

1 We thank the reviewers for their feedback. We will add the clarification sentences requested. Below are our responses:

2 **R1: no ablation study.** We conducted the requested study. Our results show that accuracy is poor if we ablate the
3 attention network or the Saccader cell. We will add the analysis to Fig Supp.3 and state: "The attention network allows
4 the Saccader to better plan for locations to attend to as it has wider receptive field (RF). Note that this wider RF does
5 not affect the interpretability of the model as the classification path RF is still limited to 77x77. Furthermore, removing
6 the Saccader cell (i.e. using the BagNet-77-lowD ordered logits policy) yields poor results compared to the Saccader."

7 **R2: computational cost and/or parameter counts.** We will add a table with the parameters count of all models. In
8 summary, our model has 88,554,409 parameters, which is nearly twice the DRAM parameters.

9 **R2: how the final prediction is made.** The prediction is made using the averaged logits across the locations.

10 **R2: Does the "location network" ... refer to the attention network, the 1-by-1 conv and the Saccader cell? Yes.**

11 **R2 and R4: other datasets.** ImageNet is an extremely large and diverse dataset that contains both coarse- and
12 fine-grained class distinctions. To achieve good performance, a method must not only distinguish among superclasses,
13 but also e.g. among the >100 fine-grained classes of dogs. Moreover, most natural image datasets have some class
14 overlap with ImageNet; few datasets are entirely disjoint. Recent work has suggested that ImageNet is a representative
15 benchmark; Kornblith et al. CVPR 2019 (<https://arxiv.org/abs/1805.08974>) showed that accuracy on ImageNet predicts
16 accuracy on other natural image classification datasets.

17 **R4: This design would fail to apply on ... pedestrian detection... cancer classification.** We agree that it would be
18 nice to apply our method to pedestrian detection and cancer classification. However, we want to stress that natural image
19 classification is an important computer vision problem. Recent advances in computer vision started with classification
20 tasks on ImageNet (e.g. Krizhevsky et al., 2012) and then future research extended these methods to other domains. We
21 will add sentences to the results to better motivate the task.

22 **R4: NASNet ... unclear whether this improve comes from the increased model capacity or higher input resolu-**
23 **tion.** We conducted an experiment on ImageNet 224 and using the Saccader-NASNet model. Our results show that the
24 accuracy is better than the Saccader model alone but worse than the Saccader-NASNet model on the high resolution
25 ImageNet 331 (we added this experiment to Figure 6). This finding demonstrates that the accuracy benefits from both
26 the increased capacity as well as the higher input resolution.

27 **R4: pre-training ... This step introduces a strong bias.** We agree with the reviewer that pre-training introduces bias.
28 However, we find that this bias is helpful in getting a better final policy. In Figure Supp3, we show that reinforcement
29 learning after pretraining location network enhances accuracy compared to starting learning without this pretraining
30 step. Also, note that the final learned policy is different than the pretraining target policy (i.e., the ordered logits policy).
31 As we show in Figure 3 and 5a, the final learned policy performs much better. Just as SGD biases neural network
32 training toward solutions that generalize well, we find that the pretraining alters the training trajectory in a way that
33 produces a better-performing model.

34 **R4: baseline ... such as Class Activation Map (CAM).** In this work, we are concerned with models with hard visual
35 attention. CAM (Zhou et al. 2016) and similar interpretability methods try to provide an explanation of the model
36 decision in a way that relies on a heuristic (e.g., that the spatial localization of features should be preserved in the final
37 feature map) rather than explicitly constraining how the network processes its input. These methods are fragile (see
38 Hooker et al. 2019), and the model's final decision may nonetheless rely on information provided by features with
39 small weights (see Jain and Wallace, 2019). Models with hard attention take a different approach by using a controller
40 that selects parts of the input to be processed by the network, which provides interpretability by design. In our work,
41 the representation network may be regarded as a network to construct CAM, with a guarantee that the receptive field is
42 limited. We will add a citation to Zhou et al. 2016.

43 **R1: no weaknesses ... have been noted.** We will add: "Although Saccader outperforms other hard attention models, it
44 still lags behind state-of-the-art feedforward models in terms of accuracy. Future research may extend the Saccader
45 model to achieve even better classification performance while maintaining the interpretability of model decisions."

46 **R1: "what" and "where".** We will add: "These are analogous to the ventral ("what") and dorsal ("where") pathways
47 that are involved in object recognition and localization, respectively in human vision (Goodale and Milner 1992)."

48 **Other improvements.** We computed errorbars (mean \pm SD) for all plots. In the DRAM, we limited the high resolution
49 to classification and the (high, mid and low) resolutions to initialize the location LSTM state, which encouraged better
50 location exploration. We also extended the DRAM pretraining to two stages on wide and limited receptive fields (120
51 epochs each), and doubled the LSTM layers size. Despite these changes, the Saccader was still better than the DRAM.
52 We improved the ResNet model accuracy in Fig 4. We corrected the Sobel and Canny baselines plots (accuracy remains
53 poor). Since Canny and Sobel results are similar, we only included the Sobel results to improve the presentation.