

000  
001       **Unsupervised Discovery and Composition of**  
002       **Object Light Fields**  
003       **–Supplemental Material–**

004  
005              Anonymous ECCV submission  
006  
007              Paper ID 4753

011       **Abstract.** In this supplementary document, we first describe details  
012       on architecture in Sec. 1, dataset construction in Sec. 2, and training  
013       protocol in Sec. 3. We then provide additional algorithm pseudocode in  
014       Sec. 4 and qualitative results in Sec. 5.

015  
016       

## 1 Architecture

  
017

021       

### 1.1 Encoder

023       **Feature Extractor** We use the same encoder as [5] — 6 convolutional layers  
024       with bilinear upsampling applied to the last 3. Each layer has kernel size 3,  
025       padding of 1, and stride of 1 except for the second and third layers which use 2.  
026       Pixel coordinates are normalized to the range of [-1,1] in both directions, leading  
027       to 4 additional input channels to concatenate with the 3 image color channels.

029  
030       **Background-Aware Slot Encoder** We use the same slot encoder parameters  
031       as [5] — a slot dimension of 128 and 3 iterations of slot competition.

033       

### 1.2 Decoder

035       **Light Field Networks** For both the foreground and background light fields,  
036       we implement them as MLP’s of 2 hidden layers with 256 hidden features and  
037       ReLU activations. They yield ray-features of dimensionality 64.

039  
040       **Color, Depth, and Visibility Networks** The color generator and depth  
041       decoder map the ray feature from an object-LFN into color and depth, and the  
042       visibility network assigns visibility weight based on the LFN’s decoded depth and  
043       relative depth; all three networks are implemented as MLPs of 3 hidden layers  
044       with 128 hidden features and ReLU activations.

**HyperNetwork** We map the slot latent codes of dimensionality 128 to the weights of each light field MLP via a hypernetwork [1], as performed by [4,2,3]. The hypernetwork is implemented as an MLP with 1 hidden layer which accepts a 128-dimensional slot latent code and predicts all the weights of the corresponding light field.

## 2 Dataset Details

### 2.1 Room-Scenes (Clevr-567, Room-Chair, Room-Diverse)

We refer the reader to the supplement of [5] as they produced these datasets and detail their creation there.

### 2.2 City-Block

The city block is constructed with one road and several buildings, obtained from an online 3D-assets website. The camera is always facing forward in the middle of the lane, with height sampled from a range of near the ground to slightly above the car height, and depth in the scene sampled from near the front to the end of the block. Context views are always captured at the same distance from each row of cars. The training scenes always consist of two rows of cars with the back row placed a fixed distance from the front row. The cars in each row have a small difference in position.

## 3 Training

We optimize our models' parameters using the ADAM solver with learning rate of  $5 \times 10^{-5}$ . We report the training supervision schedules for each dataset evaluation below.

### 3.1 CLEVR-567

We train first at a resolution of 64x64 for 150k iterations and at a resolution of 128x128 for 60k iterations. We supervise with a combination of  $L1, L2$ , and perceptual loss [6].

### 3.2 Room-Chair and Room-Diverse

We initialize the model with the weights of the CLEVR-567 model and train for 128k iterations at 64x64, then for 90k iterations at 128x128 resolution. We supervise with a combination of  $L1, L2$ , and perceptual loss [6].

### 3.3 Multi-Lane Highway

We initialize the model with the weights of the CLEVR-567 model. We then train with the  $L2$  reconstruction loss at a resolution of 64x64 for 145k iterations.

### 090 3.4 City Block

091 We pretrain the static background network on the city block dataset at 64x64  
 092 resolution for 800k iterations with  $L_2$  loss. Then, foreground slots are introduced,  
 093 initialized with weights from the CLEVR-567 model, and trained for 60k iterations  
 094 at 64x64 resolution and 50k iterations at 128x128 resolution, supervised with  $L_2$   
 095 loss.

## 097 4 Background-Aware Slot Encoder Pseudocode

100 The pseudocode for [5]’s background-aware modification to the slot attention  
 101 algorithm is presented below at the courtesy of its authors:

---

102 **Algorithm 1:** Object-centric latent inference with background-aware  
 103 slot attention.

---

104 **Input:**  $\text{feat} \in \mathbb{R}^{N \times D}$   
 105 **Learnable:**  $\mu^b, \sigma^b, \mu^f, \sigma^f$ : prior parameters,  $k, q^b, q^f, v^b, v^f$ : linear mappings,  
 106  $\text{GRU}^b, \text{GRU}^f, \text{MLP}^b, \text{MLP}^f$

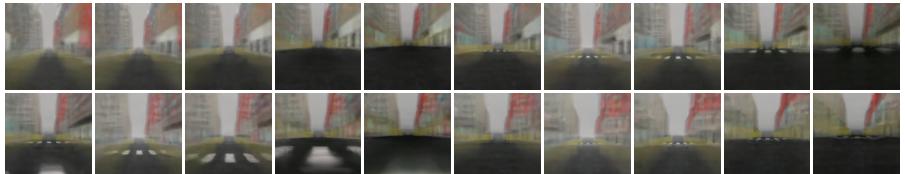
107    $\text{slot}^b \sim \mathcal{N}^b \in \mathbb{R}^{1 \times D}$   
 108    $\text{slots}^f \sim \mathcal{N}^f \in \mathbb{R}^{K \times D}$   
 109   **for**  $t = 1, \dots, T$  **do**  
 110      $\text{slot\_prev}^b = \text{slot}^b, \text{slots\_prev}^f = \text{slots}^f$   
 111      $\text{attn} = \text{Softmax}\left(\frac{1}{\sqrt{D}} k(\text{feat}) \cdot \begin{bmatrix} q^b(\text{slot}^b) \\ q^f(\text{slots}^f) \end{bmatrix}^T, \text{dim}=\text{'slot'}\right)$   
 112      $\text{attn}^b = \text{attn}[0], \text{attn}^f = \text{attn}[1:\text{end}]$   
 113      $\text{updates}^b = \text{WeightedMean}(\text{weights}=\text{attn}^b, \text{values}=\text{v}^b(\text{inputs}))$   
 114      $\text{updates}^f = \text{WeightedMean}(\text{weights}=\text{attn}^f, \text{values}=\text{v}^f(\text{inputs}))$   
 115      $\text{slot}^b = \text{GRU}^b(\text{state}=\text{slot\_prev}^b, \text{inputs}=\text{updates}^b)$   
 116      $\text{slots}^f = \text{GRU}^f(\text{state}=\text{slots\_prev}^f, \text{inputs}=\text{updates}^f)$   
 117      $\text{slot}^b += \text{MLP}^b(\text{slot}^b), \text{slots}^f += \text{MLP}^f(\text{slots}^f)$   
 118   **end**  
 119   **return**    $\text{slot}^b, \text{slots}^f$

---

## 135 5 Additional Results

### 136 137 5.1 Comparison to Baseline on Unbounded Scene

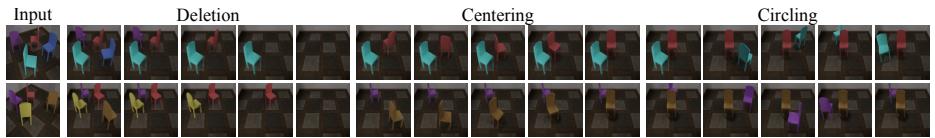
138 We evaluate the baseline model [5] on the unbounded City Block scene and  
 139 illustrate a sample in Fig. 1. Employing a volumetric decoder, the model is  
 140 forced by its memory constraints and the large bounds of the scene to sample  
 141 coarsely between the near and distant far plane. As a result of their coarse  
 142 sampling, their model is unable to learn to render any foreground elements of  
 143 the scene.



144  
 145  
 146 Fig. 1: The baseline model’s [5] reconstruction on the City Block scene. Using a vol-  
 147 umetric decoder, the model is unable to render smaller foreground objects due to its  
 148 coarse sampling through the long scene.  
 149  
 150  
 151

### 152 153 154 155 5.2 Additional Scene Manipulation Results

156 We perform additional scene manipulation demonstrations in Fig. 2.



157  
 158  
 159  
 160 Fig. 2: We demonstrate object-level scene editing tasks of object deletion, centering,  
 161 and circling, on scenes from the chairs dataset.  
 162  
 163  
 164

### 171 172 173 174 5.3 Additional Decomposition and Novel View Synthesis Results

175 We provide additional novel view synthesis results and decomposition on four  
 176 scenes from each room dataset in Fig. 3 and the composited city block scene in  
 177 Fig. 4.

178  
 179

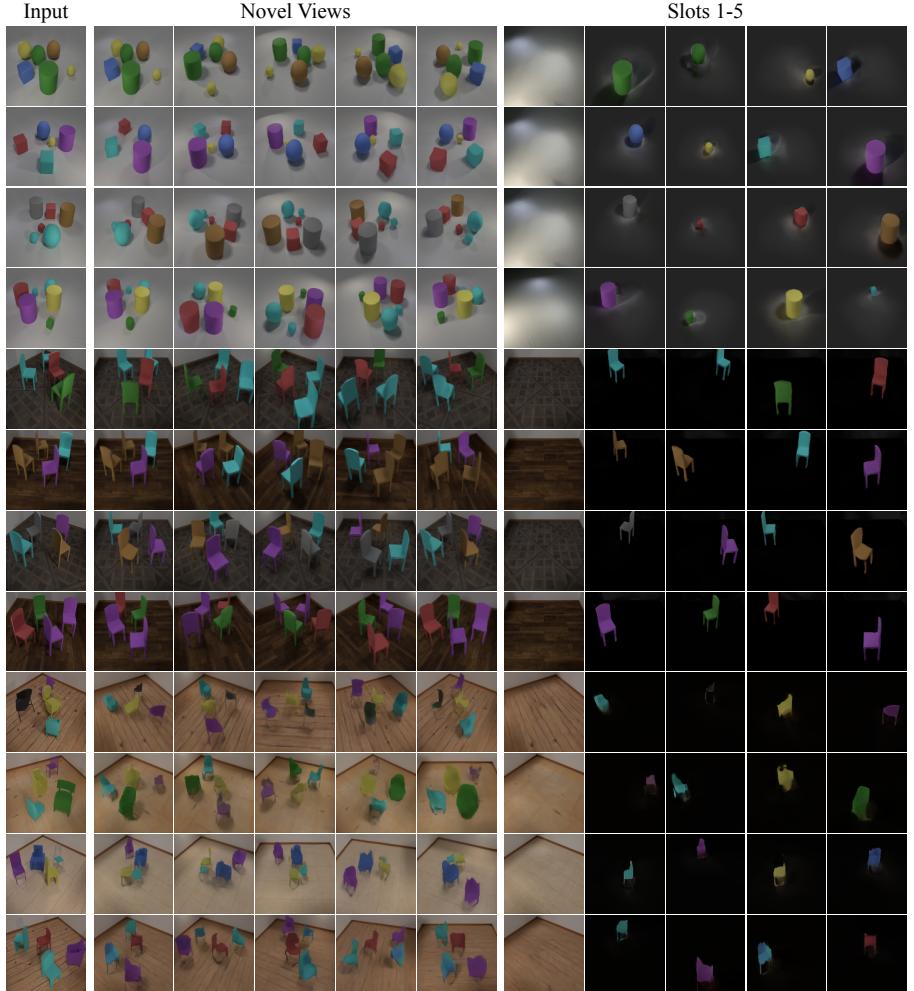
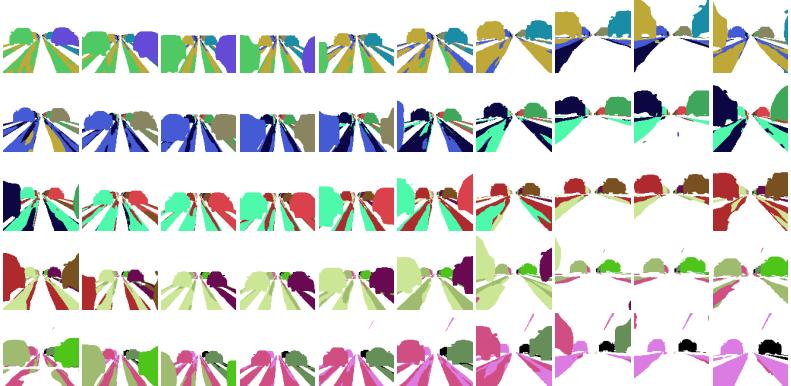


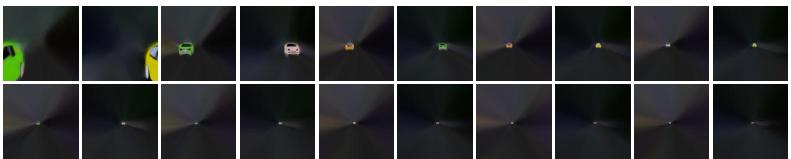
Fig. 3: We demonstrate additional qualitative results for the room-scene datasets on tasks of novel view synthesis and scene decomposition. We show five novel views at the scene level (middle) and five slots from the first novel view (right).



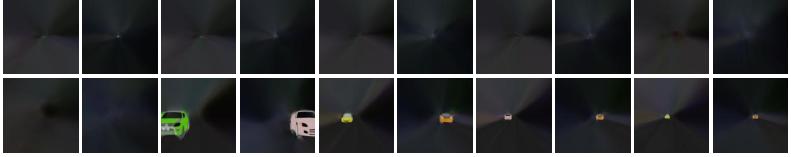
(a) Novel views of composite scene (left to right, top to bottom)



(b) Novel segmentation from composite scene (left to right, top to bottom)



(c) Object light fields at frame 4



(d) Object light fields at frame 35

Fig. 4: We render a sequence of novel views from the cross-scene composition application described (first), as well as each view's segmentation (second), and per-object contributions from two frames (third, fourth).

## 270 References

- 271 1. Ha, D., Dai, A., Le, Q.V.: Hypernetworks. In: Proc. ICLR (2017) 2
- 272 2. Sitzmann, V., Martel, J.N., Bergman, A.W., Lindell, D.B., Wetzstein, G.: Implicit  
273 neural representations with periodic activation functions. In: Proc. NeurIPS (2020)  
2
- 274 3. Sitzmann, V., Rezhikov, S., Freeman, W.T., Tenenbaum, J.B., Durand, F.: Light  
275 field networks: Neural scene representations with single-evaluation rendering. In:  
276 Proc. NeurIPS (2021) 2
- 277 4. Sitzmann, V., Zollhöfer, M., Wetzstein, G.: Scene representation networks: Con-  
278 tinuous 3d-structure-aware neural scene representations. In: Advances in Neural  
279 Information Processing Systems. pp. 1121–1132 (2019) 2
- 280 5. Yu, H.X., Guibas, L.J., Wu, J.: Unsupervised discovery of object radiance fields.  
281 arXiv preprint arXiv:2107.07905 (2021) 1, 2, 3, 4
- 282 6. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable  
283 effectiveness of deep features as a perceptual metric. In: Proceedings of the IEEE  
284 conference on computer vision and pattern recognition. pp. 586–595 (2018) 2