

1 We thank the reviewers for their valuable comments.

2 **Reviewer #1:**

3 **Q:** Optimization of objective (22) with the proximal map (21) is not discussed in detail. MATLAB code is not available.

4 **A:** Thank you for bringing it to our attention. The MATLAB code has been uploaded to the GitHub repository. It  
5 simply invokes `fminunc` to optimize (22), without providing a gradient subroutine. This leaves `fminunc` to choose its  
6 own solver which typically utilizes its own finite difference routine. The result looks good and efficient for this dataset.

7 **Q:** the paper does not explore more straightforward applications to per-layer regularization in deep neural networks.  
8 Mentioning this seemingly fundamental application, the authors do not study it just saying that it is straightforward to  
9 implement. This makes the reader wonder whether this method can actually show promising results in this setting.

10 **A:** ProxNet for multi-view learning used proximal mapping for only one layer, and so did Figure 1. In contrast,  
11 ProxLSTM in Section 4 embeded a proximal mapping into *each* step/layer of an LSTM. It will be very interesting to  
12 study the use of multiple proximal mappings at different layers for diverse purposes, e.g., to enforce equivariance in  
13 each layer of feature extractor by using the violation as the regularier  $R$ , disentanglement at a certain layer, and fairness  
14 in prediction. This requires more space than a NeurIPS paper, and we are exploring some of these combinations.

15 **Reviewer #2:**

16 Thank you for pointing out Meinhardt et al. (2017). It is relevant and we will cite it in the future versions.

17 **Q:** Proximal operators may hurt the performance if they are not instantiated correctly for the problem and model that are  
18 being considered. I do not have much intuition on how the ProxNet improvements upon other regularization approaches.

19 **A:** ProxNet is a means of enforcing deep regularization, while the regularizer  $R$  itself is to be designed by the  
20 practitioners for the specific application, e.g., CCA for multi-view learning. So ProxNet by itself does not introduce any  
21 new regularizer. Implementing it is straightforward as shown in Eq (2) for any generic  $R$ , though some parameters  
22 need to be chosen (see the next question). The purpose of the paper is to show that for regularizers defined in terms of  
23 hidden layer outputs, it is more effective to regularize **in-place** through a proximal mapping at that layer, compared with  
24 adding the regularizer to the overall objective and relying on backpropagation for optimization. By “more effective”,  
25 we have compared by using the test performance instead of the training objective value, because unlike comparing two  
26 different nonconvex solvers, ProxNet results in a different objective than regularized risk minimization. A rigorous  
27 analysis beyond the intuition of in-place regularization (to quote Reviewer 3: “the kind of idea that seems obvious in  
28 retrospect”) will be interesting in the deep context, and we plan to investigate it in the future.

29 **Q:** Does care need to be taken to properly setting the  $\epsilon$  term to ensure the derivative approximation is well-behaved?  
30 Will the finite-difference derivative approximation in L178 causes instability during learning?

31 **A:** [16] provided several heuristics for setting  $\epsilon$ . Our experiment set  $\epsilon = \delta(1 + \|(X, Y)\|_\infty) \left\| \left( \frac{\partial L}{\partial P}, \frac{\partial L}{\partial Q} \right) \right\|_\infty^{-1}$  so as to be  
32 adaptive to the magnitude of the gradient, and  $\delta$  is a small constant whose value was chosen by checking a few gradients  
33 at the beginning of training. To estimate the “true” gradient needed for the check, we used a small  $\epsilon$  and solved the  
34 proximal mapping in Eq (10) to high accuracy. With this heuristic, the approximate gradient did not cause instability in  
35 training. It is also noteworthy that the stochastic gradient computed from a mini-batch introduces noise in the first place.

36 **Reviewer #3:**

37 **Q:** Would it be helpful to anneal the RRM term over time as the ProxNet term is annealed out?

38 **A:** We conducted some additional experiments to show  
39 how an annealed  $\lambda$  influences the performance of RRM.

40 We not only tried the annealing schedule used for Prox-  
41 Net, but also other schedules. The results for multi-view  
42 learning on the Sketchy dataset are shown in the right  
43 table, where, overall, annealing has mixed influence on  
44 the vanilla RRM, but it remains inferior to ProxNet.

#class	20	50	100	125
RRM new	17.4 $\pm$ 0.8	22.5 $\pm$ 0.5	24.3 $\pm$ 0.9	26.1 $\pm$ 0.7
RRM old	15.2 $\pm$ 0.6	20.1 $\pm$ 0.4	26.8 $\pm$ 0.5	28.1 $\pm$ 0.4
ProxNet	<b>13.7</b> $\pm$ 0.3	<b>17.9</b> $\pm$ 0.5	<b>20.2</b> $\pm$ 0.3	<b>22.0</b> $\pm$ 0.4

45 **Reviewer #4:**

46 **Q:** I don't see any aspect of the robustness method for RNNs that actually relies on the recurrent structure. Am I  
47 missing something, or could it be placed on any type of neural network hidden layer?

48 **A:** The ProxLSTM in Section 4 made RNN robust by applying a proximal mapping at each layer/step for invariantization.  
49 This is innately synergistic with the recurrent structure. In other words, we achieved robustness not because of using the  
50 recurrent structure itself, but by properly invariantizing each step via embedding a proximal mapping. So this technique  
51 can be generically deployed in any type of neural network.