

1 **General Reply**

2 We thank all reviewers for their constructive comments. Overall, the reviewers appreciated the importance of the  
3 problem, the practical utility of our results, the rigorous experiments, and the clarity of our writing. At the same time  
4 we recognize that the scores were split, with R3 advocating acceptance and R1, R2 leaning towards weak reject. We are  
5 optimistic that our rebuttal can address each reviewer’s individual concerns below:

6 **Reply to R1:**

7 Thanks for your thoughtful comments on how we can improve our draft.

8 **“The stud[y] is restricted to “natural” shifts, ie. no attacks.”** In fact we do run experiments in which shifts are due to  
9 adversarial attack, finding, interestingly that here two-sample tests performed in the ambient space are most effective.

10 **“the submission is similar to a review paper.”** While we agree that our work draws upon many other papers (like  
11 a review) would, our contribution here is a rigorous empirical study with results that stem from combining existing  
12 techniques in new ways. BBSD has never been applied beyond the label shift setting and its wider applicability (in  
13 practice) had never previously been demonstrated. Moreover, the general DR + 2-sample test pipeline has never been  
14 systematically investigated. We believe that this constitutes a source of fundamental scientific insights. In short, we  
15 wish to assert that the novelty of our paper is in the science, not the introduction of new models.

16 **“the main contribution of this work should be there,”** While we agree that identifying benign vs malign examples is  
17 of paramount importance, we hope to convince you that detection itself, which is better developed in our work is a  
18 sufficiently vital problem. We hope to make more progress on detecting malign examples in future work.

19 **“the author could [summarize] the way they select the most anomalous samples”** While we considered two ap-  
20 proaches, our promising results to date consist of selecting those examples that a domain classifier assigns with highest  
21 probability to the target domain.

22 **“Better interpretation of results”** We will discuss our intuitions and interpret the results more thoroughly if accepted.

23 **Details:** *a) Sec. 3.1:* We will clarify that all of our DR + two-sampling methods are non-conventional. The label  
24 classifier applied broadly as a shift detector is the most successful technique and indeed our contribution. However,  
25 domain classifier approaches have been explored previously. *b) L 139:* We chose the Gaussian kernel for MMD  
26 following the original paper. Per your suggestion, we will add ablations to show sensitivity to kernel bandwidth. *c) L*  
27 **165-166:** *c* and *r* refer to the number of columns and rows of the contingency table. *d) L 198:* “difference” classifier  
28 refers to the domain classifier (pardon the typo). *e) Tab. 1:* Detection accuracy measures how often, across a variety  
29 of shifts, a statistically-valid detector identifies the shift. *f) Fig. 3:* For this dataset, differences between source and  
30 target data arise from a different set of angles from which the items were photographed (source:  $0^\circ - 175^\circ$ , target:  
31  $180^\circ - 355^\circ$ ). *g) Typos and reference bugs:* Thanks for the attention to detail. We will make all appropriate changes  
32 in the camera-ready if accepted.

33 **Reply to R2:**

34 *a) L 196:* You are right. This statement entered the paper due to a miscommunication among collaborators. We  
35 apologize for the error and will strike it from the camera-ready version if accepted. You are also right that this  
36 heuristic is predicated on strong assumptions (addressed below). *b) Access to labeled anomalous samples is a strong*  
37 **assumption:** Our intuition here is that while labeling all examples at deployment would be prohibitive but that a small  
38 quality team might be used to selectively label those *top anomalous* examples. Such monitoring teams are increasingly  
39 common (Google’s search quality team is a famous example). *c) Related work and other baselines:* We point out  
40 that all standard two-sample tests in the input space constitute baselines as does the domain classifier approach. Notably  
41 (Tab. 1) MMD on inputs performs worst among all tested methods. *d) Distribution of p-values in experiments:* This  
42 is a great question. We will investigate the distribution of p-values and add these results to the paper. *e) Why is Tab. 1*  
43 **averaged over MNIST and CIFAR-10:** In this table, we average over all different perturbation models and datasets  
44 because our goals was to see if there exists a shift detection method that consistently outperforms the others. We  
45 will make this rationale clearer. *f) What sort of intervals are shown in the plots:* the intervals around the (mean)  
46 p-value are  $1-\sigma$  error bars. See L 222 - 225 and the caption of Fig. 2 for details. Time steps used for testing are  
47 described in L 229. *g) Weaknesses, underlying assumptions and failure cases w.r.t. shift characterization and*  
48 **malignancy detection:** The big assumption here is that *obviously* shifted examples are somehow more likely to be  
49 misclassified. This assumption might not hold. We will work to assess the correlation of domain classifier confidence  
50 with misclassification in real datasets and to formalize settings in which the assumption is reasonable.

51 **Reply to R3:**

52 Thanks for your strong positive review for our paper. We are glad that you appreciated the usefulness of the approach  
53 and the lack of pretense. We agree that further experiments with data from other domains and extensions to handle  
54 streaming data are promising next steps.