

Self-Supervised Object Detection from Egocentric Videos Supplementary Material

Peri Akiva^{1,2} Jing Huang² Kevin J Liang² Rama Kovvuri² Xingyu Chen²
Matt Feiszli² Kristin Dana¹ Tal Hassner²
¹Rutgers University ²Meta AI

1. Implementation details

For the patch feature extractor and patch matching network, we use a ResNet-50 (random initialization) and a deformable transformer encoder to generate multi-scale patch features (see Figure 1 for example of patch indexing notation). For the object residual module, the number of clusters is $K = 16$ and feature dimension $D = 64$. For the random affine transformation \mathbb{T} , we sample random rotation and translation parameters from the ranges $(-10, 10)$ degrees and $(-15, 15)$ pixels respectively. For training, we use a batch size of 16, learning rate of 0.0001, momentum of 0.9, weight decay of 0.05, and number of negative samples of 200 patches per anchor patch. We first train the patch matching network for 8 epochs (everything else is frozen), after which the patch matching network is frozen (everything else is un-frozen) and used to predict multi-temporal patch matches for the main task objective. The patch feature extractor uses sine positional encoding, multi-scale self and deformable attention modules with 8 heads and 4 scales. The patch feature extractor is trained for 100 epochs or until convergence. During inference, we use $\eta = 50$, $\gamma = 0.25$, $\beta = 10$, and $\psi = 0.05$.

1.1. Pseudo-code

The implementation for the object residual module can be found in Pseudo Code 1, and the implementation for the patch matching module can be found in Pseudo Code 2.

2. Additional Qualitative Results

Figure 2 presents additional qualitative examples of viewing angle and illumination condition variation on the pre-smoothed cluster maps. Fig. 3, 4, and 5 present additional qualitative results for the Ego4D [1], EgoObjects [3], and COCO [2] validation sets. Figure 6 and Figure 7 present qualitative results of challenging cases and batch-wise smoothing outputs.

References

- [1] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18995–19012, 2022. 1
- [2] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 1
- [3] Lorenzo Pellegrini, Chenchen Zhu, Fanyi Xiao, Zhicheng Yan, Antonio Carta, Matthias De Lange, Vincenzo Lomonaco, Roshan Sumbaly, Pau Rodriguez, and David Vazquez. 3rd continual learning workshop challenge on egocentric category and instance level object understanding. *arXiv preprint arXiv:2212.06833*, 2022. 1

$\mathbf{z}_{s=0}^\tau$				$\mathbf{z}_{s=0}^{\tau'}$				$\mathbf{z}_{s=1}^{\tau'}$	
$z_{0,0}^\tau$	$z_{1,0}^\tau$	$z_{2,0}^\tau$	$z_{3,0}^\tau$	$z_{0,0}^{\tau'}$	$z_{1,0}^{\tau'}$	$z_{2,0}^{\tau'}$	$z_{3,0}^{\tau'}$	$z_{0,1}^{\tau'}$	$z_{1,1}^{\tau'}$
$z_{4,0}^\tau$	$z_{5,0}^\tau$	$z_{6,0}^\tau$	$z_{7,0}^\tau$	$z_{4,0}^{\tau'}$	$z_{5,0}^{\tau'}$	$z_{6,0}^{\tau'}$	$z_{7,0}^{\tau'}$		
$z_{8,0}^\tau$	$z_{9,0}^\tau$	$z_{10,0}^\tau$	$z_{11,0}^\tau$	$z_{8,0}^{\tau'}$	$z_{9,0}^{\tau'}$	$z_{10,0}^{\tau'}$	$z_{11,0}^{\tau'}$	$z_{2,1}^{\tau'}$	$z_{3,1}^{\tau'}$
$z_{12,0}^\tau$	$z_{13,0}^\tau$	$z_{14,0}^\tau$	$z_{15,0}^\tau$	$z_{12,0}^{\tau'}$	$z_{13,0}^{\tau'}$	$z_{14,0}^{\tau'}$	$z_{15,0}^{\tau'}$		
Anchor : $z_{2,1}^\tau$				Positive : $P_{2,1}^\tau = \{(6, 0), (0, 1)\}$					
				Negative : $N_{2,1}^\tau = \{(8, 0), (12, 0), (15, 0), (3, 1)\}$					

Figure 1: **Patch indexing notation example.** Anchor $z_{2,1}^\tau$ is matched with 2 positive patches highlighted in green, $z_{6,0}^{\tau'}$ and $z_{0,1}^{\tau'}$, and 4 negative patches highlighted in red, $z_{8,0}^{\tau'}$, $z_{12,0}^{\tau'}$, $z_{15,0}^{\tau'}$, and $z_{3,1}^{\tau'}$.

```

1 class ObjectResidualModule(nn.Module):
2     def __init__(self, D: int, K: int) -> None:
3         """Object Residual Module.
4         Args:
5             D (int): number of channels to encode
6             K (int): number of clusters
7         """
8         self.D, self.K = D, K
9         self.clusters = nn.Parameter(K, D, requires_grad=True)
10        self.scale = nn.Parameter(K, requires_grad=True)
11
12    def normalize(self, x):
13        """find weights of cluster centers using l2 distance
14        Args:
15            x (torch.Tensor): input features
16        Returns:
17            torch.Tensor: cluster-wise weights
18        """
19        num_clusters, in_channels = self.clusters.size()
20        batch_size = x.size(0)
21        scale = self.scale.view((1, 1, K))
22
23        x = x.expand((batch_size, x.size(1), K, D))
24        clusters = self.clusters.view((1, 1, K, D))
25
26        norm = scale*(x - clusters).pow(2).sum(dim=3)
27        return norm
28
29    def accumulate(self, theta, x):
30        """accumulation of weighted residuals
31        Args:
32            theta (torch.Tensor): learned cluster weights
33            x (torch.Tensor): input features
34        Returns:
35            torch.Tensor: accumulated residuals
36        """
37        clusters = self.clusters.view((1, 1, self.K, self.D))
38        batch_size = x.size(0)
39
40        x = x.expand((batch_size, x.size(1), self.K, self.D))
41        out = (theta*(x - clusters)).sum(dim=1)
42        return out
43
44    def forward(self, x):
45        """forward function for Surface Encoder with
46        Args:
47            x (torch.Tensor): input features
48        Returns:
49            torch.Tensor: residual features
50        """
51        theta = F.softmax(self.l2_norm(x, self.clusters, self.scale), dim=2)
52        residuals = self.accumulate(theta, x, self.clusters)
53        return residuals

```

Listing 1: Object Residual Module Pseudo-Code

```

1
2 # x: input image
3
4 def patch_matching_train(x):
5
6     rand_affine_params = affine_transformation_param_generator()
7     x_transformed = affine_transformation(x, rand_affine_params)
8
9     z1 = patch_matching_network(x)
10    z2 = patch_matching_network(x_transformed)
11
12    z1_transformed = affine_transformation(z1, rand_affine_params)
13
14
15    # Both sets of patches are spatially aligned
16    z1 = rearrange_as_patches(z1)
17    z2 = rearrange_as_patches(z2)
18
19    loss = InfoNCE(z1, z2)
20
21    return loss

```

Listing 2: Patch Matching Pseudo-Code

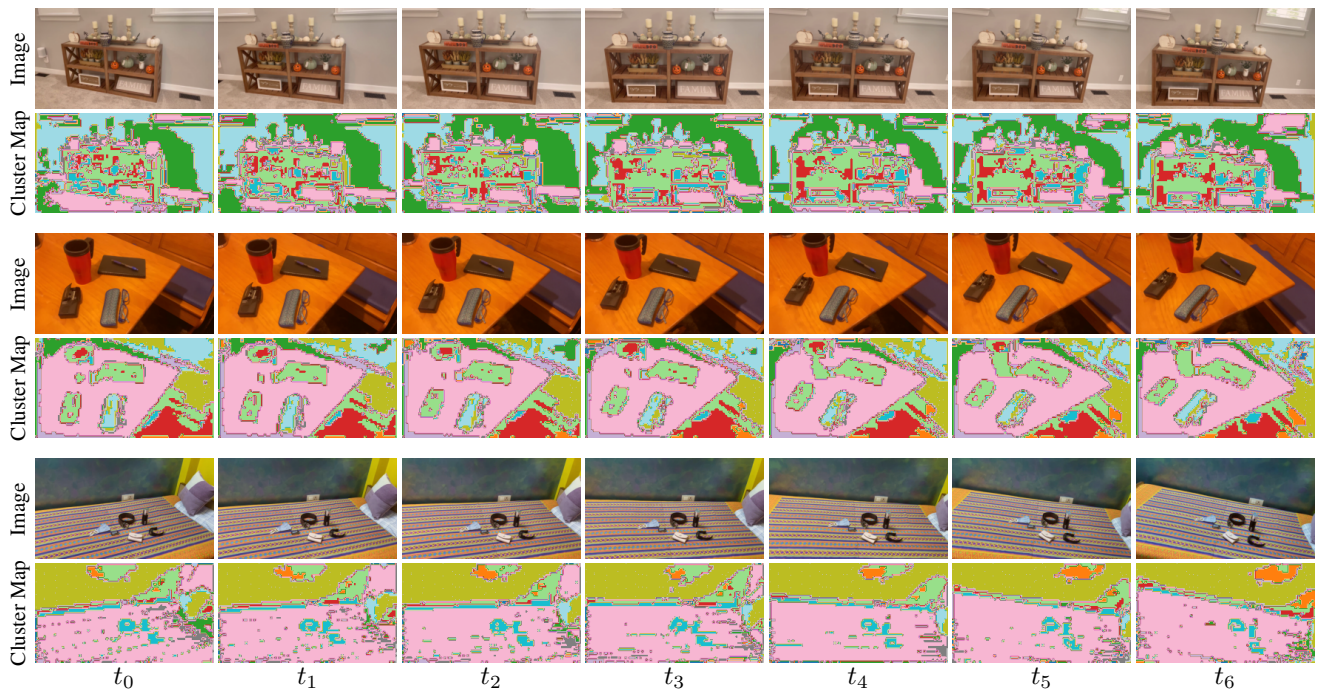


Figure 2: Additional qualitative examples of viewing angle and illumination conditions variation on the pre-smoothed cluster maps. Best viewed in color and zoomed.

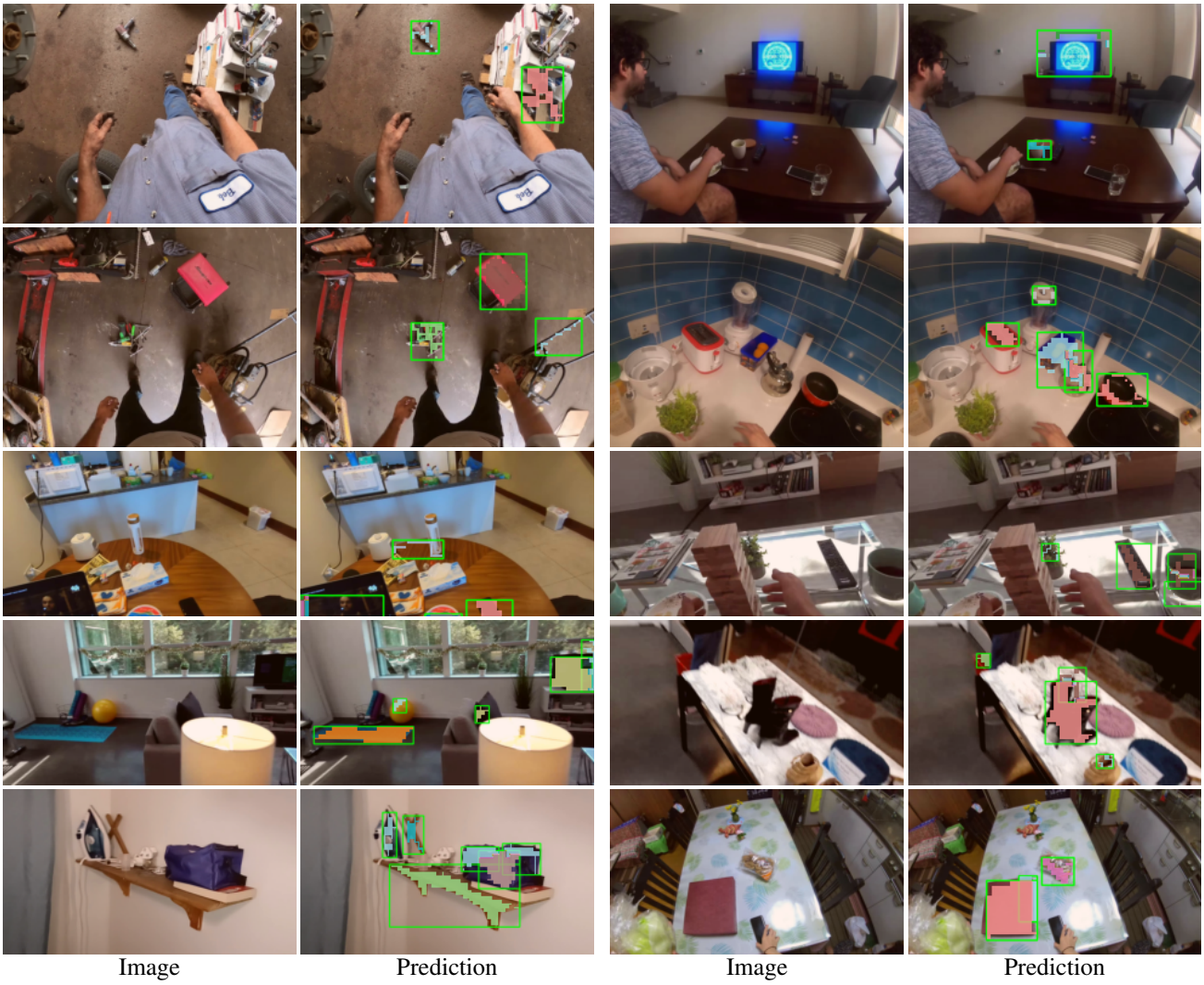


Figure 3: Additional qualitative results on Ego4D validation sets.

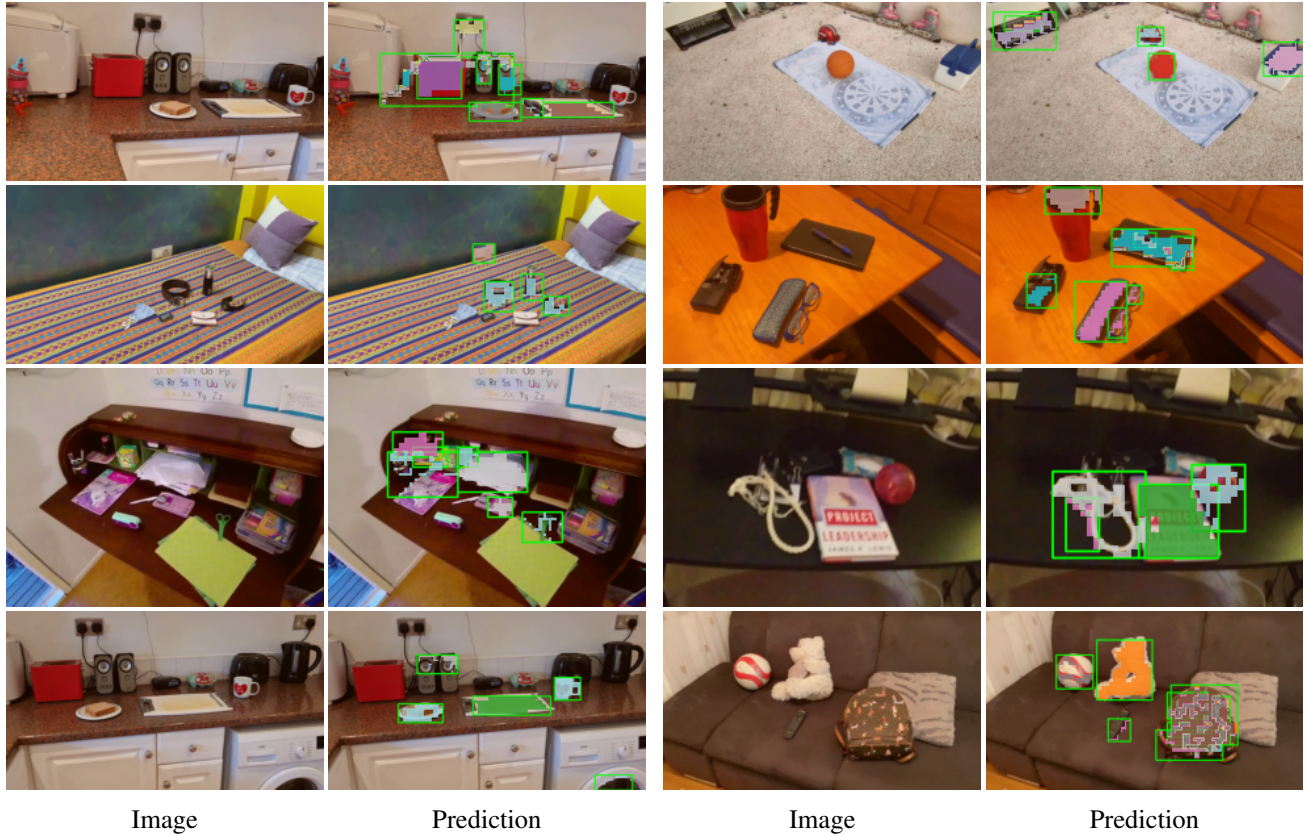


Figure 4: Additional qualitative results on EgoObjects validation sets.



Figure 5: Additional qualitative results on the COCO validation sets.

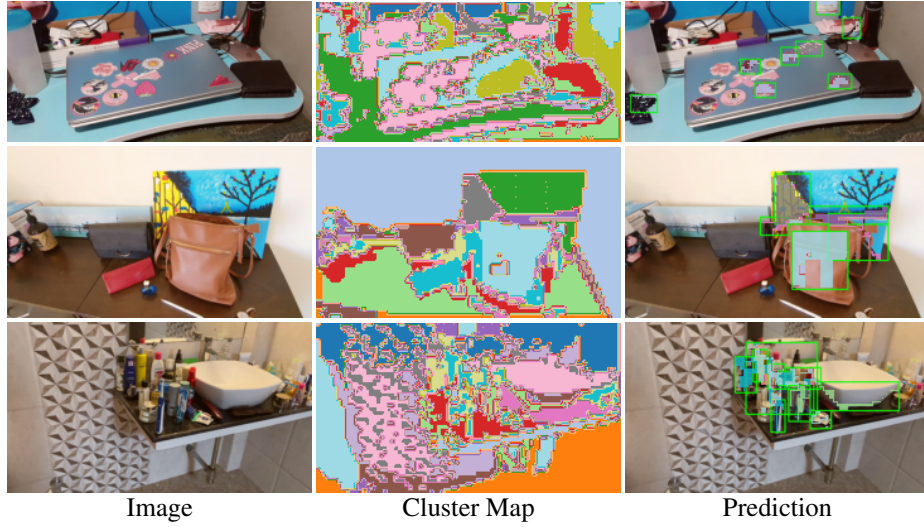


Figure 6: Qualitative examples of the challenging examples.

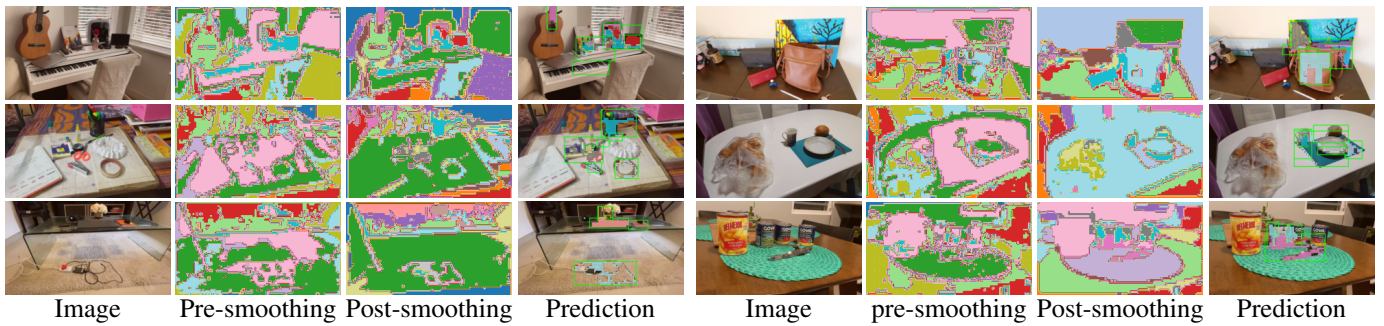


Figure 7: Qualitative examples of the smoothing operation.