1 We thank all reviewers and AC for the time and constructive comments. Below we address the main concerns, and will
2 fix other issues/typos and cite suggested references in the revised paper. Code will be released upon acceptance.

3 **R2 / R3 / R4: Sensitivity analysis on** $\gamma$**:** Evaluated on VOC07, CASD with $\gamma = 0.05, 0.075, 0.1, 0.15, 0.2$ gets mAP
4 $54.1\%, 54.7\%, 55.3\%, 55.0\%, 55.0\%$ respectively, demonstrating CASD is robust to $\gamma$.

5 **R1: Advantages of attention-based method:** In Line 41-55 and Fig 1, we motivate our attention-based method by
6 visualizing the features activated by WSOD networks. We observe that failures in WSOD are closely associated with the
7 flawed features. Non-attention-based WSOD methods typically improve pseudo-labels and loss functions, instead of the
8 features. In contrast, our method explicitly uses the structural info in attention maps, and promotes discriminative and
9 consistent features on whole objects (Fig 1). Our method is orthogonal and applicable to non-attention-based WSOD.

10 **Visualization of layer-wise attention maps:** Fig. 3 (b) visualizes the attention maps of 2nd/3rd/4th conv blocks.

11 **R2: Data splits:** We follow the data splits in [7,8,10,11,13,34,42]: For VOC 2007 and 2012, we train on the train+val
12 set and evaluate on the test set. For COCO, we train on the train set and evaluate on the val set.

13 **Clarification on hyper-parameter selection:** We agree that over-tuning hyper-parameters does not improve WSOD
14 research. However, as long as hyper-parameters are database-independent, certain choices of networks, parameters
15 and settings represent reasonable assumptions required to solve the ill-posed WSOD problem. Examples of the above
16 hyper-parameters include those in selective search, object proposal, score aggregation, NMS, scaling and flipping.
17 CASD inherits above hyper-parameters from PCL and OICR and fixes them on VOC07, VOC12 and COCO.

18 For the remaining CASD-specific hyper-parameters (Inverted Attention, choice of layers, max operation, loss weight $\gamma$),
19 we follow the routine of [7,8,10,11,13,34,42], and only tune CASD parameters on VOC07 with the VOC07 test set, and
20 then keep them fixed for VOC12 and COCO. The state-of-the-art results on VOC12 and COCO demonstrate that the
21 VOC07 hyper-parameters of CASD transfer well and are practical for new databases.

22 We share the reviewer's concern that some WSOL methods in [A] lack clear validation. However, our work and
23 [7,8,10,11,13,34,42] all select hyper-parameters on VOC07 and fix them for VOC12 and COCO. The VOC07 results,
24 although over-tuned on the test set, are still valuable for comparing performance "upper bounds" among WSOD
25 methods. The VOC12 and COCO results are NOT tuned on their test sets.

26 **Performance bounded by the proposal method:** We agree that using RPN in WSOD should beat Selective Search
27 (SS). However, CASD is a training module independent of proposal methods. We adopt the same SS as most WSOD
28 literatures only for fair comparisons. We will add both RPN-based CASD and Faster-RCNN results in the revised paper.

29 **R3: IW-CASD without classification score aggregation:** Without score aggregation and IA, IW-CASD achieves
30 $51.4\%$ mAP, clearly improving the baseline by $2.5\%$ mAP. Adding score aggregation brings $1.2\%$ mAP improvement.

31 **Clarification on** $k$ **and** $\alpha$**:** The vanilla OICR [10] has 3 OICR branches ($k = 3$), and suggests the larger $k$ the better
32 results. We only use $k = 2$ in all our experiments due to our GPU limitations. We reimplemented the OICR loss in
33 Pytorch by ourselves. On VOC07, the OICR baseline at $\alpha = 0.1$ achieves $48.9\%$ which is slightly superior to $48.3\%$
34 using the adaptive weighting policy in vanilla OICR. Thus we keep the fixed weight for the OICR loss.

35 **Clarification on method: (1).** For $L_{IW}$ and $L_{LW}$, the gradients are not back-propagated to the comprehensive
36 attention maps $A_r^{IW}$ and $A_r^{LW}$. **(2).** In Eq. (4), $A^{IW}$ should be $A_r^{IW}$. **(3).** Following Fast-RCNN, we define
37 $L_{reg}^k = \sum_{r=1,...,G^k} smooth_{L1}(t_r, \hat{t}_r)/G^k$, $G^k$ is the total number of positive proposals in $k$-th branch. $t_r$ and $\hat{t}_r$ are
38 tuples including location offsets and sizes of $r$-th predicted and gt bounding-box. **(4).** We will adopt the notation of
39 $A_r = max(A_r^{IW}, A_r^{LW})$ in revised paper. Thanks for the suggestion.

40 **R4:** We sincerely thank the detailed feedback. We will fix the typos, inconsistency, add details and improve writing.
41 **Different** $\gamma$**s for** $L_{IW}$ **and** $L_{LW}$**:** We conduct the following ablation study on VOC07. Fixing $\gamma_{IW} = 0.1$, CASD
42 gets $55.0\%, 55.5\%, 55.3\%, 54.8\%, 54.8\%$ when $\gamma_{LW} = 0.05, 0.075, 0.1, 0.15, 0.2$ respectively. Fixing $\gamma_{LW} = 0.1$,
43 CASD gets $54.3\%, 54.6\%, 55.3\%, 55.1\%, 54.9\%$ when $\gamma_{IW} = 0.05, 0.075, 0.1, 0.15, 0.2$ respectively. At $L_{IW} = 0.1$
44 and $L_{LW} = 0.075$, CASD achieves $55.5\%$ which is only $0.2\%$ better than $55.3\%$ (the best performance of single $\gamma$).
45 Thus we conclude that single $\gamma$ is a good trade-off between performance and hyper-parameter tuning.

46 **Evidence for consistency and completeness:** First, the attention maps in Fig. 1 b) compare the results of baseline
47 OICR and CASD, qualitatively demonstrating CASD gets more consistent and complete object features. More examples
48 will be added in the revised paper. Second, on the horizontal flipped VOC07 test set, CASD achieves $56.5\%$ mAP
49 which is similar to $56.8\%$ of the unflipped test set. This indicates that CASD is consistent w.r.t flipping.

50 **Result with Grad-CAM:** The layer-wise CASD with Grad-CAM gets $52.0\%$ mAP on VOC07. This is slightly worse
51 than $52.6\%$ mAP of layer-wise CASD with channel-wise average pooling+sigmoid. The latter is computationally more
52 efficient and adopted in our paper.