

1 We sincerely thank all reviewers and ACs for their time and efforts. We also appreciate the positive comments on our
2 novelty, contributions and state-of-the-art performance, *e.g.*, “impressive,” “meaningful”, “detailed ablation study”
3 and “well written and easy to understand”. We have released code in GitHub anonymously, which can be searched by
4 querying for “UWSOD”. Below we give our responses point-by-point.

5 **R1Q1 & R1Q2: More failure cases and qualitative results. The performance for the smaller objects in COCO.**

6 **A:** We will analyze some failure cases and qualitative results in the supplementary material as suggested. The *mAP* of
7 the small object in MS COCO is only 2.2%, as it is more challenge to detect small object than the large one in WSOD.

8 **R1Q3 & R3Q2: The performance of the approach if they use Selective Search (SS), Edge Boxes (EB) or MCG.**

9 **A:** Replacing SSOPG with traditional SS and EB proposals decreases 1.5% and 1.1% *mAP* on Pascal VOC 2007. And
10 MCG proposals improve the detection performance of UWSOD by 1.5% *mAP*. However, MCG takes about 34.3 s. per
11 image to generate proposals, which is time-consuming for practical applications. We will add the above results to paper.

12 **R2Q1: The proposed method is quite complex and its presentation is in some points difficult to follow.**

13 **A:** Thanks for your suggestions on presentation. We will revise the paper thoroughly under your suggestions.

14 **R2Q2: These approaches seem to add a relevant computational cost. The authors should comment on that.**

15 **A:** Our implemented WSDDN on VGG16 takes 439 ms. per images in GTX 1080TI GPU and PASCAL VOC. When
16 we sequentially add SWBBFT, SSOPG and MRRP, the training times increase to 618, 848 and 1, 581 ms., respectively.

17 **R2Q3: How the hyper-parameters are chosen without a validation set?**

18 **A:** Most hyper-parameters are set based on the previous WSOD and FSOD methods. And the hyper-parameters of our
19 components are manually chosen by intuition and experience, which are kept the same for different dataset.

20 **R3Q1: How to set T_i^{obn} , T_i^p , \bar{B}_i^{obn} , \bar{B}_i^p , and λ^{obn} , λ^p ? Are these two hyperparameters sensitive to performance?**

21 **A:** T_i^{obn} , T_i^p , \bar{B}_i^{obn} , \bar{B}_i^p are the classification and regression targets for objectness detection branch and object proposal
22 generator, respectively. They are set based on IoU threshold between proposals and the corresponding pseudo-ground-
23 truths. As we stated in the Implementation details, we set the labelling threshold λ^{obn} and λ^p to 0.5 and 0.7, respectively.
24 And varying λ^{obn} and λ^p within a range of $[-0.1, 0.1]$ only decreases the performance by about $1.3 \sim 2.1\%$ *mAP*.

25 **R4Q1: All three components considered by the authors have already been explored by existing works[11,17,72].**

26 **A:** As we emphasized in the Introduction, there are significant differences between our method and the existing works.
27 In particular, WSOD methods in [11,17] relied on external object proposal algorithms or additional instance-level
28 supervision to learn detectors, while work in [72] focused on encoding traditional image descriptions in fully-supervised
29 learning. To our best knowledge, our work is the *first* to propose learnable object proposals (SSOPG) without external
30 modules or additional supervision and aggregate multi-scale in-network contextual information (MRRP) for WSOD
31 task. And SWBBFT improves instance refinement from a new perspective of the trade-off between precision and recall
32 requirements in different branches. Our method has both practical and methodological contributions to facilitate the
33 development of this area. Moreover, detaching detectors from external time-consume modules makes WSOD more
34 capable of handling thousands of real-world categories and taking advantage of large-scale weak annotations.

35 **R4Q2: Why [17-20] cannot obtain such huge performance gain by using similar network modules (SWBBFT)?**

36 **A:** Actually, careful designs of instance refinement module have shown significant improvement of performance in
37 many works. The instance refinement proposed in OICR [8] improves WSDDN by 6.4% *mAP*. Recent methods
38 in [17,18,20,35] proposed various strategy to improve the instance refinement module, which achieved gains of
39 5.4%, 7.0%, 8.9% and 6.0% *mAP*, respectively, compared to vanilla instance refinement in OICR. We attribute the
40 improvement from SWBBFT to a sequence of effective refinement branches of increasing quality of pseudo-ground-
41 truths and a well balance between positive and negative samples, which are neglected in the previous works.

42 **R4Q3: It is not clear without meaningful proposals, how to train WSDDN, SWBBFT and SSOPG?**

43 **A:** We further exhibit how our UWSOD works in the learning process. First, WSDDN formulates WSOD as multiple
44 instance learning and captures the target object from a large set of proposals. Therefore, we enforce SSOPG to
45 output redundancy object proposals to ensure high recall, which makes WSDDN capable of selecting positive samples.
46 Although the output proposals of initial SSOPG are messy, WSDDN still has high-probability of finding informative
47 proposals, which either contain discriminative object part or cover the entire object loosely. Second, with the output of
48 WSDDN, a sequence of effective refinements in SWBBFT bootstraps the quality of predicted bounding boxes, which
49 has been demonstrated to improve performance in many existing works [8,17-20,35-36]. Third, a self-supervised
50 learning is leveraged to train SSOPG with supervision distilled from SWBBFT, which in turn improves WSDDN results.