1 We thank all reviewers for their valuable feedback. In particular, the reviewers have supported our:

2 **research and ideas** (**R4**-promising research direction & appealing, **R5**-novel, **R2**-fresh and convincing),

3 **experiments** (**R2**-thorough and well-conducted, **R5**-carefully designed) as well as

4 **presentation** (**R1**-clear, **R2**-well-written, **R5**-well written and easy to understand).

5 We address all of the concerns below. In particular, we provide *new comparisons* to past works, address the question of

6 performing this *research in simulation* and place this within the context of past works in *interactive perception*.

8 **R1**: **Segmentation and mass are not the end goal of interaction.** We totally agree. However, addressing more

9 complex tasks requires obtaining reasonable performance on these building blocks, which itself is very challenging in a

10 self-supervised setting. Complex downstream tasks including rich manipulations are part of our ongoing research.

11 **R1**: **Comparison with [2, 45].** We now provide results for a new baseline [45] - augmenting Mask-RCNN with their

12 Robust Set Loss using their public implementation. Across the suite of metrics, this helps a little beyond Table-1(i), but

13 still inferior to our approach, e.g. 22.3 (theirs) vs 28.06 (ours) on NovelSpaces Box $AP^{0.5}$, row (e). [2] addresses the

14 problem of poking an object from one point to another – very different from us, and not valid to compare to. However,

15 we now provide results for our model trained on their dataset (we do not poke, but use their pokes instead) and with

16 almost no hyperparameter tuning obtain 39.1 for Box $AP^{0.5}$ – validating that our method shows promise on real world

17 data. The higher number (39 vs 28) also indicates the relative complexity of the THOR scenes.

18 **R1**: **All in a simulated environment and simple scenes** Please see *R4 Research/Assumptions in simulation* below.

19 **R2**: **Rule-based decisions.** These were made to make the learning process more manageable. Future extensions of this

20 work will address more generic architectures and learning paradigms, to remove some these handcrafted designs.

21 **R2**: **Prioritized replay and curriculum learning.** Prioritized replay and curriculum learning encourage sampling

22 high-loss and low-loss examples, respectively. Combining these two approaches discourages sampling medium-loss

23 examples. This resulted in high gains for our method.

24 **R2**: **Fully supervised Mask-RCNN for novel objects.** It is trained with groundtruth mass and segmentation masks

25 (obtained from THOR) for seen objects (no category information is used). Novel categories are not in training images.

26 **R4**: **Research in simulation.** The end goal of this research is to train interactive agents in the real world. Using a

27 rich, large scale and variable simulator allows us to have faster research iterations via fast training (no mechanical

28 constraints), improve generalization (scene and object variability compared to [2, 45]) whilst ensuring safety. Future

29 work includes (a) deployment on a Locobot robot (b) exploring simulation-to-real transfer via real world fine tuning,

30 along the lines of navigation research in embodied AI (e.g., Habitat ICCV19, CARLA CORL17, iGibson ICRA20).

31 **R4**: **Assumptions in simulation.** Our setup presently makes some assumptions. However, past real world works also

32 use highly simplified setups. E.g. [45] assumes: (1) objects on a flat surface (2) fixed distance of surface to camera, (3)

33 fairly uniform background, (4) objects going out of the field of view after interaction, (5) fixed camera view, (6) no

34 exogenous motion. Our work relaxes (1)-(4) by providing variable objects, surfaces, rooms and backgrounds.

35 **R4**: **Interactive perception baselines.** One can cluster past works in this area into two parts: (a) passive methods that

36 use objects motion to segment them (b) active approaches that jointly learn to act and segment. Papers in **a** cannot be

37 thought of as baselines – they are orthogonal to our method. When our agent learns to interact with objects, the resulting

38 motion can be exploited by these works to result in better segmentation. Adding such techniques on top of our proposed

39 method is very interesting and we leave it for future work. Also, papers mentioned by R4 use *supervised* components.

40 Papers in **b** are very relevant, and we provide new comparisons above. Please see response to R1. However, our method

41 is designed to not just provide object masks but also object attributes – an improvement over past algorithms.

42 **R4**: **Additional supervision decreases the quality.** Without supervision the masks and interaction points are noisy.

43 We believe this serves as a form of data augmentation for our method and helps the model to better generalize.

44 **R4**: **Majority vote is > 33%.** The metric is mean per class accuracy. So if everything is predicted as the majority class,

45 the accuracy is 100% for that class and 0 for the other 2 classes. Hence, majority vote accuracy is (100+0+0)/3 = 33.3.

46 **R4**: **Intuition for mass prediction & how to use motion cues.** The model probably learns to associate the size and

47 texture of the objects to the mass. To incorporate motion cues, we can enable our training pipeline during inference. It

48 can predict the interaction points and infer the mass based on the motion caused by interaction.

49 **R4**: **Class imbalance & Objectness prior.** We circumvent the class imbalance issue by learning the objectness score

50 by focal loss only on pixels selected for interaction (line 238). This improves successful interactions from 4% to 60%

51 during training. We did not use objectness priors (e.g., [56,33]) since they use supervised data.

52 **R4**: **Mass prediction.** It is per object. We predict an interaction point for objects and a mass map for the entire image.

53 The mass for an object is the mass we predict for the pixel corresponding to the interaction point of that object.

54 **R5**: **Non-rigid objects & force directions.** Great challenging suggestions and will be considered for future work.

55 **R5**: **Forces.** L112-4 is about the point that the force is applied to. L184-5 is about the direction of the force vector,

56 which is not necessarily in the direction of the ray from the camera center to the point.