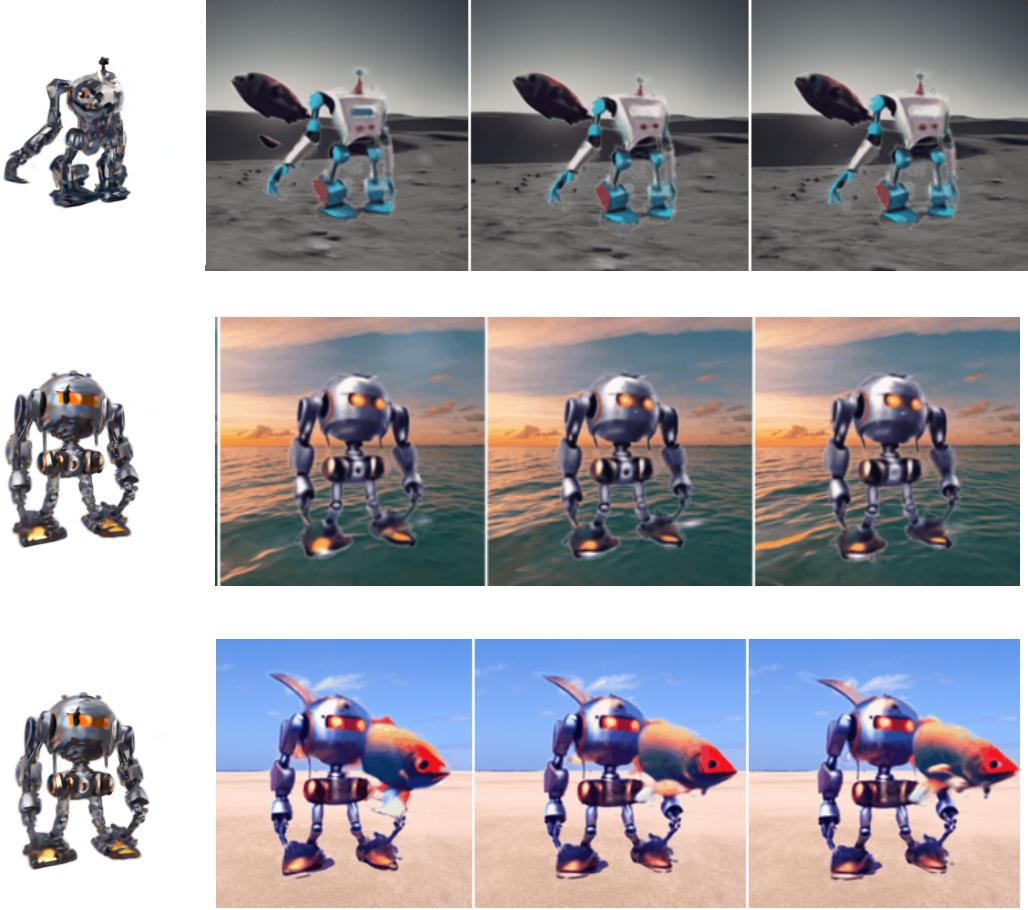


Figure 10: Emergent temporal semantics are observed after harmonization. For example, we see light reflection in the first row, water immersion in the second row, and sand resistance in the third row.

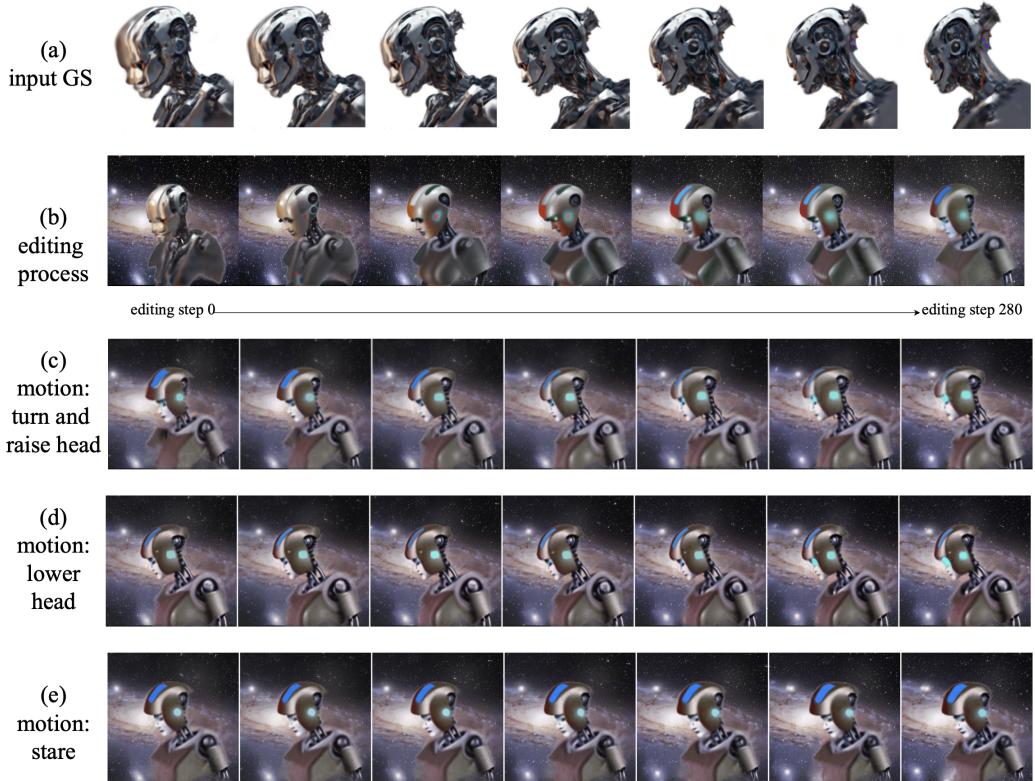


Figure 11: Expanding of Fig. 7(a) shows input from a rough 3D GS model layered onto a universe background with visual results of the editing process; (b) demonstrates the edited 3D GS attached to different motion styles.

399 **NeurIPS Paper Checklist**

400 **1. Claims**

401 Question: Do the main claims made in the abstract and introduction accurately reflect the
402 paper's contributions and scope?

403 Answer: [Yes]

404 Justification: Our main claims are made based on the real experimental results.

405 Guidelines:

- 406 • The answer NA means that the abstract and introduction do not include the claims
407 made in the paper.
- 408 • The abstract and/or introduction should clearly state the claims made, including the
409 contributions made in the paper and important assumptions and limitations. A No or
410 NA answer to this question will not be perceived well by the reviewers.
- 411 • The claims made should match theoretical and experimental results, and reflect how
412 much the results can be expected to generalize to other settings.
- 413 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
414 are not attained by the paper.

415 **2. Limitations**

416 Question: Does the paper discuss the limitations of the work performed by the authors?

417 Answer: [Yes]

418 Justification: We discussed the limitations of our method at the end of the main part.

419 Guidelines:

- 420 • The answer NA means that the paper has no limitation while the answer No means that
421 the paper has limitations, but those are not discussed in the paper.
- 422 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 423 • The paper should point out any strong assumptions and how robust the results are to
424 violations of these assumptions (e.g., independence assumptions, noiseless settings,
425 model well-specification, asymptotic approximations only holding locally). The authors
426 should reflect on how these assumptions might be violated in practice and what the
427 implications would be.
- 428 • The authors should reflect on the scope of the claims made, e.g., if the approach was
429 only tested on a few datasets or with a few runs. In general, empirical results often
430 depend on implicit assumptions, which should be articulated.
- 431 • The authors should reflect on the factors that influence the performance of the approach.
432 For example, a facial recognition algorithm may perform poorly when image resolution
433 is low or images are taken in low lighting. Or a speech-to-text system might not be
434 used reliably to provide closed captions for online lectures because it fails to handle
435 technical jargon.
- 436 • The authors should discuss the computational efficiency of the proposed algorithms
437 and how they scale with dataset size.
- 438 • If applicable, the authors should discuss possible limitations of their approach to
439 address problems of privacy and fairness.
- 440 • While the authors might fear that complete honesty about limitations might be used by
441 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
442 limitations that aren't acknowledged in the paper. The authors should use their best
443 judgment and recognize that individual actions in favor of transparency play an impor-
444 tant role in developing norms that preserve the integrity of the community. Reviewers
445 will be specifically instructed to not penalize honesty concerning limitations.

446 **3. Theory Assumptions and Proofs**

447 Question: For each theoretical result, does the paper provide the full set of assumptions and
448 a complete (and correct) proof?

449 Answer: [NA]

450 Justification: Our works does not include theoretical results.

451 Guidelines:

- 452 • The answer NA means that the paper does not include theoretical results.
- 453 • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- 455 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 456 • The proofs can either appear in the main paper or the supplemental material, but if
- 457 they appear in the supplemental material, the authors are encouraged to provide a short
- 458 proof sketch to provide intuition.
- 459 • Inversely, any informal proof provided in the core of the paper should be complemented
- 460 by formal proofs provided in appendix or supplemental material.
- 461 • Theorems and Lemmas that the proof relies upon should be properly referenced.

462 **4. Experimental Result Reproducibility**

463 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
464 perimental results of the paper to the extent that it affects the main claims and/or conclusions
465 of the paper (regardless of whether the code and data are provided or not)?

466 Answer: [Yes]

467 Justification: We described our experimental results in the Experiments section.

468 Guidelines:

- 469 • The answer NA means that the paper does not include experiments.
- 470 • If the paper includes experiments, a No answer to this question will not be perceived
471 well by the reviewers: Making the paper reproducible is important, regardless of
472 whether the code and data are provided or not.
- 473 • If the contribution is a dataset and/or model, the authors should describe the steps taken
474 to make their results reproducible or verifiable.
- 475 • Depending on the contribution, reproducibility can be accomplished in various ways.
476 For example, if the contribution is a novel architecture, describing the architecture fully
477 might suffice, or if the contribution is a specific model and empirical evaluation, it may
478 be necessary to either make it possible for others to replicate the model with the same
479 dataset, or provide access to the model. In general, releasing code and data is often
480 one good way to accomplish this, but reproducibility can also be provided via detailed
481 instructions for how to replicate the results, access to a hosted model (e.g., in the case
482 of a large language model), releasing of a model checkpoint, or other means that are
483 appropriate to the research performed.
- 484 • While NeurIPS does not require releasing code, the conference does require all submis-
485 sions to provide some reasonable avenue for reproducibility, which may depend on the
486 nature of the contribution. For example
 - 487 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
488 to reproduce that algorithm.
 - 489 (b) If the contribution is primarily a new model architecture, the paper should describe
490 the architecture clearly and fully.
 - 491 (c) If the contribution is a new model (e.g., a large language model), then there should
492 either be a way to access this model for reproducing the results or a way to reproduce
493 the model (e.g., with an open-source dataset or instructions for how to construct
494 the dataset).
 - 495 (d) We recognize that reproducibility may be tricky in some cases, in which case
496 authors are welcome to describe the particular way they provide for reproducibility.
497 In the case of closed-source models, it may be that access to the model is limited in
498 some way (e.g., to registered users), but it should be possible for other researchers
499 to have some path to reproducing or verifying the results.

500 **5. Open access to data and code**

501 Question: Does the paper provide open access to the data and code, with sufficient instruc-
502 tions to faithfully reproduce the main experimental results, as described in supplemental
503 material?

504 Answer: [Yes]

505 Justification: We include our code in the supplementary material submission.

506 Guidelines:

- 507 • The answer NA means that paper does not include experiments requiring code.
- 508 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 509 • While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- 510 • The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 511 • The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 512 • The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- 513 • At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- 514 • Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

515 **6. Experimental Setting/Details**

516 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

517 Answer: [Yes]

518 Justification: We described our experimental settings and details in the supplementary materials.

519 Guidelines:

- 520 • The answer NA means that the paper does not include experiments.
- 521 • The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- 522 • The full details can be provided either with the code, in appendix, or as supplemental material.

523 **7. Experiment Statistical Significance**

524 Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

525 Answer: [No]

526 Justification: Our quantitative results are not accompanied by error bars, confidence intervals, or statistical significance tests.

527 Guidelines:

- 528 • The answer NA means that the paper does not include experiments.
- 529 • The authors should answer “Yes” if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- 530 • The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- 531 • The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- 532 • The assumptions made should be given (e.g., Normally distributed errors).

- 556 • It should be clear whether the error bar is the standard deviation or the standard error
 557 of the mean.
 558 • It is OK to report 1-sigma error bars, but one should state it. The authors should
 559 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
 560 of Normality of errors is not verified.
 561 • For asymmetric distributions, the authors should be careful not to show in tables or
 562 figures symmetric error bars that would yield results that are out of range (e.g. negative
 563 error rates).
 564 • If error bars are reported in tables or plots, The authors should explain in the text how
 565 they were calculated and reference the corresponding figures or tables in the text.

566 **8. Experiments Compute Resources**

567 Question: For each experiment, does the paper provide sufficient information on the com-
 568 puter resources (type of compute workers, memory, time of execution) needed to reproduce
 569 the experiments?

570 Answer: [Yes]

571 Justification: We described our experimental settings and required resources in the supple-
 572 mentary materials.

573 Guidelines:

- 574 • The answer NA means that the paper does not include experiments.
- 575 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
 576 or cloud provider, including relevant memory and storage.
- 577 • The paper should provide the amount of compute required for each of the individual
 578 experimental runs as well as estimate the total compute.
- 579 • The paper should disclose whether the full research project required more compute
 580 than the experiments reported in the paper (e.g., preliminary or failed experiments that
 581 didn't make it into the paper).

582 **9. Code Of Ethics**

583 Question: Does the research conducted in the paper conform, in every respect, with the
 584 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

585 Answer: [Yes]

586 Justification: Our work does not violate respects listed in the NeurIPS Code of Ethics.

587 Guidelines:

- 588 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 589 • If the authors answer No, they should explain the special circumstances that require a
 590 deviation from the Code of Ethics.
- 591 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
 592 eration due to laws or regulations in their jurisdiction).

593 **10. Broader Impacts**

594 Question: Does the paper discuss both potential positive societal impacts and negative
 595 societal impacts of the work performed?

596 Answer: [No]

597 Justification: Our work does not use photorealistic contents, which limits its potential for
 598 societal impact, as it avoids issues like disinformation, privacy violations, and generating
 599 fake profiles.

600 Guidelines:

- 601 • The answer NA means that there is no societal impact of the work performed.
- 602 • If the authors answer NA or No, they should explain why their work has no societal
 603 impact or why the paper does not address societal impact.
- 604 • Examples of negative societal impacts include potential malicious or unintended uses
 605 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
 606 (e.g., deployment of technologies that could make decisions that unfairly impact specific
 607 groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: Our work does not pose risk of unsafe misuse.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We properly cited the source of the generated images/videos we used to conduct our work.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

- 661 • If this information is not available online, the authors are encouraged to reach out to
662 the asset's creators.

663 **13. New Assets**

664 Question: Are new assets introduced in the paper well documented and is the documentation
665 provided alongside the assets?

666 Answer: [Yes]

667 Justification: We include a anonymized URL to document all released new assets from our
668 work.

669 Guidelines:

- 670 • The answer NA means that the paper does not release new assets.
671 • Researchers should communicate the details of the dataset/code/model as part of their
672 submissions via structured templates. This includes details about training, license,
673 limitations, etc.
674 • The paper should discuss whether and how consent was obtained from people whose
675 asset is used.
676 • At submission time, remember to anonymize your assets (if applicable). You can either
677 create an anonymized URL or include an anonymized zip file.

678 **14. Crowdsourcing and Research with Human Subjects**

679 Question: For crowdsourcing experiments and research with human subjects, does the paper
680 include the full text of instructions given to participants and screenshots, if applicable, as
681 well as details about compensation (if any)?

682 Answer: [NA]

683 Justification: Our work does not involve crowdsourcing nor research with human subjects.

684 Guidelines:

- 685 • The answer NA means that the paper does not involve crowdsourcing nor research with
686 human subjects.
687 • Including this information in the supplemental material is fine, but if the main contribu-
688 tion of the paper involves human subjects, then as much detail as possible should be
689 included in the main paper.
690 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
691 or other labor should be paid at least the minimum wage in the country of the data
692 collector.

693 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
694 Subjects**

695 Question: Does the paper describe potential risks incurred by study participants, whether
696 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
697 approvals (or an equivalent approval/review based on the requirements of your country or
698 institution) were obtained?

699 Answer: [NA]

700 Justification: Our work does not involve crowdsourcing nor research with human subjects.

701 Guidelines:

- 702 • The answer NA means that the paper does not involve crowdsourcing nor research with
703 human subjects.
704 • Depending on the country in which research is conducted, IRB approval (or equivalent)
705 may be required for any human subjects research. If you obtained IRB approval, you
706 should clearly state this in the paper.
707 • We recognize that the procedures for this may vary significantly between institutions
708 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
709 guidelines for their institution.
710 • For initial submissions, do not include any information that would break anonymity (if
711 applicable), such as the institution conducting the review.