

1 **Evaluation.** Our initial experiments only used the first 500 images from the MS-COCO data set, as was done in [3].
 2 In the tables below we use the full test set with the best performing reduction based on the initial experiments. As
 3 **R3** and **R4** suggest, we will use the full test set to update the comparison Table 1 of Sec. 6 towards the final manuscript.

4 **Additional Experiments.** As **R3** suggested, we provide additional results for different architectures. While the
 5 median-smoothed RCNN has better clean performance as measured by F1 score, its certified precision recall is strictly
 6 lower highlighting the trade-off between clean accuracy and robustness as discussed in the literature.

	Precision	Recall	F1	Certified Precision	Certified Recall
YOLOv3	91.80%	20.88%	34.02%	29.65%	11.43%
Faster RCNN	86.58%	24.63%	38.35%	17.69%	10.85%
Mask RCNN	85.50%	25.14%	38.85%	17.42%	10.85%

7 **Other Performance Metrics.** As **R4** suggested, we examine the average precision (AP) for both the plain and
 8 certified detectors; we report the results for YOLOv3 below, and plan to include Faster RCNN and Mask RCNN in the
 9 final manuscript. Since varying the objectness threshold changes the base detector (f), we reevaluate the smoothed
 10 detector at objectness thresholds $\{0.1, 0.2, 0.4, 0.6, 0.8\}$ and calculate the area under the steps to lower bound true AP.
 11 As for inference speed, the smoothing paradigm is inherently costly, as we use 2000 perturbations for Monte Carlo
 12 estimation. We leave it to future work to improve these important metrics as needed in practice.

	AP@50	Certified AP@50
YOLOv3	32.0%	4.2%

13 **Comparison to Prior Work.** We thank **R4** for bringing relevant prior work to our attention. While we cannot claim
 14 to be the first adversarial defense for object detection, we maintain that we provide the first certified defense. The
 15 performance of the certified defense approach is currently so weak compared to the adversarial training approach that
 16 we do not think a meaningful quantitative comparison can be done. We will make sure to discuss the relation to this prior
 17 work as follows: "our certified radius is 0.36 in terms of ℓ_2 -norm whereas Zhang et al. (ICCV'19) achieved robustness
 18 radius of $8/255$ in terms of the stronger ℓ_∞ -norm threat model." As **R2** suggests, we will distinguish certification and
 19 defense as follows: "Several methods of obtaining robustness certificates for classification problems have been proposed
 20 [19, 21, XX]. In addition, [9, 11, 24, 25, 29, YY] proposes methods to both defend the model while enabling better
 21 certificates;" we will make sure to include the citations suggested by **R2**.

22 **Tailored Empirical Attacks.** As **R2** suggested, we implemented a DAG attack against our best performing model.
 23 The DAG attack is modified to include Monte Carlo sampling to increase the strength of the attack. We take 20 PGD
 24 steps and draw 5 random samples to estimate the gradient of the smoothed model. Surprisingly, the smoothed model
 25 is quite robust within the desired radius. The DAG attack was only able to decrease recall by 1.1%. This illustrates
 26 that the bound we obtained is likely quite loose with respect to the true robustness of the object detector, and we leave
 27 improvements of the robustness certificate as future work.

28 **Certifiable radius.** **R2** rightly points out that as of right now the certifiable radius is too low for real-world ap-
 29 plications. We emphasize that certified robustness is a challenging domain and none of the existing methods yield
 30 practical certificates even for classification problems. For example, while the SOTA certified defense by Salman et al.
 31 (NeurIPS'19) achieved 68.2% certified accuracy for $\|\epsilon\|_\infty < 2/255$ on CIFAR-10, the empirical approach of Xie et al.
 32 (CVPR'19) can achieve similar empirical robustness but at the ImageNet scale and with larger radius. That said, our
 33 work leverages a principled smoothing approach to provide the first non-trivial certificates for architectures as complex
 34 as object detectors.

35 **Generality of the Techniques.** **R1** remarks that the sorting and bucketing techniques proposed in Section 5 may be
 36 too specific for object detection. We note that the proposed techniques are potentially applicable to certifying other
 37 networks through reductions to a regression formulation. This is particularly relevant for tasks that have variable length
 38 outputs, such as key points detection, instance segmentation, or image captioning. Viewed in the broader context of
 39 adversarial robustness for ML models, computer vision continues to provide exemplar problems and we hope our work
 40 on object detection will help advance both the theory and practice of this important field.

41 **Other Clarifications and Corrections.** **R2** correctly points out an issue with Equation 3. We already fixed this issue
 42 in the supplemental materials, replacing min/max with inf/sup. We will also clarify the definition of the worst-case
 43 bounding box as "the box with coordinates satisfying the certified upper and lower bounds which realizes the lowest