

1 Thank you for the thorough reviews.

2 1. *New insights about detectors (all reviewers) and which detectors are better (reviewer 1)*. At least six new insights:

3 i) Existing object detectors largely fail to generalize well not because of a generic problem with performance but
4 because of specific phenomena (lack of invariance to background, viewpoint and rotation) that they are brittle to in a
5 way that humans are not. When we selected a subset of images in ObjectNet which had the same biases as those in
6 ImageNet, detector performance on ObjectNet reached that of humans, just as it is in ImageNet. This indicates that
7 focusing on structural methods to handle these phenomena in detectors can have significant payoff.

8 ii) We demonstrate that there is a large performance gap with respect to humans even over object classes that are
9 part of ImageNet. Human performance over ObjectNet is around 95% on a free-form text answer question, which is
10 considerably higher than networks' performance. This is in stark contrast to other work which indicates that object
11 detectors are approaching human-level performance.

12 iii) Our results indicate that data alone will not fix these problems. We systematically fine-tuned on ObjectNet using a
13 ResNet-152 trained on ImageNet. Fine-tuning was with 1-64 per class images in steps of powers of 2. Performance
14 increased by $\approx 2\%$ when doubling the data. This shows that even when networks learn biases that lurk in ObjectNet
15 they remain brittle with respect to the controls used.

16 iv) ObjectNet helps explain why a large gap in performance is seen between developers and users of object detectors in
17 a way that no other modern dataset does. On other datasets fine-tuning recovers the vast majority of the performance
18 very quickly, which is not the case with ObjectNet, and which is not observed, for example, in robotics applications.

19 v) Large-scale data can be collected in a way that is less biased, not just for object detection but for machine learning in
20 general. Bias is rarely addressed despite affecting all modern machine learning methods.

21 vi) In a sense, detectors are roughly the same. The gap between ObjectNet and ImageNet is not closing with newer
22 models. 1% gained on ImageNet is also 1% gained on ObjectNet. The community as a whole is making steady
23 measured progress and a breakthrough awaits.

24 2. *Occlusion and clutter, introduced biases due to people holding objects. (reviewer 1)* Without a doubt, ObjectNet
25 has its own biases which include the fact that certain object shapes are likely to be held when imaged from certain
26 angles. Had we removed more biases, model performance would have likely further decreased. Not all biases are the
27 same though: some are easy to learn and others are difficult. Our fine-tuning results show that learning about how
28 specific objects appear in ObjectNet does not help performance much — far less than the large performance jump seen
29 in other datasets. What we do show is the importance of three particular biases that humans are known to be extremely
30 resilient to from an early age, but networks fail to understand. These overwhelm other types of biases that might help
31 performance. The reviewer astutely points out that we do not consider occlusion or clutter. In the next version of the
32 dataset we intend to do so.

33 3. *Real-world biases are not necessarily bad, as it can reflect real-world distributions. These biases exist in the world,
34 and leveraging them is useful for task performance in many cases. (all reviewers)* Is it true that real-world biases can
35 be useful and they do increase average-case performance. At the same time, relying on them can be very dangerous.
36 Applications of computer vision to robotics, autonomous cars, surveillance, and other domains require high performance
37 in unusual situations. Indeed, that's where robustness is most needed because those are likely to be situations that
38 were not encountered during the development phase. Additionally, many applications of computer vision can have
39 adversarial opponents — humans that attempt to fool systems by defacing lane markings or hiding dangerous objects
40 from detectors. Most datasets that exist today, ObjectNet being an exception to this, focus on the usual case and do not
41 shed light on the performance one will encounter in unusual cases.

42 4. *Limitations (all reviewers), simulation and AR (reviewers 2 and 3)*. ObjectNet like all methods has many limitations:
43 it consists of indoor objects, the objects have to be available to many people, objects must be mobile, must not be too
44 large or small or fragile or dangerous. We cannot ask users to damage, paint, cut, or otherwise permanently change
45 objects.

46 In simulation one can create datasets that minimize bias although current methods produce images which are on the
47 whole easy for object detectors even when bias is removed. Foreground/background separation is easier in rendered
48 images and there are fewer sources of noise. In addition, there is a kind of variety and bias that is not easy to control in
49 simulation: shape. There are many possible shapes for a chair but in simulation one will have access to only a small
50 number of variations of static models. Methods exist to synthesize objects of different shapes for the same class but
51 these generally require specifying class boundaries or generative shape models — itself an open and difficult problem.

52 AR can be useful for collecting data and will play a role as we ask users to pose more complex scenes, like those which
53 include occlusion. We are excited to bring an awareness of more systematic experimental design methods, and show
54 how these can identify the most important phenomena for different machine learning tasks, can help characterize the
55 state of models and compare models in more fine-grained ways, and ultimately help develop new types of models.