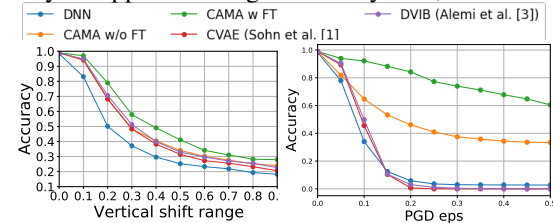


1 We thank the reviewers for their time and insightful comments on our paper. We respond to each reviewer(R) below. In
2 addition, we will make the code publicly available, together with the paper.

3 **R1. The proposed model is a form of conditional variational autoencoder (CVAE)** : We believe this is a misunderstanding,
4 and our contribution goes far beyond that. Our contributions are in three-fold: (1) a causal view on the robustness of
5 neural networks and the discussion on valid artificial manipulations; (2) the CAMA model as an instance of causal
6 consistent generative models, including versions for both single modality data and generic causal graphs; (3) fine-tuning
7 approaches to improving model robustness to unseen manipulations. R1’s comments cover only part of contribution (2),
8 i.e., the CAMA model for single modality data. R3 also pointed out "The proposed fine-tuning phase to learn unseen M
9 intervention is novel and sets this paper apart from previous work like Narayanaswamy et al." We discussed in detail
10 how our work differs from CVAE-type of models in appendix B, and will clarify on this further in revision.

11 **R1.decomposition necessity and CVAE comparisons**: We emphasise that CVAE and other mentioned work can
12 only be applied to single modality data, which is a special case discussed in our paper. For empirical comparisons,



13 we present on the left of below figure the robustness of CVAE
14 (Sohn et al. [1]) & DVIB (Alemi et al. [3]) to vertical shifts
15 and against PGD attacks, which clearly shows the advantage of
16 CAMA, especially the fine-tuning ability. Similar results have
17 been reported in [25] of our paper’s reference, and our work
18 provides extra advantages due to the use of a causally consistent
19 model and the fine-tuning method motivated by causal reasoning.

20 **R2. training data format... additional information**: The training procedure only requires knowledge on whether the
21 training data is clean (in such case we set $m = 0$) or manipulated (potentially with unknown manipulation). In the
22 manipulated case we do not require explicit label for M and perform inference on it instead. For test data, the labels for
23 both Y and M are not available, only the input X is given, similar to the standard test setting.

24 **R2. Comparisons to IRM**: First, IRM only considers single modality data. Moreover, IRM aims to find representations
25 that are invariant to environmental changes, which does not necessarily provide robustness to *unseen* manipulations.
26 Imagine training a model with IRM using clean MNIST and MNIST shifted 50% to the left, then the model is very
27 likely to fail when tested on MNIST shifted 50% to the right. Will clarify in revision.

28 **R3. (adversarial) data augmentation**: Deep CAMA also benefits from adversarial training (Figure 10). Both CAMA
29 and discriminative model are robust to *known* adversarial attacks that are observed in adversarial training. However, the
30 key advantage of CAMA over discriminative models is in its robustness to *unseen* manipulation. First, Figure 1 shows
31 that training discriminative DNNs with one manipulation can hurt its robustness to another unseen manipulation. This
32 is not the case for CAMA (see the overlapping green & orange curves in Figure 7(b)). Second, with fine-tuning (which
33 is another main contribution of our work), CAMA can be significantly more robust to unseen manipulations.

34 **R3. concerns on line 52 statement**: In causality literature, causal factorizations correspond to modular/independent
35 mechanisms [34,38] of our paper’s reference; building ML models consistent with these causal factorizations also
36 improves robustness [19] of our paper’s reference. Will clarify further by expanding the arguments.

37 **R3. Bayes’ rule prediction clarity**: As shown on the RHS of eq. (6), we sample m from $q(m|x)$ and sample z from
38 $q(z|x, y, m)$ with the previously sampled m . In the paper, we use 1 sample for m and K samples of z associated with
39 each m sample. Given the sampled m and z instances, we can compute the terms inside the log for each $y = c$ (as an
40 approximation to $p(x, y = c)$), and apply softmax to obtain the predictive distribution. Will clarify further on it.

41 **R3. parameter size**: The network size detail is presented in appendix D for all experimental settings. In appendix C
42 (Figure 16), we present the performance of discriminative models with larger network sizes and the result is similar to
43 the one with smaller network sizes.

44 **R4. computation cost**: The time complexity of CAMA is in the same order of a regular VAE. For predictions, fine-tuning
45 requires a small amount of additional time, as only a small fraction of data is needed for fine-tuning (Figures 8 & 20).

46 **R4. "Limitations of Adversarial Robustness"**: We do not intend to claim CAMA’s robustness to all possible manipu-
47 lations *straightaway after training*. So the theory of the mentioned work (and many others) is applicable to CAMA
48 before fine-tuning. However, this line of theoretical work evaluates the robustness of a classifier trained on clean data
49 only. Rather, CAMA brings extra advantages by fine-tuning, which enables adaptation to the unknown manipulation in
50 test time. This adaptation is required whenever a new manipulation is present, allowing CAMA to “learn” to be robust to
51 a growing number of manipulations.

52 **R4. Network architecture bias**: Indeed RCNN-like networks can be more robust to shifts by design. However, Fig. 15
53 & 16 shows that the use of CNN (instead of MLP) does not fully address the over-fitting issue to seen manipulations.
54 Moreover, architecture design typically provides robustness to specific manipulations. Instead, our goal is to make the
55 model robust to (infinity number of) unseen manipulations, which cannot be archived by architecture design only.

56 **R4. limit mis-specification**: In practice, the causal graphs are either provided by domain experts or obtained by running
57 causal discovery algorithms. For the former, it requires working closely with domain experts to refine the causal
58 hypotheses. For the latter, choosing a suitable causal discovery algorithm for the application at hand would be critical,
59 and there exist approaches to evaluating the performance of causal discovery algorithms.