

Supplemental document for: Learning 3D object-centric representation through prediction

3

4

January 25, 2024

5

The major codes of the OPPL model are provided together with this document.

6

1 Additional demonstration of OPPL’s inference and prediction

8

9 1.1 Visualization of outputs by the learned networks and prediction for the next frame

10

As visualized in Figure S1, Frame 2 (column 1) is the basis of prediction (frame 1 is not shown). Frame 3 (ground truth in column 2) is predicted (column 3) by combining two approaches: warping-based prediction (column 4) and imagination (column 5). The warping-based prediction is made based on the optical flow predicted from the inferred depth (column 8), object segmentation (the last 4 columns), inferred object location and velocity, and the camera’s motion. In this experiment, the maximum possible number of visible objects is set to 3.

11

As expected it can be seen that warping-based prediction generates sharper images (column 4) than imagination-based prediction (column 5). This is because the latter is based on learned statistical regularity of environments and is only needed for small regions in the image (thus receiving less teaching signal). This relaxation of the need to encode fine details of a scene in the imagination network is advantageous because the goal of object perception is primarily to obtain an estimation of 3D information and coarse features of an object. The warping-based prediction also does not require encoding fine details of each object. The warping weights (column 6) is the total contribution of warping-based prediction towards the final prediction in each pixel.

12

29 1.2 Depth perception

We demonstrate a few example images, their ground truth depth and the depth inferred by OPPL. Our network can capture the global 3D structure of the

30

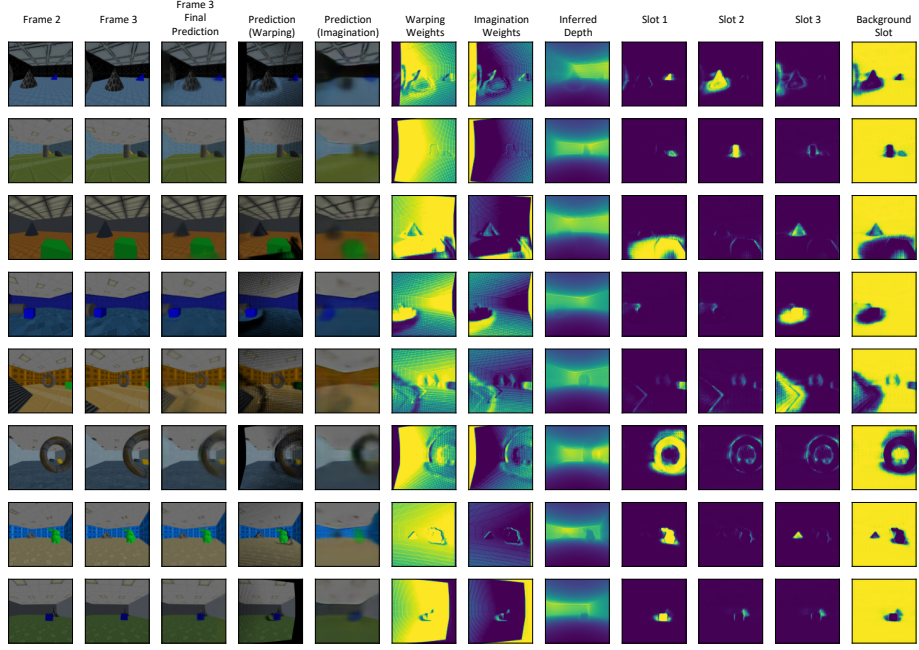


Figure S1: OPPL's internal segmentation, depth inference and prediction for triplets in our dataset

scene, although the estimation of the depth of the objects tends to be biased towards the depth of the background behind them, consistent with the scatter plot in the main paper. We suspect that by introducing more translational motion of the camera parallel to the image plane will allow the depth network to learn more accurate depth inference for objects, as such camera movement induces strong motion parallax effect. We see evidence for this in our results on the GQN dataset, as discussed in the next section.

We would like to emphasize that the depth inference is performed on single images without the need of motion cues, unlike other unsupervised models such as O3V and Savi++ that need videos at both training and test time to learn and infer 3D information.

1.3 Performance on GQN dataset

We evaluated our model on the ROOMS dataset using the rendering code provided by [5]. The dataset is based on the GQN rooms-ring-camera dataset [3] but we allowed all motion speeds to be uniformly sampled from an interval instead of being chosen from narrow Gaussian distributions as in [5]. It appears that the depth images (column 8) are sharper and exhibit less bias for objects' depth towards longer distance than in Figure S1 (on our custom dataset). This is potentially because in this dataset, the dominant motion of the camera is

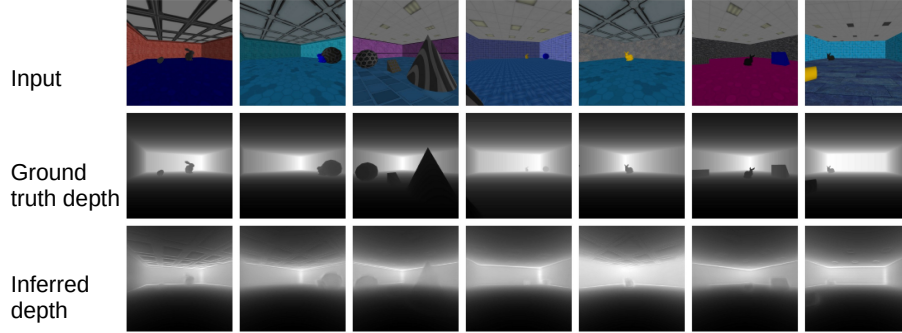


Figure S2: Comparison between ground truth depth and inferred depth

parallel to the image plane (along the ring while looking to a direction close to the center of the room), consistent with our hypothesis above. This visualization also demonstrates that although we set a maximum number of slots (objects) visible to the camera (3 in these experiments), the network is able to output empty object masks when the actual number of visible objects is smaller than the number of slots. In other words, it is not required for the network to know how many objects are in the scene. Only a maximum number of slots needs to be set for the LSTM within the object extraction network.

As shown below, the depth inference of OPPL on GQN dataset reaches superior performance of $r=0.99$ between ground truth depth and inferred depth (excluding pixels in the sky, because the sky's distance is effectively infinity in the ground truth). A scatter plot of inferred vs. ground truth depth for pixels excluding the sky is shown below (the vertical stripes are likely due to rounding issue in saving depth map and down-sampling).

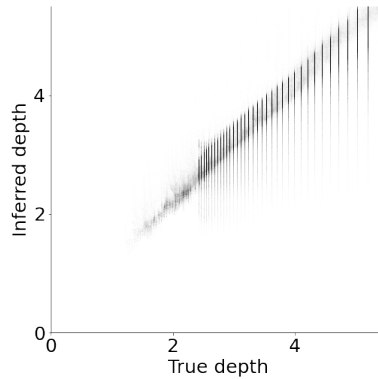


Figure S4: Performance of depth inference of OPPL on GQN dataset.

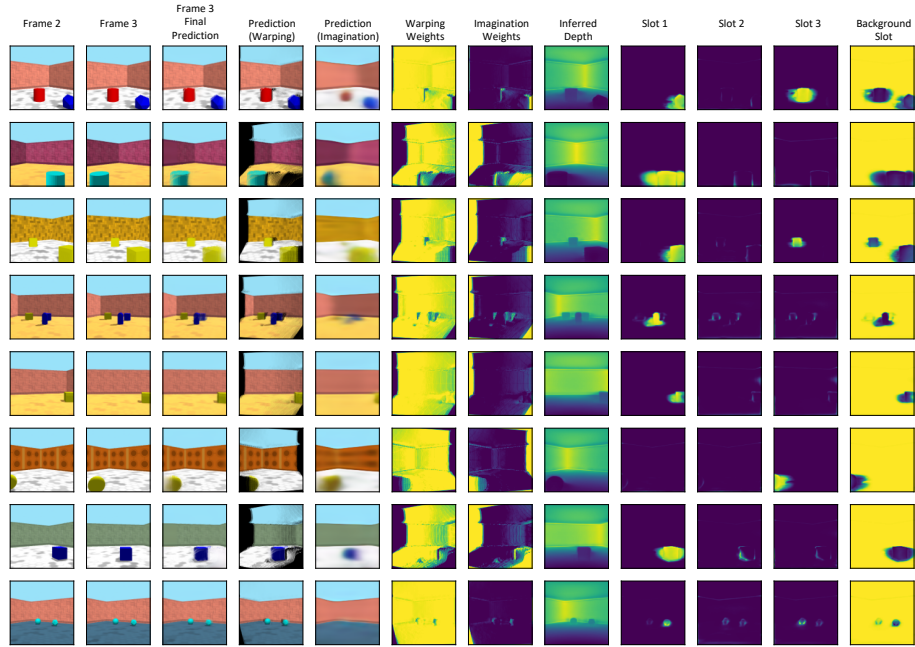


Figure S3: Visualization of the two approaches of prediction (warping and imagination), their combination weights, inferred depth, and probabilistic segmentation maps from OPPL on scenes of moving objects in a simulated room based on the ring environment of GQN[3], based on code from [5].

2 Pseudo code of the OPPLP framework

Algorithm 1 Object Perception by Predictive LEarning (OPPLE)

initialize $f_{\theta_{\text{obj}}}$, $h_{\theta_{\text{depth}}}$, $g_{\theta_{\text{imag}}}$

Input: images $\mathbf{I}^{(t-1)}, \mathbf{I}^{(t)}, \mathbf{I}^{(t+1)} \in \mathbb{R}^{w \times h \times 3}$, self-motion
 $\mathbf{v}_{\text{obs}}^{(t-1)}, \omega_{\text{obs}}^{(t-1)}, \mathbf{v}_{\text{obs}}^{(t)}, \omega_{\text{obs}}^{(t)}$

Output: prediction $\mathbf{I}'^{(t+1)}$, segmentation $\pi_{1:K+1}^{(t-1)}, \pi_{1:K+1}^{(t)}$, objects' codes $\mathbf{z}_{1:K}^{(t-1)}, \mathbf{z}_{1:K}^{(t)}$, objects' locations and poses $\hat{\mathbf{x}}_{1:K}^{(t-1)}, \mathbf{p}_{\phi_{1:K}}^{(t-1)}, \hat{\mathbf{x}}_{1:K}^{(t)}, \mathbf{p}_{\phi_{1:K}}^{(t)}$

for $\tau = \{t-1, t\}$ **do**

scene code $e^{(\tau)} \leftarrow \text{U-NetEncoder}_{f_{\theta_{\text{obj}}}}(\mathbf{I}^{(\tau)})$

object code $\mathbf{z}_{1:K}^{(\tau)}$, location $\hat{\mathbf{x}}_{1:K}^{(\tau)}$, pose $\mathbf{p}_{\phi_{1:K}}^{(\tau)} \leftarrow \text{LSTM}_{f_{\theta_{\text{obj}}}}(e^{(\tau)})$

background code $\mathbf{z}_{K+1} = 0$

depth $\mathbf{D}^{(\tau)} \leftarrow h_{\theta_{\text{depth}}}(\mathbf{I}^{(\tau)})$

segmentation mask $\pi_{1:K+1}^{(\tau)} \leftarrow \text{Softmax}([\text{U-NetDecoder}_{f_{\theta_{\text{obj}}}}(\mathbf{I}^{(\tau)}, \mathbf{z}_{1:K}^{(\tau)}), 0])$

end for

object matching scores $r_{kl} \leftarrow \frac{\text{RBF}(\mathbf{z}_k^{(t)}, \mathbf{z}_l^{(t-1)})}{\sum_{m=1}^{K+1} \text{RBF}(\mathbf{z}_k^{(t)}, \mathbf{z}_m^{(t-1)})}, k, l \in 1 : K + 1$

for $k \leftarrow 1$ to K **do**

object motion $\hat{\mathbf{v}}_{1:K}, \omega_{1:K} \leftarrow r_{k,l}, \hat{\mathbf{x}}_k^{(t)}, \hat{\mathbf{x}}_l^{(t-1)}, \mathbf{p}_{\phi_k}^{(t)}, \mathbf{p}_{\phi_l}^{(t-1)}, \mathbf{v}_{\text{obs}}^{(t-1)}, \omega_{\text{obs}}^{(t-1)}, l = 1 : K + 1$

object-specific optical flow $\mathbf{v}_k^{(t)}, \omega_k^{(t)}, \mathbf{D}^{(t)}, \hat{\mathbf{x}}_k^{(t)}$

end for

warping-based prediction $\mathbf{I}'_{\text{warp}}^{(t+1)}, \mathbf{W}_{\text{Warp}} \leftarrow$

$\text{Warp}(\mathbf{I}^{(t)}, \text{optical flow}_{1:K+1}, \pi_{1:K+1}^{(\tau)})$

imagination $\mathbf{I}'_{\text{imagine}}^{(t+1)} \leftarrow g_{\theta_{\text{obj}}}(\mathbf{I}^{(t)} \odot \pi_{1:K+1}^{(t)}, \log(\mathbf{D}^{(t)}) \odot \pi_{1:K+1}^{(t)}, \mathbf{v}_{\text{obs}}, \omega_{\text{obs}}, \hat{\mathbf{v}}_{1:K}, \hat{\mathbf{x}}_{1:K})$

final image prediction: $\mathbf{I}'^{(t+1)} \leftarrow \mathbf{I}'_{\text{warp}}^{(t+1)}, \mathbf{I}'_{\text{imagine}}^{(t+1)}, \mathbf{W}_{\text{Warp}}$

update parameters: $\theta \leftarrow \theta - \gamma \nabla_{\theta} [|\mathbf{I}'^{(t+1)} - \mathbf{I}^{(t+1)}|^2 + \text{regularization loss}]$,
 $\theta \in [\theta_{\text{obj}}, \theta_{\text{depth}}, \theta_{\text{imag}}]$

3 Network training and dataset

We trained the three networks jointly using ADAM optimization [6] with a learning rate of $3e-4$, $\epsilon = 1e-6$ and other default setting in PyTorch, with a batch size of 20. 55 epochs were trained on the dataset. We set $\lambda = 1.0$.

The model was implemented in PyTorch and trained on NVidia RTX 6000 for about 8 days. We will release the dataset upon publication of the manuscript.

In consideration that other models originally demonstrated on simple datasets with homogeneous colors may need to increase their size for our dataset, we tested MONet[1] and slot attention[7] with both the original network sizes and

increased sizes. For MONet, we tested the original number of channels ([32, 32, 64, 64]) for the hidden layers of encoder of the component VAE (IoU=0.12, ARI-fg=0.27), and increased them to [32, 64, 128, 256, 256] with one more layer while also increasing the channel number in its attention network from 64 to 128 (MONet-bigger as reported in the main text: IoU=0.20, ARI-fg=0.36). The decoder layers' sizes were increased accordingly. For slot attention, we tested a variant which increased the number of features in the attention component from 64 to 128 (slot-attention-128, which took 15 days to train) in addition to the original setting (slot-attention) and found the performance to be similar. For our model, we tuned sizes of networks based on performance on the training set.

Images in the dataset are rendered in a custom Unity environment at a resolution of 512×512 and downsampled to 128×128 resolution for training. Images are rendered in sequence of 7 time steps for each scene. In each scene, a room with newly selected textures and objects is created. The starting location and orientation of the camera and objects are initialized randomly. The camera moves with random steps and pivots with random angles between consecutive frames; each object moves with random constant velocity drawn independently. All possible valid sequential triplets out of the recorded 7 frames (e.g., frames 1,2,3; 1,3,5; 1,4,7, etc.) form our training samples, excluding the frame sets where objects collide with each other or the camera. For our dataset we use 90k scenes for training which give us 275k valid triplets.

4 Details of prediction methods

4.1 Prediction of new location and pose of an object

We first explain how to use the inferred 3D spatial states of an object in two frame to estimate its velocity and predict its location in the next frame. Following the notation in 2.1 of the main text, in order to calculate the instantaneous velocity of object k at $t - 1$ from the perspective of the camera at t (with a static reference frame), we need to first calculate the 3D location that object k observed at $t - 1$ should appear at t to the camera if it is static relative to the background: $\hat{x}_k^{t-1 \rightarrow t} = \mathbf{M}_{-\omega_{\text{obs}}}^{(t-1)}(\hat{\mathbf{x}}_k^{(t-1)} - \mathbf{v}_{\text{obs}}^{(t-1)})$, where $\hat{\mathbf{x}}_k^{(t-1)}$ is the inferred 3D location of object k at $t - 1$; $\mathbf{v}_{\text{obs}}^{(t-1)}$ is the velocity of the camera from $t - 1$ to t viewed from a static reference frame centered at its original location; and $\mathbf{M}_{-\omega_{\text{obs}}}^{(t-1)}$ is the rotation matrix
$$\begin{bmatrix} \cos \omega_{\text{obs}}^{(t-1)} & \sin \omega_{\text{obs}}^{(t-1)} & 0 \\ -\sin \omega_{\text{obs}}^{(t-1)} & \cos \omega_{\text{obs}}^{(t-1)} & 0 \\ 0 & 0 & 1 \end{bmatrix}$$
 induced by the rotational speed of the camera from $t - 1$ to t . Thus, the spatial offset from any object j at $t - 1$ to object k at t , as viewed from the camera, is $\hat{\mathbf{v}}_{j \rightarrow k}^{(t)} = \hat{\mathbf{x}}_j^{(t)} - \hat{\mathbf{x}}_j^{t-1 \rightarrow t}$. We calculate the velocity of object k as the weighted average of the spatial offset against all candidate objects detected at $t - 1$ using the matching scores r_{kj} between object k and those candidate objects: $\hat{\mathbf{v}}_k^{(t)} = \sum_{j=0}^K r_{kj} \hat{\mathbf{v}}_{j \rightarrow k}^{(t)}$

The spatial location where the object k observed at t would be at $t + 1$ can be

114 predicted as $\mathbf{x}'_k^{(t+1)} = \mathbf{M}_{-\omega_{\text{obs}}}^{(t)}(\hat{\mathbf{x}}_k^{(t)} + \hat{\mathbf{v}}_k^{(t)} - \mathbf{v}_{\text{obs}}^{(t)})$, with $\mathbf{v}_{\text{obs}}^{(t)}$ and $\mathbf{M}_{-\omega_{\text{obs}}}^{(t)}$ defined
 115 similarly as above, but for camera movement between t and $t + 1$.

116 As the matrix multiplication involved in this prediction is almost impossible
 117 to be hard-wired in the brain, we also tested the possibility of learning to
 118 approximate this computation with a fully-connected neural network, as is
 119 further discussed in section 5 of this Supplementary.

120 As described in the main text, we use $\mathbf{p}_{\phi_k}^{(t)} = [p_{\phi_{k,1}}^{(t)}, p_{\phi_{k,2}}^{(t)}, \dots, p_{\phi_{k,b}}^{(t)}]$ to repre-
 121 sent the probabilities that the pose of object k at t falls into b equally spaced
 122 bins of angles in $(0, 2\pi)$, where $\sum_i p_{\phi_{k,i}}^{(t)} = 1$. With this notation, if we con-
 123 sider b possible discrete rotational speeds $-\frac{(b-2)\pi}{b}, -\frac{(b-4)\pi}{b}, \dots, 0, \dots, \frac{b\pi}{b}$, the
 124 probability that the rotational speed $\omega_{j \rightarrow k}^{(t)}$ is equal to any γ_1 among these
 125 discrete values if the same object is assigned index j at $t - 1$ and k at t is
 126 $p(\omega_{j \rightarrow k}^{(t)} = \gamma_1) \propto \sum_{l,m} p_{\phi_{j,l}}^{(t-1)} p_{\phi_{k,m}}^{(t-1)} \mathbb{1}(\frac{(m-l)2\pi}{b} \in \{\gamma_1 - 2\pi, \gamma_1, \gamma_1 + 2\pi\})$. We then
 127 calculate the likelihood of rotational speed for object k similarly by weighting
 128 the likelihoods calculated for all possible candidate pairs j, k with their matching
 129 scores.

130 The probability that the pose of object k is equal to any possible discrete
 131 yaw angle bin γ_2 at $t + 1$ can be predicted through $p'(\phi_k^{(t+1)} + \omega_{\text{obs}}^{(t)} = \gamma_2) =$

$$132 \sum_{\gamma_1, \omega} p(\omega_{j \rightarrow k}^{(t)} = \omega) p(\phi_k^{(t)} = \gamma_1), \text{ where } \omega_{\text{obs}}^{(t)} \text{ is the angular velocity}$$

$$133 \gamma_2 - \gamma_1 \in \{\omega - 2\pi, \omega, \omega + 2\pi\}$$

134 of the observer. Since we model $\phi_k^{(t)}$ and $\omega_k^{(t)}$ as probability distributions
 135 on discrete bins of angle while ω_{obs} can take continuous values, we convert
 136 $p'(\phi_k^{(t+1)} + \omega_{\text{obs}})$ to $p'(\phi_k^{(t+1)})$ on the same set of bins as $\phi_k^{(t)}$ by interpolation
 137 with a Von Mises kernel.

138 For the purpose of obtaining a point estimation of the rotational speed of
 139 ω_k to be used in predicting pixel-wise optical flow, we multiply the estimated
 140 probability for each angle bin with a vector pointing to the direction of that angle
 141 and take the angle of the sum of all these products (resembling the population
 142 vector code commonly observed in the brain for angular variables such as arm
 reaching direction or visual motion direction[8]).

143 4.2 Camera intrinsic

144 The Depth Perception network uses visual features in the image $\mathbf{I}^{(t)}$ to infer the
 145 distance of all pixels to the camera (depth) $\mathbf{D}^{(t)} \in \mathbb{R}^{w \times h} = h_{\theta}(\mathbf{I}^{(t)})$. With the
 146 inferred depth $D^{(t)}(i, j)$ for a pixel at any coordinate (i, j) in the image and the
 147 focal length f of the camera, the pixel's 3D location as viewed by the camera
 148 can be determined as $\hat{\mathbf{m}}_{(i,j)}^{(t)} = \frac{D^{(t)}(i,j)}{\sqrt{i^2 + j^2 + f^2}} \cdot [i, f, j]$. Here, we take the pixel
 149 coordinate of the center of an image as $(0, 0)$ and pixel indices are assumed to
 150 be symmetric around the center of the image: $|i| \leq \frac{w-1}{2}, |j| \leq \frac{h-1}{2}$. Conversely,
 151 given a 3D coordinate location x, y, z relative to the camera, the pixel coordinate
 152 it appears on the image can be calculated as well: $[i, j] = \frac{f}{y}[x, z]$. And its depth
 153 $d = \sqrt{x^2 + y^2 + z^2}$. This is the knowledge of camera intrinsic, which we assume

is known. However, since these are relatively simple mapping functions, we think it is possible to learn to approximate them with neural networks together with other networks in future works.

4.3 Weights in warping-based prediction and combination of two streams of predictions

The details of the procedure of predicting where a pixel at t would land at $t + 1$ on the image has been detailed in the main text. Here we describe the calculation of the weights when drawing the colors for every pixel at $t + 1$.

As the object attribution of each pixel is not known but is inferred by $f_{\text{obj}}(\mathbf{I}^{(t)})$, it is represented for every pixel as a probability of belonging to each object and the background $\pi_k^{(t)}$, $k = 1, 2, \dots, K + 1$. Therefore, the predicted motion of each pixel should be described as a probability distribution over $K + 1$ discrete target locations $p(\mathbf{m}_{(i,j)}^{(t+1)}) = \sum_{k=1}^{K+1} \pi_{kij}^{(t)} \cdot \delta(\mathbf{m}_{k,(i,j)}^{(t+1)})$, i.e., pixel (i, j) has a probability of $\pi_{kij}^{(t)}$ to move to location $\mathbf{m}_{k,(i,j)}^{(t+1)}$ at the next time point, for $k = 1, 2, \dots, K + 1$. With such probabilistic prediction of pixel movement for all visible pixel $(i, j)^{(t)}$, we can partially predict the colors of the next image at the pixel grids where some original pixels from the current view will land nearby by weighting their contribution:

$$\mathbf{I}_{\text{Warp}}^{(t+1)}(p, q) = \begin{cases} \frac{\sum_{k,i,j} w_k(i, j, p, q) \mathbf{I}^{(t)}(i, j)}{\sum_{k,i,j} w_k(i, j, p, q)}, & \text{if } \sum_{k,i,j} w_k(i, j, p, q) > 0 \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

We define the weight of the contribution from any source pixel (i, j) to a target pixel (p, q) as

$$w_k(i, j, p, q) = \pi_{kij}^{(t)} \cdot e^{-\beta \cdot D'_k{}^{(t+1)}(i, j)} \cdot \max\{1 - |i'_k{}^{(t+1)} - p|, 0\} \cdot \max\{1 - |j'_k{}^{(t+1)} - q|, 0\} \quad (2)$$

Here, $(i'_k{}^{(t+1)}, j'_k{}^{(t+1)})$ is the pixel coordinate corresponding to the predicted new 3D location $\mathbf{m}_{k,(i,j)}^{(t+1)}$ for pixel (i, j) from the current frame. The first term $\pi_{kij}^{(t)}$ in the definition incorporates the uncertainty of which object a pixel belongs to. The second term $e^{-\beta \cdot D'_k{}^{(t+1)}(i, j)}$ resolves the issue of occlusion when multiple pixels are predicted to move close to the same pixel grid by down-weighting the pixels predicted to land farther from the camera. These last two terms mean that only the source pixels predicted to land within a square of of 2×2 pixels centered at any target location (p, q) will contribute to the color $\mathbf{I}_{\text{Warp}}^{(t+1)}(p, q)$. The depth map $\mathbf{D}'_{\text{Warp}}{}^{(t+1)}$ can be predicted by the same weighting scheme after replacing $\mathbf{I}^{(t)}(i, j)$ with each predicted depth $D'_k{}^{(t+1)}(i, j)$ assuming the pixel belongs to object k .

Note that the formula above guarantees that the weights from all contributing original pixels sum to 1. However, because we assume that only a portion $\pi_{kij}^{(t)}$ of a pixel might land nearby a target pixel, the sum of all the portions of the

pixels that would land in the 2×2 square near a target pixel may be less than 1, in which case, the above prediction is combined with the prediction based on imagination to form the final prediction for such target pixels. Therefore, the weight on warping-based prediction is calculated as $\mathbf{W}_{\text{warp}}(p, q) = \min\{1, \pi_{kij}^{(t)} \cdot \max\{1 - |i_k^{(t+1)} - p|, 0\} \cdot \max\{1 - |j_k^{(t+1)} - q|, 0\}\}$.

5 Jointly learning rules of rigid body motion and apparent movement of objects induced by ego-motion

As mentioned above, the expected apparent change of 3D location of a static object induced by the motion of the camera itself, and the expected movement of all pixels on an object given the object’s 3D location and velocity and pixels’ locations (rigid-body motion) are assumed to be known in our main model. However, it is unnatural to assume these are hard-wired in the brain. To test whether they can be jointly learned together with all other networks by optimizing for the same objective, we approximated them with neural networks.

We used a series of 3 fully-connected layers, each with 16 units, with a residual connection from input to output, to learn the function that maps from an object’s location $\hat{\mathbf{x}}_k^{(t-1)}$, the camera’s horizontal movement velocity $\mathbf{v}_{\text{obs}}^{(t-1)}$ and rotational speed ω_{obs} to the 3D location where the object would appear after the camera’s motion: $\hat{\mathbf{x}}_k^{t-1 \rightarrow t}$. We then used two convolutional networks, each with channel sizes [64,64,32,4] and a residual connection from input to output layer, to approximate the prediction of pixels’ 3D motion. One network is dedicated for the background and the other for objects. The network for predicting pixels belonging to objects takes as input a pixel’s inferred 3D location, the location of the object it belongs to, the camera’s motion parameters, the object’s inferred 3D location and its velocity and rotational speed, and then predicts the new 3D location of the pixel in the next frame. The network for predicting background pixels’ motion only takes the pixels’ inferred 3D locations and camera’s motion parameters as input.

By jointly optimizing all networks for the same objective function, we found that the relaxed model obtains a segmentation performance of ARI-fg=0.58, IoU=0.44, which is on par with the full model (ARI-fg=0.58, IoU=0.45) (Figure S5). Further we found that the objects’ 3D locations are also still inferred well when the objects are segmented well (Figure S6). The correlation between ground truth polar angle and inferred polar angle for objects segmented with IoU>0.5 (mainly red dots) is 0.79 (while the correlation of all objects regardless of whether they can be segmented correctly by the model is 0.62). The correlation between the ground truth distance and the inferred distance for objects with IoU>0.5 is 0.48 (while the correlation without such filtering is -0.01). In regard to depth inference, we find that this network struggles to learn the background depth, which is reflected in the decrease in the overall correlation to the ground truth

229 depth ($r = 0.33$)

230 As discussed in the main text, this suggests that the assumption of rigid body
 231 motion may not be strictly necessary for segmentation and 3D localization of
 232 objects together with depth perception. Note that we did make an assumption
 233 that objects move with inertia (keeping their speed) as is true in the data. Future
 234 works may test whether this assumption still allows learning all the capacities
 235 when objects actually change speed (e.g, due to collision or gravity).

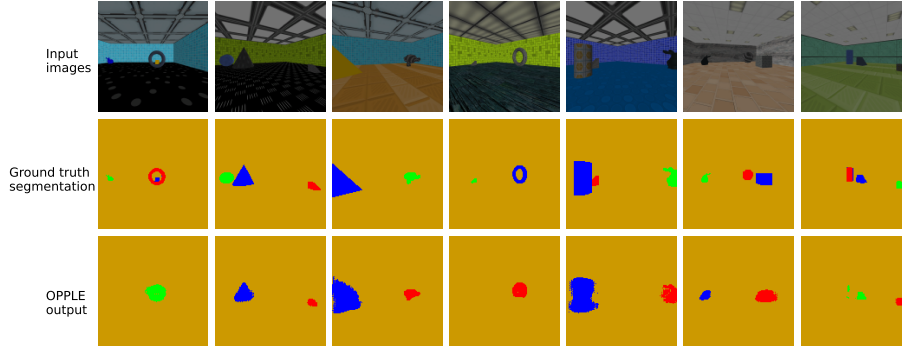


Figure S5: Segmentation performance for relaxed model (the rule of rigid body motion and apparent object motion induced by egomotion are learned).

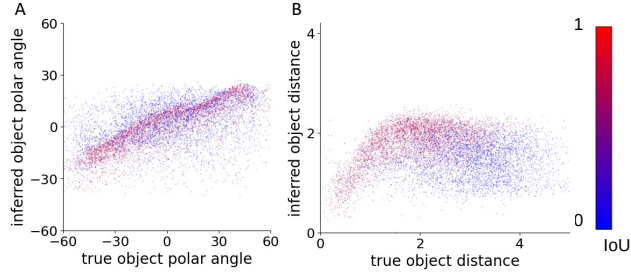


Figure S6: Performance of object 3D localization color-coded by segmentation performance for all objects in the relaxed model.

236 6 Performance of OPple on MOVi dataset

237 The recently published MOVi dataset [4] poses more challenges to OCRL models.
 238 It introduces more complex and realistic object motion following gravity and
 239 including interaction between objects. We therefore also test OPple on the
 240 most challenging version of MOVi dataset (MOVi-E) tested in a previous work
 241 SAVi++ [2]. We generated scenes including 3 objects following the default data
 242 generation procedure. In order to incorporate 3D object rotation, we made
 243 modification to the object extraction network such that it outputs cosines and

244 sines of each Euler angle characterizing the current pose (α, β, γ) of each object
 245 viewed by the camera. These values are used to compute a rotation matrix
 246 corresponding to the pose. Then, using the rule of apparent object movement
 247 induced by camera motion that incorporates 3D rotation, the pose of any object
 248 at $t - 1$ that would appear at t if it did not move is calculated. This allows
 249 the calculation of the rotation matrix for any pair of detected objects between
 250 $t - 1$ and t . Using the same weighting introduced in Section 4.3, the expected
 251 rotation for each object from t to $t + 1$ is calculated with an assumption that
 252 the rotation from $t - 1$ to t carries over.

253 On testing data of MOVi-E, OPPLÉ achieved a slightly lower segmentation
 254 performance (ARF-fg=0.38 and IoU=0.44) than on our custom dataset. Figure
 255 S7 illustrates the segmentation performance on MOVi-E dataset. It can be seen
 256 although overall the union of all segmentation masks often cover the majority
 257 of the true object masks, the model often attribute different objects into the
 258 same mask. We think the lower performance in segmentation may be due to the
 259 difficulty of learning a to assign a pose to each object relative to a canonical
 260 pose that is consistent to all objects. Additionally, because all objects follow
 261 gravitational force and fall at a similar time, their motions are often synchronized
 262 and the speed may thus appear close to each other, which encourages our network
 263 to attribute them to the same objects.

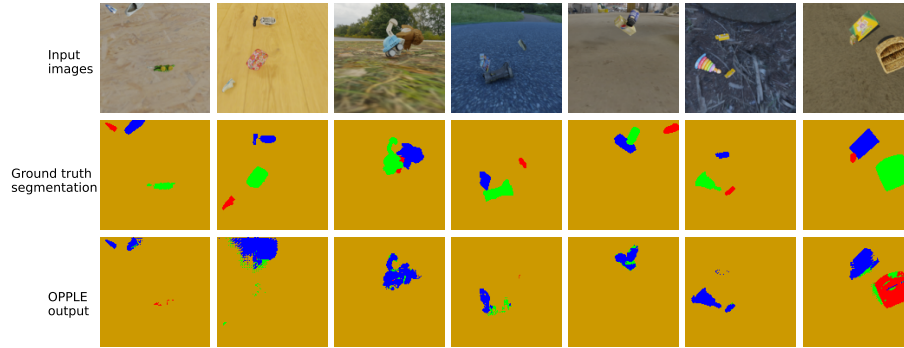


Figure S7: OPPLÉ's segmentation performance on 3-object version of MOVi-E dataset.

264 **7 Depth perception performance of O3V model** 265 **on our dataset**

266 As shown below, the O3V model generates blurrier depth maps for our dataset
 267 and often fails to capture the relative distance relationship between objects and
 268 background, even though it has the advantage of using 3 consecutive frames to
 269 infer 3D scene structure. In contrast, our model can make inference of depth
 270 based on a single image and produce more accurate depth estimation (Figure
 271 S1).

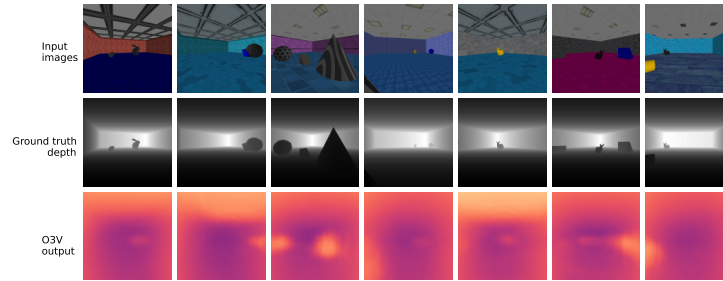


Figure S8: Example frames (first row), the ground truth depth map (second row) and the depth map inferred by O3V (third row).

References

- [1] Christopher P Burgess, Loic Matthey, Nicholas Watters, Rishabh Kabra, Irina Higgins, Matt Botvinick, and Alexander Lerchner. Monet: Unsupervised scene decomposition and representation. *arXiv preprint arXiv:1901.11390*, 2019.
- [2] Gamaleldin F Elsayed, Aravindh Mahendran, Sjoerd van Steenkiste, Klaus Greff, Michael C Mozer, and Thomas Kipf. Savi++: Towards end-to-end object-centric learning from real-world videos. *arXiv preprint arXiv:2206.07764*, 2022.
- [3] SM Ali Eslami, Danilo Jimenez Rezende, Frederic Besse, Fabio Viola, Ari S Morcos, Marta Garnelo, Avraham Ruderman, Andrei A Rusu, Ivo Danihelka, Karol Gregor, et al. Neural scene representation and rendering. *Science*, 360(6394):1204–1210, 2018.
- [4] Klaus Greff, Francois Belletti, Lucas Beyer, Carl Doersch, Yilun Du, Daniel Duckworth, David J Fleet, Dan Gnanapragasam, Florian Golemo, Charles Herrmann, et al. Kubric: A scalable dataset generator. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3749–3761, 2022.
- [5] Paul Henderson and Christoph H Lampert. Unsupervised object-centric video generation and decomposition in 3d. *Advances in Neural Information Processing Systems*, 33:3106–3117, 2020.
- [6] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [7] Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg Heigold, Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with slot attention. *arXiv preprint arXiv:2006.15055*, 2020.
- [8] Alexandre Pouget, Peter Dayan, and Richard Zemel. Information processing with population codes. *Nature Reviews Neuroscience*, 1(2):125–132, 2000.