



4D-Editor: Interactive Object-level Editing in Dynamic Neural Radiance Fields via Semantic Distillation

Supplementary Material

1. Implementation Details 001

002 We use Principal Component Analysis (PCA) method to ex-003 tract 64 most important semantic features from DINO [1] ViT-b8 model. We use 64 sampled points on each ray in 004 005 volume rendering.

During the editing process, we set different threshold α 006 and number of K-Means clusters N_k according to different 007 008 sizes of target objects: $\alpha = 0.4 \pm 0.1$, $N_k = 3 \pm 2$ for small 009 objects and $\alpha = 0.7 \pm 0.1$, $N_k = 20 \pm 15$ for small ones. 010 The value of β should be approximately 1/3 of α . Actually, in our experiments, all these values can be set loosely within 011 Recursive Selection Refinement. 012

2. Interactive GUI 013

We design a user-friendly graphical user interface (GUI) 014 015 that facilitates interactive editing of dynamic scenes. The 016 editing procedure involves three steps: 1) Selecting a refer-017 ence frame, 2) Marking target objects with strokes in distinct colors, and 3) Configuring editing parameters, includ-018 ing editing operations, threshold α , exploration range β and 019 020 recursion depth K. These steps are depicted in Fig. ??.

3. Model Structure 021

3.1. Structures of Hybrid Radiance Fields 022

023 We directly adopt RobustNeRF [2] and DynamicNeRF [4] to represent a dynmaic scene, their methods are displayed 024 in Fig. 9. 025

3.2. Structures of Hybrid Semantic Fields 026

We design semantic fields G^s and G^d , to represent seman-027 028 tic information of the static and dynamic components in the 029 scene, respectively. Both semantic fields are modeled using 030 an 8-layer multi-layer perceptron (MLP). As illustrated in Fig. 1, G^s takes the position **x** as input after position encod-031 032 ing, while G^d incorporates both the position **x** and the time 033 t.

3.3. Volume Rendering on Semantic Features 034

The spatial semantic features themselves are represented 035 036 by a 64-dimensional vector. We apply volume rendering 037 to these features and select the first 3 dimensions to gen-038 erate RGB visualizations. As shown in Fig. 2, upon re-039 constructing the spatial semantic information using seman-040 tic fields, our approach demonstrates superior preservation



Figure 1. Structures of Semantic Fields.



Figure 2. Visualization of Semantic Features. After semantic features distillation, we obtain enhanced semantic features in 4D space (For additional comparisons, refer to the demo videos in the supplementary material).

of both multi-view and spatio-temporal consistency of se-041 mantic information in the 4D space, compared to the rela-042 tively coarse-grained feature map generated by the DINO 043 model [1]. Notably, on object edges, the transition of se-044 mantic information appears remarkably smooth. 045

4. Visualization of Recursive Selection Refine-046 ment

Figure 3 shows the flow of semantic feature matching and 048 target selection(semantic segmentation). Given sampled 049 points $p_1 \dots p_8$, we categorize them based on feature dis-050 tances: valid points p_4 p_5 , impossible points p_1 p_2 p_8 , pos-051 sible points $p_3 p_6 p_7$. Then we introduce a random offset on 052

047



3DV 2025 Submission #*****. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

K	$new \mathcal{V}$	\mathcal{P}	all \mathcal{V}	IoU	Acc
1	2.84M	828.13k	2.84M	72.79	79.81
5	20.69k	120.73k	3.18M	76.53	85.39
20	426	6.55k	3.23M	80.85	92.70
30	114	1.64k	3.24M	78.63	93.45
40	8	169	3.24M	80.11	93.24
50	1	56	3.24M	79.21	93.88

Table 1. Effects of maximum recursion number K.

possible points and recalculate feature distances, repeating 053 054 until these points are judged as valid or impossible. Here, 055 $p_3 p_6$ are considered as valid points after 1 or 2 iterations re-056 spectively, while p_7 excluded after 2 iterations. Ultimately, the selected points are p_3 , p_4 , p_5 , p_6 . We will repeat this 057 process for all sampling points on each ray during the ren-058 dering process, as shown in Algorithm 4. As for the same 059 batch of sampling points, we can apply different random 060 offset step s at one time (e.g., s = (1/50, 1/100, 1/150)) to 061 control different ranges of selection on target. 062

In RSR, we initially get N possible points near the edges 063 of the object and M valid points. Given the feature points 064 around the object's edges, we determine the probabilities 065 of a point with an offset, being classified as valid with a 066 possibility of μ or **possible** with a possibility of λ . Note 067 that $\mu + \lambda < 1$ due to the presence of invalid points, where 068 069 λ also indicates the probability of progressing to the next 070 round. For each point, the expected probability of being 071 valid is $\lambda + \lambda \mu + \lambda^2 \mu + \dots$ This is a geometric progression, 072 the sum of which is $\lambda/(1-\mu)$. Finally, the expected number 073 of points that will be selected is $M + N\lambda/(1-\mu)$.



Figure 3. **Process of Recursive Selection Refinement.** Although DINO model generates coarse semantic feature maps which are used as the training set, our reconstruction and volume rendering of the 4D spatial semantic information results in finer semantic feature maps.

5. Analysis on Object Segmentation

In this section, we begin by conducting a qualitative and
quantitative comparison between N3F [3] and our method in
terms of object segmentation performance. Subsequently,

nput: $U, \mathcal{A}, K, \alpha, \beta, \gamma, s, k = 0$
Dutput: valid point index set W
Def RSR($\mathcal{U}, \mathcal{A}, k, \alpha, \beta, \gamma, s$):
₩ ← []
if $k = K$ or $\mathcal{U} = \emptyset$: return \mathcal{W}
<pre>// Calculate feature distances</pre>
$\mathcal{D} \leftarrow \{d(\gamma, u_i), u_i \in \mathcal{U}\}$
// Store valid points & indexes
$\mathcal{V} \leftarrow \{d(\gamma, u_i) \leq \alpha, d(\gamma, u_i) \in \mathcal{D}\}\$
\mathcal{W} .append(index of \mathcal{V} from \mathcal{A})
// Dismiss impossible points
$Q \leftarrow \{d(\gamma, u_i) > \alpha + \beta, d(\gamma, u_i) \in \mathcal{D}\}$
// Get possible points and indexes
$\mathcal{P} \leftarrow \{\alpha < d(\gamma, u_i) \le \alpha + \beta, d(\gamma, u_i) \in \mathcal{D}\}$
$\mathcal{A}' \leftarrow \text{index of } \mathcal{P} \text{ from } \mathcal{A}$
// Apply random offsets
$\mathcal{P}' \leftarrow offset(\mathcal{P}, randn(0, s))$
$\mathcal{W}' \leftarrow RSR\left(\mathcal{P}', \mathcal{A}', k+1, \alpha, \beta, \gamma, s\right)$
\mathcal{W} .append(index of \mathcal{V}' from \mathcal{A})
mataram (11/

Figure 4. Recursive Selection Refinement



Figure 5. **Results of Object Segmentation.** Compared with N3F, we achieve finer object segmentation, which contains fewer artifacts especially for edge areas.

we evaluate the impact of the hyperparameter β (explo-	078
ration range) on Recursive Selection Refinement.	079

5.1. Evaluation on Segmentation Accuracy

We fine-tune the threshold of N3F for object segmenta-081 tion tasks and dispaly its best results in the right column 082 of Fig. 5. Since our method utilizes Recursive Selection 083 Refinement to achieve more precise selection of target ob-084 jects, surpassing N3F in both qualitative and quantitative 085 assessments (indicated by higher Acc and IoU values in Ta-086 ble. 2). In our experiments, we observe that our method per-087 forms especially well in scenes containing significant ob-088 ject movement (e.g., a girl starting from the left side and 089 jumping to the right side in the Rollerblade dataset). Con-090 versely, N3F produces numerous broken artifacts in such 091 scenes. Fig. 6 presents additional results of object segmen-092 tation in 4D space. 093

5.2. Analysis on exploration range β

The selection of the exploration range parameter β can be loose and flexible. Fig. 3 demonstrates that the segmentation quality remains consistent despite using different val-

3DV #*****

080

094

115

116

129

130

131

132

133

134

135

136

Table 2. **Quantitative Analysis on Object Segmentation.** We choose several scenes from DAVIS dataset and use its true motion masks as ground truth.

Scene	Metric	N3F	Ours
Breakdance Flare	Acc ↑	89.96	94.12
Dreukaance Frare	IoU ↑	84.97	83.09
Rollerblade	Acc ↑	83.38	93.10
Konerblade	IoU ↑	71.55	81.56
Hika	Acc ↑	94.80	95.16
IIIKE	IoU ↑	85.73	89.88



Figure 6. Object Segmentation in 4D Space.

Table 3. Effects of exploration range β . We set the parameters as follows: K = 3, $N_k = 8$, and $\alpha = 0.6$. Additionally, we experiment with different values of β to demonstrate that it has no impact on the segmentation quality. Consequently, β can be loosely set during editing operations.

β	0.05	0.08	0.1	0.12	0.15	0.18	0.2
Acc ↑	93.10	93.42	93.81	93.66	93.94	94.02	93.94
IoU↑	81.56	81.26	80.90	80.77	80.21	79.65	79.26

1098 ues of β . This finding supports the notion that it is the recursion depth parameter K that significantly enhances segmentation accuracy, while β merely serves as an auxiliary factor.

102 6. Scene Editing

This section focuses on recoloring and two complex editingoperations: transformation and composition(Fig. 8).

105 6.1. Recoloring

106 Recoloring on HSL is showed in Fig. 7.

107 6.2. Transformation Details

Table 4 shows three types of geometric transformation operations (*Shift, Scale*, and *Mirror*) and a time-variant operation (*Reverse*). The time-variant transformation allows
editing in the temporal dimension. For example, we can reverse the trajectory of moving objects in a time series while
maintaining the others (*e.g.*, Fig. 8b shows two boys: the
boy in the red box follows the real trajectory to the right







Figure 7. Recolor the balloon.

(f) lightness ↑

Table 4. Transformation Functions.			
	Mapping Function		
Shift	$(x, y, z) : (x, y, z) + \delta$		
Scale	$(x, y, z) : (x, y, z) \times scale$		
Mirror	(x,yz):(-x,-y,z)		
Reverse	(x, y, z, t) : (x, y, z, -t)		

Table 5. Final σ , c after Transformation.

	,
$u_i \in \mathcal{T}'$	$(\sigma, c) \leftarrow (\sigma^t c^t)$
$u_i \in \mathcal{T},$	$(\sigma, c) \leftarrow (\sigma^o c^o)$ if reserved
$u_i \notin \mathcal{T}'$	$\sigma \leftarrow 0$ if not reserved
$u_i \notin \mathcal{T},$	$(\sigma, c) \leftarrow (\sigma^o c^o)$
$u_i \notin \mathcal{T}'$	

side, but the boy in the blue box moves towards the opposite direction and reverses to the left side).

In dynamic scenes, we lack explicit knowledge regarding 117 the distribution of different objects in both the current space 118 and the space after transformation. We can only judge the 119 selected target region by calculating the feature distance. 120 Therefore, we perform two rounds of calculations for all 121 sampled points. In the first round, we use initial positions 122 of all sampled points, dividing them into \mathcal{T} (inside the ob-123 ject) and S (outside) according to Section 3.3. In the second 124 round, we use the positions after transformation, dividing 125 edited space into \mathcal{T}' and \mathcal{S}' . Additionally, for each sampled 126 point u_i , we obtain its original attributes c^o and σ^o , as well 127 as c^t and σ^t after transformation. 128

To preserve the original object, it is kept intact. Alternatively, if removal is desired, we set $\sigma = 0$. In cases where there are overlapping areas between the transformed object and the original object, we consistently apply c^t and σ^t to assign the properties to these points. The settings for the final density and color can be found in Table 5 provided below.

6.3. Composition Details

Assuming the presence of N objects derived from distinct dynamic scenes, which have been filtered from the original NeRF and are represented using masks in 4D space. Trans-139 3DV #*****

169

170



(b) Time reversal

Figure 8. **Composition & Time Reversal.** We maintain spatial-temporal consistency across the entire time series when combining different scenes or applying time-variant transformations to individual objects within the same scene.

140 formation can also be applied here to prevent object overlap

during composition. The composition is performed at eachsample point for all segmented parts collectively:

(43)
$$(\sigma', \mathbf{c}') = \left(\sum_{i=1}^{N} \sigma_i \ mask_i, \sum_{i=1}^{N} c_i \ mask_i\right)$$
(1)

144 where $mask_i$ denotes the membership of the current point 145 to the *i*th object. Additionally, blending weights for dy-146 namic objects must be recalculated:

$$b' = \sum_{i=1}^{N} b_i \; mask_i \tag{2}$$

In Fig. 8a, we create a novel scene by separately using
moving objects from the *Balloon2* dataset and backgrounds
from the *Playground* dataset, and then compositing them
together. Subsequently, the rendering formula can be employed to compute the pixel color. The computational complexity of the combination increases in proportion to the
number of objects.

155 References

147

- [1] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou,
 Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on com- puter vision*, pages 9650–9660, 2021. 1
- [2] Yu-Lun Liu, Chen Gao, Andreas Meuleman, Hung-Yu Tseng,
 Ayush Saraf, Changil Kim, Yung-Yu Chuang, Johannes Kopf,
 and Jia-Bin Huang. Robust dynamic radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision*and Pattern Recognition, pages 13–23, 2023. 1, 4
- [3] Vadim Tschernezki, Iro Laina, Diane Larlus, and Andrea
 Vedaldi. Neural feature fusion fields: 3d distillation of selfsupervised 2d image representations. In 2022 International



(b) Dynamic NeRF [4]

Figure 9. Hybrid Radiance Field Reconstruction of Dynamic Scenes. The structures of these two methods are directly adopted from their respective theses.

Conference on 3D Vision (3DV), pages 443–453. IEEE, 2022.

[4] Jae Shin Yoon, Kihwan Kim, Orazio Gallo, Hyun Soo Park, and Jan Kautz. Novel view synthesis of dynamic scenes with globally coherent depths from a monocular camera. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5336–5345, 2020. 1, 4
[4] Jae Shin Yoon, Kihwan Kim, Orazio Gallo, Hyun Soo Park, and Pattern Recognition, pages 5336–5345, 2020. 1, 4