

Supplementary Material for

Catalyst4D: High-Fidelity 3D-to-4D Scene Editing via Dynamic Propagation

001 A. Video Demo

002 We provide demo videos on our project page showing the
003 dynamic editing capabilities of our method. The videos
004 include comparative results against baseline approaches on
005 novel test viewpoints, as well as demonstrations of image-
006 guided global style transfer and text-guided local edits en-
007 abled by the proposed Catalyst4D. We encourage readers to
008 view these videos, as they offer a clear illustration of the
009 high visual fidelity and editing consistency achieved by our
010 approach.

011 B. Additional Comparison Results

012 B.1. Additional Quantitative Comparison

013 Standard CLIP and VBench scores may fail to capture fine-
014 grained 3D-to-4D propagation due to their global nature,
015 we further adopt **EditScore** (instruction-following accu-
016 racy and visual fidelity) [6] and **VE-Bench** (local tempo-
017 ral consistency) [8], better aligned with human perception.
018 Tab. B.1 on 16 prompts(multi-camera) from the main paper
019 shows Catalyst4D consistently surpasses prior methods.

020 B.2. Additional Quantitative Ablation

021 We clarify that the main paper ablation (Tab. 2) is cumula-
022 tive: w/o AMG is the “w/o both” baseline, and w/o CUAR
023 retains only AMG. Tab. B.2, evaluated with new metrics,
024 further confirms both AMG and CUAR are essential, out-
025 performing DeformNet-Guide (w/o both).

026 B.3. Comparison on Background Preservation

027 Although the compared CTRL-D [3] appears visually plau-
028 sible, it still introduces undesired modifications in non-
029 edited regions due to the limitations of the underlying 2D
030 diffusion model. As illustrated in Fig.B.1, while CTRL-
031 D edits the character, it also causes noticeable deviations
032 in non-target objects such as the dog on the stool and ob-
033 jects on the table compared to the original scene. By con-
034 trast, our method constrains target dynamic Gaussians di-
035 rectly through 3D editing gradients, enabling more precise
036 and localized editing in dynamic scenes.

037 B.4. Comparison on the Monocular Dataset

038 We further evaluate text-driven 4D editing on the monocu-
039 lar HyperNeRF [7] dataset. Since Instruct 4D-to-4D does
040 not support this dataset, we primarily compare our method

Table B.1. Quantitative comparison using EditScore and VE-Bench.

Metric	IN4D	I4DGS	CTRL-D	Ours
EditScore↑	4.034	5.618	4.326	7.375
VE-Bench↑	0.155	0.256	0.163	1.080

Table B.2. Ablation studies on AMG and CUAR modules.

Metric	w/o both	w/ AMG	w/ CUAR	w/ both
EditScore↑	2.713	5.112	5.005	7.375
VE-Bench↑	0.059	0.452	0.580	1.080

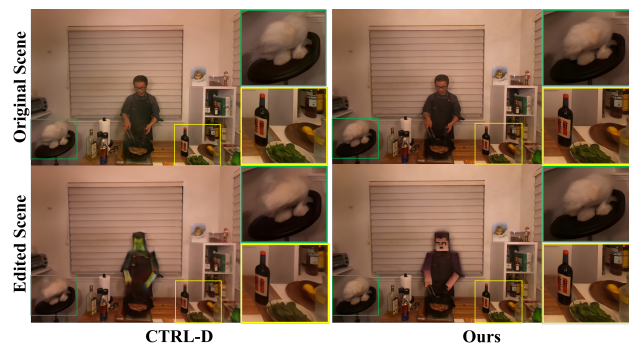


Figure B.1. Qualitative comparison of localized editing. In contrast to CTRL-D, which introduces inconsistencies in non-edited regions, our method achieves more precise and localized editing by constraining dynamic Gaussians via 3D editing gradients.

041 against Instruct-4DGS and CTRL-D. The qualitative com-
042 parisons are shown in Fig. B.2, where the top and bottom
043 rows of each example correspond to different timesteps.
044 Instruct-4DGS struggles to localize the target object, result-
045 ing in edits being incorrectly applied to irrelevant regions
046 of the scene. CTRL-D, which reconstructs the 4D represen-
047 tation from 2D diffusion outputs, suffers from the inherent
048 modality gap, producing blurry results and noticeable tempo-
049 ral inconsistency. As illustrated in the “Glacial extruder”
050 example in the lower half of Fig. B.2, the appearance of
051 the edited extruder varies substantially across timesteps. In
052 contrast, our method directly propagates a high-fidelity 3D
053 edit from the first frame to all subsequent frames, yielding
054 clearer textures, stronger temporal coherence, and signifi-
055 cantly improved overall 4D editing quality.

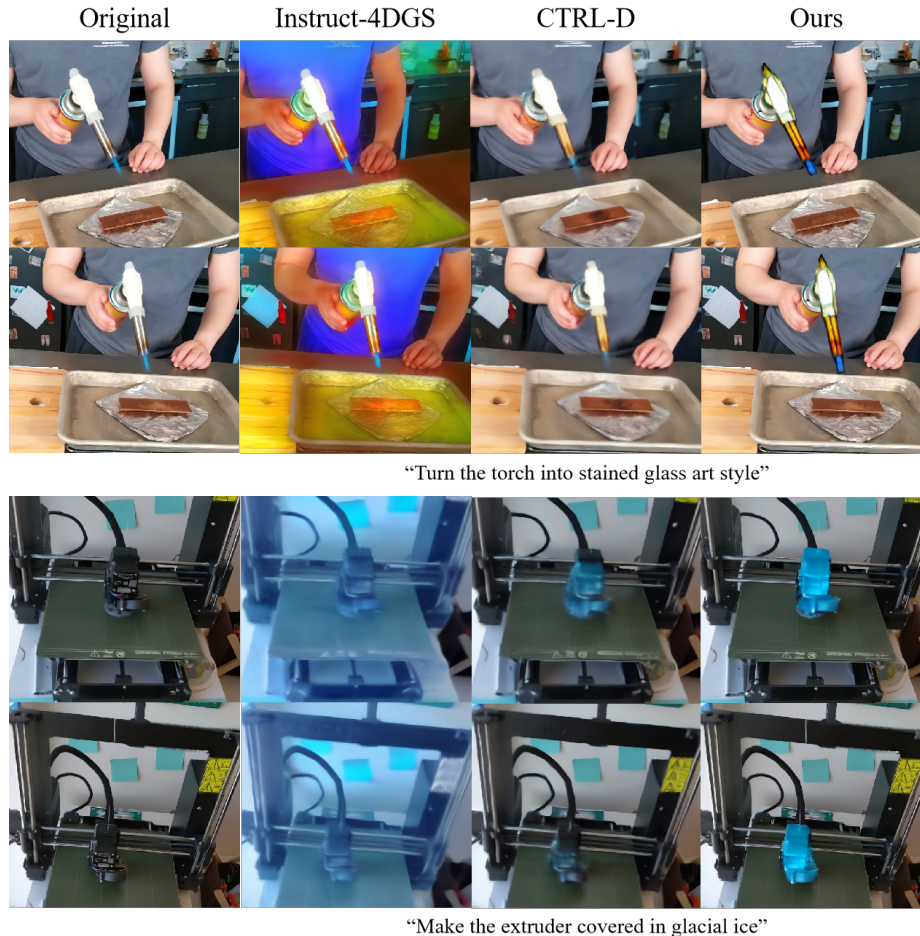


Figure B.2. Comparison on HyperNeRF Dataset. The results show that Instruct-4DGS fails to correctly localize edits and CTRL-D produces blurry and temporally inconsistent results, whereas our method demonstrates significant advantages in both visual quality and temporal consistency.

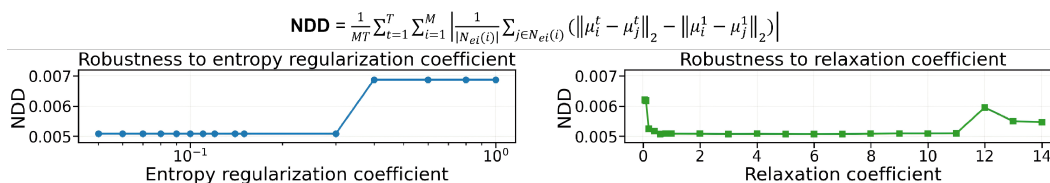


Figure B.3. Effect of Sinkhorn regularization on motion transfer.

056 B.5. Robustness of Sinkhorn Algorithm

057 Our outlier-filtered anchors yield a well-conditioned cost
 058 matrix, enabling stable Sinkhorn convergence. To assess
 059 robustness under extreme topology changes (e.g. *Sear-*
 060 *steak to Minecraft*), we measure Neighborhood Distance
 061 Deviation (NDD), which captures the temporal preservation
 062 of local anchor geometry and directly reflects corres-
 063 pondence stability. Fig. B.3 confirms NDD stability across
 064 wide Sinkhorn ranges ($\lambda_0 \in [0.05, 0.3]$, $\lambda_{1,2} \in [0.6, 11]$).

B.6. Comparison of Global Style Transfer

We further compare our method with CLIPGaussian [4] and
 CTRL-D on the task of image-guided global style transfer.
 To establish a controlled comparison with CTRL-D, which
 requires an image prompt, we first stylize the first frame 3D
 scene using SGSST [2], our 3D baseline. A rendered image
 from this stylized 3D scene is then provided to CTRL-D as
 its style reference. Consequently, our method and CTRL-D
 begin with nearly identical style information but leverage

065

066

067

068

069

070

071

072

073



Figure B.4. Comparative results for global style transfer. Unlike the over-smoothed results of CTRL-D and the chaotic textures from CLIP-Gaussian, our method produces finer visual textures, better preserves the original scene’s geometry, and demonstrates a color distribution more consistent with the reference style image.

074 distinct technical approaches. Our method propagates the
 075 style from the 3D scene across the sequence (a 3D-to-4D
 076 approach), while CTRL-D uses its DreamBooth-based fine-
 077 tuning to reconstruct the 4D effect from the 2D style image
 078 (a 2D-to-4D approach). This allows for a direct assessment
 079 of the two pathways, minimizing the influence of external
 080 stylization models.

081 As shown in Fig. B.4, the global style transfer results
 082 produced by CTRL-D appear superficially similar to ours,
 083 an expected outcome given that CTRL-D’s style reference
 084 is a rendered image from our 3D stylization baseline. How-
 085 ever, our method achieves noticeably finer textures and
 086 higher visual clarity. CTRL-D, which relies on a 2D dif-
 087 fusion model, often generates results with blurriness and
 088 over-smoothing due to the inherent 2D-to-4D reconstruc-
 089 tion gap. Meanwhile, CLIPGaussian tends to introduce
 090 chaotic textures, as reflected in the distorted facial region of
 091 the character in the example shown in the fourth column of
 092 Fig. B.4. In contrast, our approach preserves the underlying
 093 4D geometry more effectively by propagating SGSST [2]-
 094 based stylization (our 3D baseline) consistently across the
 095 sequence. Additionally, our results exhibit a color distribu-
 096 tion that more faithfully matches the reference style image.

097 C. Implementation Details

098 C.1. Optimal Transport-based Anchor Matching

099 Given the anchor set $A_{\text{src}} = \{\mathbf{p}_i^{\text{src}}\}_{i=1}^n$ of the original dy-
 100 namic Gaussians \mathcal{G}^1 in the first frame and the anchor set
 101 $A_{\text{edit}} = \{\mathbf{p}_j^{\text{edit}}\}_{j=1}^m$ of the edited Gaussians $\mathcal{G}_{\text{edit}}^1$, we for-

102 mulate the anchor correspondence as an unbalanced optimal
 103 transport (UOT) problem and solve it using the Sinkhorn al-
 104 gorithm [1].

105 We define a distance matrix D using a Welsch robust
 106 penalty:

$$D_{ij} = 1 - \exp\left(-\frac{\|\mathbf{p}_i^{\text{src}} - \mathbf{p}_j^{\text{edit}}\|_2^2}{2\beta^2}\right), \quad (1) \quad 107$$

108 where $\beta = \gamma \cdot d_{\text{med}}$, d_{med} is the median of all pairwise
 109 anchor distances, and we empirically set $\gamma = 0.05$.

110 The UOT objective is written as:

$$\begin{aligned} & \min_{P \in \mathbb{R}_+^{n \times m}} \sum_{i=1}^n \sum_{j=1}^m D_{ij} P_{ij} - \lambda_0 \sum_{i,j} P_{ij} \log P_{ij} & 112 \\ & + \lambda_1 \text{KL}\left(\mathbf{P}\mathbf{1}_m \parallel \frac{1}{n}\mathbf{1}_n\right) + \lambda_2 \text{KL}\left(\mathbf{P}^T\mathbf{1}_n \parallel \frac{1}{m}\mathbf{1}_m\right), & 113 \end{aligned} \quad (2)$$

114 s.t. $P_{ij} \geq 0, \forall i \in \{1, \dots, n\}, j \in \{1, \dots, m\}$, where \mathbf{P} is
 115 the transport plan matrix, $\lambda_0 = 0.1$, $\lambda_1 = 1.0$, $\lambda_2 = 1.0$
 116 are regularization weights, and KL denotes the Kullback-
 117 Leibler divergence.

118 From the optimal plan \mathbf{P}^* , we obtain the correspondence
 119 corr by assigning each edited anchor to the source anchor
 120 with the highest transport mass:

$$\text{corr}_j = \arg \max_{i \in \{1, \dots, n\}} P_{ij}^*. \quad (3) \quad 121$$

122 C.2. Deformation of Edited Dynamic Gaussians

123 Notably, when computing the deformation of edited Gaussians, unlike the additive form used for positional deformation $\Delta\mu$ and rotational deformation $\Delta\mathbf{q}$, scaling deformation $\Delta\mathbf{s}$ is defined multiplicatively to preserve proportional change:

$$\begin{aligned}
 \mu^t &= \mu^1 + \Delta\mu^t, \\
 \mathbf{q}^t &= \mathbf{q}^1 + \Delta\mathbf{q}^t, \\
 \mathbf{s}^t &= \mathbf{s}^1 \cdot \Delta\mathbf{s}^t.
 \end{aligned}
 \tag{4}$$

129 where $\Delta\mu^t$, $\Delta\mathbf{q}^t$, $\Delta\mathbf{s}^t$ are the deformation quantities of
130 edited Gaussians from frame 1 to frame t .

131 C.3. Additional Experimental Details

132 **Anchor Construction.** Since $\mathcal{G}_{\text{edit}}^1$ and \mathcal{G}^1 are largely spa-
133 tially aligned, we compute a shared bounding sphere and
134 sample an identical set of rays to construct anchor points for
135 both Gaussian clouds. Specifically, we sample 300,000 rays
136 from the minimum bounding sphere and use $k = 2$ near-
137 est neighbors when forming local anchor neighborhoods.
138 For Gaussians not covered by any anchor region, we assign
139 them to the nearest anchor using Euclidean distance.

140 **Loss Design.** For foreground edits, we use standard
141 L1+SSIM to propagate the effect faithfully. Background
142 supervision relies on L1 only, as CUAR freezes non-edited
143 Gaussians and L1 sufficiently prevents artifacts; adding
144 SSIM yields negligible gains while increasing cost.

145 **Appearance Refinement Parameters.** Eq. 17 exploits un-
146 certainty contrast between artifact-prone and stable regions.
147 The threshold ϵ trades quality for efficiency (higher for
148 faster refinement, lower for broader coverage) and remains
149 stable over a reasonable range. To flexibly control the spa-
150 tial coverage of the artifact mask, we set the threshold pa-
151 rameter ϵ within the range $[0.5, 2.0]$. During refinement,
152 the loss weights are fixed as $\eta = 0.2$ and $\zeta = 0.3$, applied
153 to the foreground and background objectives described in
154 Sec. 4.2. These settings balance artifact correction with
155 preservation of non-corrupted regions.

156 **Evaluation Protocol and Implementation.** For quantita-
157 tive evaluation, we render novel test viewpoints and com-
158 pute the CLIP similarity by averaging scores over 30 uni-
159 formly sampled frames from each test video. Temporal con-
160 sistency is evaluated on the full rendered sequence to mea-
161 sure stability across viewpoints and timesteps.

162 For dynamic scene reconstruction with Swift4D, we
163 adopt effective rank regularization [5] to suppress needle-
164 like artifacts. Following the settings of 4DGS [9], we
165 constrain the edited dynamic Gaussians to share consistent
166 opacity and color attributes across all frames during opti-
167 mization.

The training breakdown is as follows: Anchor con-
struction (<30s), Sinkhorn solver (~ 15 s), Motion Guid-
ance (~ 1 min), and CUAR (25–35 min). All experiments
are implemented in Python 3.7.12 using PyTorch 1.13.1 on
Ubuntu 22.04. Our method is trained on a single NVIDIA
A100 Tensor Core GPU.

174 C.4. Editing Prompts and 3D Baseline Configura- 175 tions

176 Here, we provide the detailed editing configurations used in
177 our experiments, including the text prompts and the corre-
178 sponding 3D editing baselines adopted for each edit.

180 For the results shown in Fig. 1, the editing instructions
181 and 3D baseline methods are:

- “Turn the torch into Pop art style” with DreamCatalyst, 182
- “Make the torch carved from a flawless emerald” with
DGE, 183
- “Turn him into a Claymation character” with DGE, 184
- “Make the extruder covered in glacial ice” with Dream-
Catalyst, 185
- “Make the extruder look like toast” with DGE, 186
- “Turn him into JoJo’s Bizarre Adventure anime style”
with DGE. 187

188 For the edits in Fig. 4, we use:

- “Turn his clothes into a football player outfit” with
DreamCatalyst, 189
- “Turn his hat into a newsboy cap” with DGE, 190
- “Turn him into a Minecraft character” with DGE, 191
- “Make them look like Marble roman sculptures” with
DreamCatalyst, 192
- “Turn him into Deadpool” with DGE. 193

194 For the results of our method reported in Tab. 1, all 3D
195 edits are performed using DGE. The applied prompts are:

- **Sear-steak.** “Turn him into a Minecraft character”, “Turn
him into Super Mario”, 196
- **Coffee-martini.** “Turn his hat into a newsboy cap”,
“Turn him into JoJo’s Bizarre Adventure anime style”, 197
- **Trimming.** “Change the leaves to autumn leaves”, “Turn
him into Deadpool”. 198

199 For the results in Tab. 2, the editing setups are:

- “Turn him into a Minecraft character” with DGE, 200
- “Turn him into Batman” with DreamCatalyst, 201
- “Turn him into JoJo’s Bizarre Adventure anime style”
with DGE. 202

203 D. Limitation

204 While Catalyst4D effectively extends high-quality 3D edit-
205 ing capabilities to dynamic 4D scenes, it inherits certain de-
206 pendencies from the underlying components or pretrained
207 models. 208

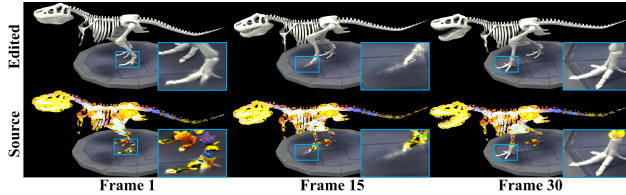


Figure D.1. Failure case (motion jitter illustration).

219 First, the temporal coherence of our method is influenced
 220 by the quality of the initial 3D edits. If the 3D editing re-
 221 sults lack sufficient spatial consistency, this may propagate
 222 and affect the spatio-temporal consistency in the final dy-
 223 namic outputs. Moreover, as our method does not modify
 224 the deformation network or re-optimize Gaussian densities,
 225 it relies on the stability of the underlying 4D reconstruc-
 226 tion. Under severe reconstruction noise (e.g., point jitter
 227 or low-opacity Gaussians), motion guidance may be locally
 228 disrupted. In Fig. D.1 (*D-NeRF, trex*), background Gaus-
 229 sians drift into the edited foreground, causing local artifacts,
 230 reflecting a shared limitation of current 4D reconstruction
 231 models. Nonetheless, such limitations are not fundamental
 232 and can be mitigated as upstream 3D editing and 4D recon-
 233 struction techniques continue to advance. Our framework is
 234 fully compatible with future improvements in both areas.

235 References

- 236 [1] Marco Cuturi. Sinkhorn distances: Lightspeed computation
 237 of optimal transport. In *NeurIPS*, 2013. 3
- 238 [2] Bruno Galerne, Jianling Wang, Lara Raad, and Jean-Michel
 239 Morel. Sgsst: Scaling gaussian splatting style transfer. In
 240 *CVPR*, pages 26535–26544, 2025. 2, 3
- 241 [3] Kai He, Chin-Hsuan Wu, and Igor Gilitschenski. Ctrl-d: Con-
 242 trollable dynamic 3d scene editing with personalized 2d diffu-
 243 sion. In *CVPR*, pages 26630–26640, 2025. 1
- 244 [4] Kornel Howil, Piotr Borycki, Tadeusz Dziarmaga, Marcin
 245 Mazur, Przemysław Spurek, et al. Clipgaussian: Universal and
 246 multimodal style transfer based on gaussian splatting. *arXiv preprint arXiv:2505.22854*, 2025. 2
- 248 [5] Junha Hyung, Susung Hong, Sungwon Hwang, Jaeseong Lee,
 249 Jaegul Choo, and Jin-Hwa Kim. Effective rank analysis and
 250 regularization for enhanced 3d gaussian splatting. In *NeurIPS*,
 251 pages 110412–110435, 2024. 4
- 252 [6] Xin Luo, Jiahao Wang, Chenyuan Wu, Shitao Xiao, Xiyan
 253 Jiang, Defu Lian, Jiajun Zhang, Dong Liu, et al. Editscore:
 254 Unlocking online rl for image editing via high-fidelity reward
 255 modeling. *arXiv preprint arXiv:2509.23909*, 2025. 1
- 256 [7] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T
 257 Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-
 258 Brualla, and Steven M Seitz. Hypernerf: A higher-
 259 dimensional representation for topologically varying neural
 260 radiance fields. *arXiv preprint arXiv:2106.13228*, 2021. 1
- 261 [8] Shangkun Sun, Xiaoyu Liang, Songlin Fan, Wenxu Gao, and
 262 Wei Gao. Ve-bench: Subjective-aligned benchmark suite for

- text-driven video editing quality assessment. In *AAAI*, pages
 7105–7113, 2025. 1 263
 264
 [9] Guanjun Wu, Taoran Yi, Jiemin Fang, Lingxi Xie, Xiaopeng
 Zhang, Wei Wei, Wenyu Liu, Qi Tian, and Xinggang Wang.
 4d gaussian splatting for real-time dynamic scene rendering.
 In *CVPR*, pages 20310–20320, 2024. 4 265
 266
 267
 268