

1 We thank all the reviewers for recognizing the contribution of our work and providing their valuable comments. We address the
 2 questions below and will release the trained model together with code as suggested by reviewers to contribute to the community.

3 **Shared comment on gradient-masking.** We first address this by conducting analysis following reviewers’ suggestions from
 4 four complementary perspectives: §1 black-box attack, §2 stronger white-box attack, §3 loss variations and §4 gradient-free
 5 attack. The results have verified that *the improved robustness is indeed due to model improvement instead of gradient masking.*

6 **§1 Results under black-box attacks.** For the black-box attacks, the results for
 7 PDG100, CW20 and CW100 in the original submission are incorrect due to our
 8 own fault: *they are results under white-box attacks.* We spot this error after sub-
 9 mission deadline but were unable to upload the correct version. *We sincerely*

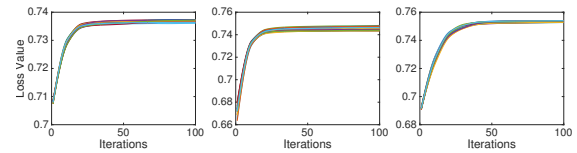
B-Attack	PDG20	PDG100	CW20	CW100
Undeferred	89.0	88.7	88.9	88.8
Siamese	81.6	81.0	80.3	79.8

10 *apologize for this misleading mistake and any additional efforts required from all the reviewers because of it.* The correct
 11 black-box results are shown in the table above. This black-box results, together with the white-box results (Table 1 in the main
 12 paper and §2 below) suggest that gradient masking is not present in the model and the improved robustness is indeed due to the
 13 inherent improvement of the model itself.

14 **§2 Results against stronger white-box attacks.** We have used random starts
 15 for all the white-box evaluations here as well as in the main paper. Here we run
 16 the evaluations 5 times and report the lowest performance and (max–min) as in
 17 the table on the right. The accuracy of our model under very strong white-box
 18 attacks still outperforms Madry model (44.8/45.4) by a large margin.

Attack Iteration	PGD500	PGD1000	CW500	CW1000
min over 5 runs	66.8	66.4	59.0	58.8
max–min	0.3	0.3	0.3	0.2

19 **§3 Loss plots.** The loss achieved with PGD adversary against our
 20 model increases in a fairly consistent way and plateaus rapidly for 20
 21 different runs with random starts, with relatively small variance (right).



22 **§4 Results against gradient-free attacks.** We evaluated our model
 23 using a gradient-free black-box attack based on greedy local search,¹
 24 and result is 89.9, further assuring the absence of gradient masking.

25 **To Reviewer 1:**

26 **Clarity and miscellaneous comments.** *i)* Your understanding of Alg 1 is correct. We will add more explanations on this for clarity.
 27 *ii)* $\mu = \{u_i\}$ represents the “empirical probability vector” over the support of $\{\mathbf{x}_i\}$. They are set to be uniform as shown in the
 28 Sinkhorn and IPOT algorithms in supplementary file, in the absence of prior knowledge. We have used this representation to make the
 29 presentation general enough to incorporate prior knowledge when available. We will make this clear and mention the Sinkhorn and
 30 IPOT sooner in revision. *iii)* Label smoothing parameter 0.5 is set without hyper-parameter tuning. The investigation on smoothing
 31 parameter was done afterwards. We didn’t aim to achieve the best performance by hyper-parameter search but instead to show the
 32 general applicability of our approach under a broad range of hyper-parameter settings. *iv)* The gap between PGD and the CW-variant
 33 is smaller under stronger attacks (§2) and we attribute the remaining gap to the nature of our model, where a one step unsupervised
 34 adversary is used for training, different from the multi-step supervised adversary typically used in Madry model. *v)* “Emphasizing 1
 35 attack iteration earlier” is a great suggestion. We will update this to avoid confusions as happened to **Reviewer 2.**

36 **Disentangling of distance and coupling.** Thanks for the great suggestion. We have preliminarily investigated the disentangling
 37 of distance and inter-sample coupling in our main paper as you have already noticed in Sec. 5.2 using the `identity matching`.
 38 Further investigation on it (esp. the coupling) as suggested by the reviewer is interesting and we plan to work on it as our next steps.

39 **To Reviewer 2:**

40 **Computational concern.** The number of iterations $T=1$ is used for our model as mentioned in line 225. We apologize for the
 41 confusion and will make it more clear as also suggested by **Reviewer 1.** Given that T is typically set to 7 in conventional PGD
 42 adversarial training (e.g. Madry), our approach does not take advantage of extra computation compared to conventional PGD training.

43 **Random targeted baseline.** We have experimented with random-targeted adversarial training as suggested by the reviewer. It
 44 achieves accuracy of 49.9/48.5 under PGD100/CW100 (min over 5 runs), outperformed by our model with a large margin (§2).

45 **Understanding of feature-scattering.** Conventional adversarial examples are *decision boundary oriented* (Fig.2), making the
 46 effective manifold for training deviate from the original due to *tilting* and *shrinking*, hindering performance (line 36-52), with label
 47 leaking as one manifesting phenomenon. Feature scattering is *inter-sample structure* oriented and promotes data diversity without
 48 drastically altering the structure of the manifold. We plan to conduct rigorous theoretical analysis of the proposed model as next step.

49 **To Reviewer 3:**

50 **Batch size.** As shown in the table on the right, larger batch size leads to better performance as it
 51 facilitates feature matching. Batch size of 60 is used for our model in the paper. Batch sizes larger
 52 than 60 lead to similar results. This observation is similar to other applications with embedded OT matching such as OT-GAN [48].

batch size	40	50	60	70	80
PDG100	58.7	62.7	68.4	68.2	68.4

53 **Label smoothing.** Label smoothing is necessary for our model. Our model (with
 54 1-step adversary) achieves compromised results without it, compared to standard
 55 PGD adversarial training (e.g. Madry) with 7-steps adversaries. This is an expected
 56 result as feature scattering makes the feature distributions more diffused (see Fig 1 in supplementary file), thus the corresponding
 57 label should ideally be “diffused” as well, in a spirit similar to *mixup* [70], which is achieved with label smoothing approximately in
 58 this work. Better schemes for joint treatment of feature and label scattering is an interesting topic and is left as our future work.

smooth para.	0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8
Madry	44.8	46.4	46.7	46.1	46.2	46.1	45.6	46.8	47.4
Ours	35.3	56.2	59.2	61.8	58.9	68.4	70.1	71.2	69.7

59 **Choice of distance.** Using cosine distance avoids introducing additional tuning parameter as the features are normalized before
 60 computing the distance. This and the usage of logits are just design choices. Other distance measures and intermediate features can
 61 be used together with our framework as well. We will explain this in the updated paper. As suggested by the reviewer, we will also
 62 introduce label leaking earlier in the introduction for clarity. We will release our trained model together with code as suggested.

¹N. Narodytska and S. Kasiviswanathan. Simple black-box adversarial attacks on deep neural networks, CVPRW17