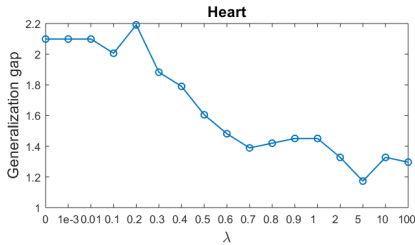1 We thank all the reviewers (**R1**-**R4**) for their insightful comments and invaluable feedback on our manuscript. As
2 the reviewers mentioned, our work shows the following strengths. (1) The proposed decomposition theorem and the
3 resulting bound are interesting (R1,R2) and new (R3,R4); (2) The theoretical study of non-convex metric learning is
4 important (R4), and the technique is general and of interest in establishing similar time-generalization tradeoff (R1,R3);
5 (3) Our algorithm is novel (R4), intuitive, and adequately uses the theoretical insights provided by the bound (R1);
6 (4) The theoretical analysis is complete (R2,R3), and the paper is easy to read (R1,R2,R3). We respond to the key
7 comments below but will address all feedback in the final version.

8 **[R3] Motivation.** Our motivation is to provide an approach to deriving a new generalization bound that considers the
9 parameters related to optimization, such as the number of iterations, which cannot be achieved by current uniform
10 convergence bounds. Then as summarized by **R1**,"the theoretical insights provided by the bound are adequately used
11 (smooth loss and classifiers, early stopping scheme, enforced small Lipschitz constants)" in the proposed algorithm.
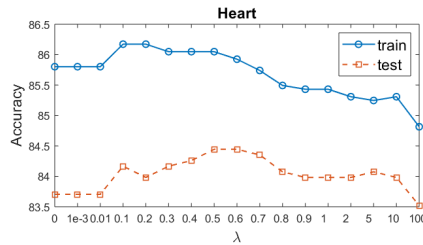
12 **[R3] Significance of the proved results.** As pointed out by **R1**, "the proposed bound seems to be new from a metric
13 learning point of view. Furthermore, it is general enough to be of potential interest for a broader range of researcher".
14 As suggested by **R2**, "the theoretical techniques may be of interest in establishing similar time generalization trade-off
15 bounds for non-convex problems, where the traditional route of uniform convergence does not make much sense" and
16 "the paper overcomes a common hurdle of requiring convex formulation or even optimization accuracy". In Sec.3.3.4,
17 we also provide discussion on the implications of Theorem 3. Greater stress on the significance of the proved results
18 and their wider applications will be added to Sec.3.3.4, Abstract and Introduction to make the significance more clear.

19 **[R3] Compared with generalization bounds of Cao et al. [6].** The bounds obtained in [6] is based on the framework
20 of uniform convergence. However, as pointed out by **R2**, when considering "time-generalization trade-off bounds", i.e.,
21 the generalization bounds considering the iteration number, "the traditional route of uniform convergence does not make
22 much sense". The approach proposed by our paper is fundamentally different from uniform convergence and it takes
23 optimization parameters into account. We will, on top of discussing the shortcomings of uniform convergence, highlight
24 the difference between our approach and their approaches including [6] in Introduction to make this more clear.
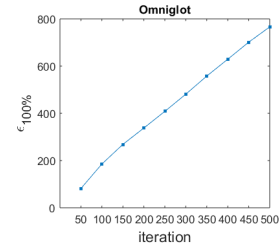
25 **[R1] Study the influence of the parameter $\lambda$; [R3] Performance of the model determined by loss term or regu-**
26 **larization term?** Figs.(a) and (b) show the influence of $\lambda$ on the generalization gap, training accuracy and test accuracy.
27 The generalization gap decreases with $\lambda$, which is consistent with our theoretical result that constraining the norms of
28 $L, x, r$ gives smaller Lipschitz constants, thereby tightening the bound. As the training accuracy generally decreases
29 with $\lambda$ as well, the test accuracy is highest when $\lambda = 0.5$. More detailed analysis will be added into the final version.



(a) Effect of $\lambda$ on the generalization gap.    (b) Effect of $\lambda$ on training and test accuracy.    (c) Concentration of deep ML.

30 **[R2, R4] Explore with deep metric learning (ML) algorithms; [R1] Extension to batch gradient descent.** We
31 have investigated the impact of training iterations on the concentration of parameters of deep ML and present the result
32 in Fig.(c). The network consists of three convolutional blocks and one fully connected (FC) layer and is trained with
33 mini-batch gradient descent (batch size=100); $\epsilon_{100\%}$ is calculated from the parameters of the FC layer over 10 rounds.
34 Similar to Fig.2 in the main paper, we see again that $\epsilon_{100\%}$ increases along training, indicating that the variance of
35 learned parameters becomes larger. While the empirical result is encouraging, we acknowledge that the proved theorem
36 cannot guarantee the learnability of mini-batch/stochastic gradient descent algorithms due to the randomness of training
37 instances introduced in each batch. This is indeed important and will be added into the Broader Impact Section.

38 **[R4] The method is only applicable to binary classification problems**. Our theoretical framework can be readily
39 generalized to multi-class classification. Theorem 2 and Lemma 2 are proved regardless of the number of classes.
40 Lemma 3 needs to be revised by considering the uniform convergence property of the hypothesis class defined over a
41 multi-class risk $R$. This has been discussed in some references. For example, based on the risk defined on p185 of
42 [Ref 1], uniform convergence is guaranteed by Theorem 8.1 on p187. More details will be added into the final version.
43 [Ref 1] Mohri, M., Rostamizadeh, A., & Talwalkar, A. (2018). Foundations of Machine Learning.

44 **[R2] Advantage gained by learning the representative instances.** We will add experiments by awarding other
45 methods, such as LMNN, the luxury of choosing representative points, as suggested by **R2**.