

1 We thank all the reviewers for their insightful comments and suggestions. We will update the final version accordingly.

2 **R#1.** We thank the reviewer for pointing us to the two relevant references (we will discuss them in the final version),
3 that use low-rank representations to get robustness to specific empirical attacks like PGD. We remark that in our work
4 low rank representations are used in a different way. To provide the improved *certified* robustness guarantees, we take
5 advantage of good low-rank representations in extending the randomized smoothing approach appropriately.

6 (a) *On using non-linear projections:* We agree that introducing a non-linear dimensionality operation may lead to
7 classifiers that are more robust against PGD/FGSM style attacks, as observed in Sanyal et al. Our setting is quite
8 different since our aim is to produce classifiers with *certified* accuracy guarantees. Introducing a projection operation at
9 an intermediate layer completely breaks down the theoretical analysis of certified accuracy. Hence, we are restricted
10 to using linear projections in order to provide provable guarantees for our classifier. We believe that the reviewer’s
11 suggestion of using/analyzing non-linear projections is an excellent direction for future work.

12 (b) *On training the linear projection vs. using a fixed projection:* We agree that in certain cases such as text data where
13 the input representation is not fixed, training the linear projection along with the network could be beneficial and in
14 fact necessary. We are currently exploring this direction. For vision datasets that we used in the paper, we indeed had
15 experimental results with simultaneously training the projection with the network parameters. The results we obtained
16 were similar to using a fixed projection and we did not see any significant advantage. We chose to present the simpler
17 approach to convey the core idea clearly. We will include these results in the final version.

18 (c) *Choice of r :* We plot the PCA reconstruction error as a function of r and choose a value of r in a certain range where
19 the error is not too high (less than 3%). See Fig. 5. There are multiple choices of r that work equally well.

20 (d) *Training complexity:* The complexity of Algorithm 1 is comparable to the complexity of training a smoothed
21 classifier as in the work of [SYL⁺19]. The PCA step incurs a one time preprocessing cost and the projection step at
22 the beginning simply corresponds to adding a linear layer to an existing ResNet architecture. As an example, on the
23 CIFAR-10 dataset, for $\epsilon = 0.25$, training the classifier of [SYL⁺19] takes on average 21.27 seconds per epoch, whereas
24 Algorithm 1 takes 21.29 seconds per epoch on average. The same behavior holds across different parameter settings.

25 (e) *On trading off certified accuracy vs. natural accuracy:* Notice that for most values of ϵ , as Fig. 2 shows, we suffer
26 almost no loss in accuracy at small radii. Additionally, as we state in Line 141 (Page 4), for a large range of values for
27 the robustness radius, our method gets much better natural accuracy for a desired robust accuracy and radius. In practice,
28 we may not know the radius of adversarial perturbation (and the ideal choice of ϵ) beforehand, hence sacrificing a small
29 amount of accuracy at small radii for a significant gain at higher radii is a desirable tradeoff.

30 **R#2.** In Line 142, thanks for catching the typo: yes we meant blue instead of yellow. We now address the other points.

31 (a) *Scaling PCA to large datasets:* One can perform PCA on a smaller random sample to get the projection. Alternately,
32 one can also train the network and the projection operator simultaneously. Formally, reparameterize $\Pi = UU^T$ and
33 augment the loss function with two terms: 1) Reconstruction error on the **mini batch** namely, $\|x - Ux\|^2$; and 2)
34 $\|U^T U - I\|_F^2$ to encourage U to be orthonormal. Our experiments with this approach on CIFAR-10/100 show results
35 similar to those reported in the paper. Further, this approach naturally scales to large datasets.

36 (b) *Certified accuracy vs. r & hyperparameters:* Changing r by small amounts had negligible effect on the certified
37 accuracy (we tried values of r close to 200). The smaller we can keep r , the more robustness we can achieve. As
38 mentioned in our response to R#1 (1c), the reconstruction error (see Fig. 5) dictated our choice of $r = 200$. This choice
39 made apriori to the training stage works well across different settings, suggesting that choice of hyperparameters can be
40 decoupled. However, we did not run extensive experiments for end to end training for many different values of r . The
41 performance of our algorithm is smooth in λ (see Fig. 1). Furthermore, in all our experiments $\lambda = 0.5$ worked very well
42 pointing to the fact that there are generic settings of λ that one can often use. We will clarify more in the final version.

43 (c) *On feature dimensionality:* We do not know if linear projection forces the intermediate representations to be even
44 lower dimensional. Forcing intermediate representations to be low dimensional may certainly lead to empirical benefits.
45 However, it seems challenging to use this to get certified accuracy guarantees.

46 **R#3.** Please see responses to R#1 for the effect of PCA on natural accuracy (1a) and training cost (1d). Using
47 autoencoders is an interesting idea. Currently we do not know how to use them to get certified accuracy guarantees.

48 **R#4.** Regarding *new SDP algorithm* we respectfully disagree with the opinion that the approach is straightforward.
49 Please note two main contributions: (i) the algorithm is practical and novel within the broad MWU-framework (with an
50 important change to the weight update), (ii) the accompanied theoretical guarantee gives a significant improvement
51 over the current SOTA, for an even broader class of general quadratic programs (as discussed in Secs. 3, A, and D). As
52 experimentally shown in Sec. F, the algorithm is much faster than existing solvers, and may be of independent interest.

53 (a) *MWU vs DCT:* We are somewhat confused by the reviewer’s question about comparing MWU and DCT. These are
54 not two competing approaches. For certified ℓ_∞ guarantees in the DCT domain, these are used together. The MWU
55 based algorithm is used as a subroutine to find a good robust low-rank representation in the DCT domain (the MWU
56 based algorithm is crucial for obtaining certified guarantees).