

One-sample Guided Object Representation Disassembling (Paper ID 2171)

1 We sincerely thank all reviewers for the constructive comments. We provide short responses here due to the page limit.

2 **To Reviewer #1**

3 **Q1: a) How does the representation length influence disassembling? b) Is there a common set of parameters?**

4 **A1:** Smaller length yields better representation disassembling, while larger length leads to better image reconstruction.
5 For images with complex scenes, the representation length should be set larger. b) Yes. They are given in Line 199.

6 **Q2: Will the method still work on multiple object image? Q3: What is the exact goal of the dual-swap step?**

7 **A2:** In this case, features of multiple objects will be disassembled into the same part. Visual results are shown in Fig. I.

8 **A3:** The dual-swap step is devised for generating self-supervised information for the unannotated samples.

9 **Q4: a) Number of supervised samples in DSD. b) Why are the results of MoNet [6] and IODINE [12] different?**

10 **A4:** a) Around 5,000 samples augmented with annotated one-samples are used. b) MoNet [6] and IODINE [12] learn
11 object representations and corresponding masks. The fixed masks lead to the non-swappable defects of [6,12].

12 **Q5: a) Why choosing two-layer MLP as classifier? b) Does the classifier affect the disassembling effect?**

13 **A5:** a) It is a simple classifier and trainable by gradient descent. b) Yes. Simpler but not too simple classifiers better
14 reflect the properties of the representation, which favors the disassembling.

15 **To Reviewer #2**

16 **Q1: Extension to multi-object setting. Q2: Disentangling features of scene objects seems a goal of the method.**

17 **A1:** Thanks for the comment. The method is indeed able to handle images with multiple objects, in which case features
18 of multiple objects will be extracted together and denoted by the same part of the representation. Fig. I below shows the
19 visual results of multi-object samples in COCO and miniImageNet (a newly added dataset). More results will be added
20 to the revision. **A2:** In fact, disentangling features of objects is not our goal. Rather, our method aims to disassemble
21 features between one-sample guided objects and the background. Features of guided objects are denoted by one part of
22 the representation, while features of unguided objects, which are treated as the background, are denoted using the
23 remaining part of the representation. We will clarify this in the revision.

24 **Q3: Writing issues: unclarity of the introduction and conclusion; non-standard or unintuitive terminologies.**

25 **A3:** Thanks for the nice suggestions. We will clearly list the main goals in introduction and amend terminologies.

26 **Q4: Improving related works. Q5: Meaning of prior methods not working on complicated backgrounds.**

27 **A4:** Thanks. We will revise this part and provide clearer distinctions between our method and prior ones. **A5:** We
meant that, prior methods like MoNet [6] work on synthetic images only but not real-world ones. We will clarify this.



28 Figure I: Visual results reconstructed with original or object-swapped representations on multi-object images.

29 **To Reviewer #3**

30 **Q1: How robust is the method to the parameters? Q2: Will the method still work for multi-object images?**

31 **A1:** Please refer to A1(b) of R#4. **A2:** Yes. Fig. I shows the results on multi-object images.

32 **Q3: Why choosing only ten categories from COCO? Q4: The blurred results in Fig. 3. Other architectures?**

33 **A3:** The sample numbers in these ten categories are balanced. **A4:** Small representation length favors disassembling
34 but harms image reconstruction, leading to the blurred results. We tested the method with other architectures, yet results
35 showed that changing architecture barely helps. More labeled samples help improve the blurred results.

36 **Q5: Connection with 'fuzzy logic'? Q6: Image size? Q7: How to find the best representation length?**

37 **A5:** No. **A6:** SVHN: 32×32 ; Others: 64×64 . **A7:** Setting larger representation length for a more complicated image.

38 **To Reviewer #4**

39 **Q1: Are samples of COCO cropped to patches and resized? Q2: The sensitivity of hyper-parameters.**

40 **A1:** Thanks. We only resized samples to 32×32 (SVHN) or 64×64 (others), but did not crop them. **A2:** The vital
41 parameters (β, ρ) and δ respectively control the disassembling on annotated and unannotated samples. The (AMS,
42 AIS) for setting (β, ρ, δ) to be (1, 100, 5), (10, 1000, 5), and (10, 1000, 50) are (14.32, 7.83), (13.69, 6.51), and
43 (13.52, 6.31), respectively. Reasonably large β and ρ favor disassembling. More results will be added to the revision.

44 **Q3: Showing results on miniImageNet. Q4: Classification equations are expressed in an unfamiliar way.**

45 **A3:** Fig. I(b) shows visual results; the AMS and AIS scores for (S-AE, DSD, MoNet, IODINE, Ours) are respectively
46 (14.45, 13.36, 10.43, 17.49, 6.98) and (7.81, 5.76, 9.82, 18.93, 4.03). **Q4:** Our loss requires computing each value in
47 the predicted class vector. To be consistent, therefore, equations are written in an expanded form. We will clarify this.

48 **Q5: Explaining why supervised methods (S-AE) are worse than semi-supervised methods (One-GORD).**

49 **A5:** Thanks. With only labels and a two-layer MLP classifier, S-AE is not able to effectively extract the entire
50 discriminant features. However, dual-swap and object reconstruction of our method favor representation disassembling.

51 **Q6: Try ResNet as the decoder of MoNet or IODINE to deal with complex backgrounds. Q7: λ_{cla}^d ? and typos.**

52 **A6:** Thanks. Still, the modified MoNet and IODINE do not work. One reason is that the Gaussian-distribution
53 assumption facilitates decomposition but limits the representation's capacity. **A7:** It should be λ_{rec}^d . We will fix typos.