1 We thank the reviewers for their insightful and constructive comments.

2 **[R1] Multi-task not novel.** We would like to clarify that our contribution is far beyond multi-task learning: a unified
3 and general framework with 3 mechanisms to take advantage of verification tasks to help regression methods. In fact,
4 the multi-task learning only contributes +0.4 mAP gains, while others contribute +1.5 mAP gains (see Tab. 4).

5 **[R1] several works using keypoints, e.g. CentripetalNet and RetinaFace.** First, our method is beyond using
6 corner/keypoint verification: it generally exploits various verification tasks to help regression methods, e.g. the
7 foreground verification task, which is not exploited by CentripetalNet and RetinaFace. Second, even considering
8 corner verification alone, our work is significantly different from CentripetalNet and RetinaFace: CentripetalNet is an
9 improvement of CornerNet with a better corner matching mechanism, which is purely verification-based. In contrast,
10 our method focuses on how to combine corner verification into regression methods to improve object localization.
11 RetinaFace requires explicit keypoint annotations and only exploits multi-task learning, while our method does not
12 require additional annotation and is far beyond pure multi-task learning. We will discuss these works in revision.

13 **[R2] Add more description about RepPoints; still predict n sample points?** Thanks for the suggestion and we will
14 rewrite it in the revision. Yes, the understanding is correct.

15 **[R2] Within-box verification and FCOS's centerness.** We conduct experiments to include the verification task of
16 FCOS's centerness. On a 40.5 AP baseline (+corner in Tab. 3 of this paper), an additional centerness branch achieves
17 40.6 mAP which has almost no gains over the baseline, while our additional foreground branch has +0.5 mAP gains
18 (41.0 mAP vs. 40.5 mAP). On a 41.0 mAP baseline (+corner+foreground in Tab. 3 of this paper), an additional
19 centerness branch achieves 41.0 mAP, which has no gains over the baseline. Actually, although FCOS's centerness and
20 our within-box verification share some similarity, they are designed for different roles: while FCOS's centerness aims
21 for adjusting object scores and weighting the regression loss, the within-box verification aims to complement the main
22 regression branch. We will add discussions in the revision.

23 **[R2] Explain why corner verification produces more accurate localization results.** As described in Line 32-36,
24 there may be two probable intuitive explanations: corner points represent the exact spatial extent of bounding box; each
25 feature in corner point verification is well aligned to the corresponding point.

26 **[R2] Why CornerNet outperforms RepPoints v2 in AP90.** It is difficult to rigorously analyze all the factors for the
27 performance gap in AP90 due to too many implementation differences. Nevertheless, we think there are probably two
28 main factors: 1) our backbone is different (Res50 v.s HG-104), and some works have shown that ResNet architectures
29 perform significantly worse than Hourglass ones in detecting corner points, e.g. 30.2 (R101) vs. 38.4 (HG-104) in
30 Tab. 4 of CornerNet; 2) our highest resolution (C3) is lower than that of CornerNet (C2). A possible direction towards
31 better AP90 is to use HG architectures and higher resolution, but it would significantly change the main branch of
32 RepPoints (built on ResNet-FPN-C3-C7 structure), and we will leave it as our future works.

33 **[R2] Why higher resolution does not yield better performance.** The corner heatmap is used in both feature-level
34 fusion and result-level joint inference. For feature-level fusion, we compared C3 heatmap with C3-C7 heatmaps, where
35 we find the latter performs better. We hypothesize it is because more positive samples by C3-C7 benefit the training. For
36 result-level joint inference, higher resolution does yield better performance that all levels using C3 heatmap performs
37 better than that using the corresponding heatmap level. We will add the discussion in revision.

38 **[R3] Terminologies.** Thanks for the suggestion. We will carefully examine and revise them.

39 **[R3] Ablation about hyperparameter $r$.** $r = 1, 2, 3, 4$ produce mAP of $41.0, 40.8, 40.5, 40.2$, respectively, indicating
40 $r = 1$ performs best. We will add this ablation in revision.

41 **[R4] Details of joint inference.** Intuitively, for a regressed corner point $p^t$, Eq. (2) means searching a point with the
42 highest corner heatmap score in a neighborhood of $r$ as its final location. Regressed points on all pyramidal levels use
43 the C3 corner heatmap for searching, and $r=1$ corresponds to a 8-pixel neighborhood. We will polish the description.

44 **[R4] Doubts about definitions in Tab.1.** We agree with the reviewer's comment that the relatively smaller displacement
45 regression targets in RetinaNet does not have benefit over RepPoints (or FCOS) which regress the box extent from the
46 center point. However, this has no contradiction with our categorization of RetinaNet as an verification+regression
47 framework and RepPoints as a regression framework: there are pre-defined anchors (coarse localization hypothesis) in
48 RetinaNet, while there is no such coarse hypothesis for RepPoints (anchor-free). Moreover, our Paragraph 2 delivers
49 the same information as the reviewer's comment, which motivates us to clarify the difference between two verification
50 methods (CornerNet vs. RetinaNet, see Paragraph 3) and to further propose our method (Paragraph 4).

51 **[R4] Computational costs and inference speeds.** FLOPs comparison is described in Line 265-271. For real inference
52 speed, the speed of RepPoints v1 is 12.7 FPS (img/s) using ResNet-50 on a Titan XP GPU, while that of RepPoints
53 v2 is 10.1 FPS. With a ResNeXt-101-DCN backbone, the speeds are 4.3 FPS v.s. 3.8 FPS for RepPoints v1 and v2,
54 respectively. We will add inference speed comparisons in the revision.