

1 **To Reviewer 1** Thanks for your comments and questions. It seems you misunderstood some key points and details.
2 Hope our explanation below could help to clarify some misunderstandings and confusion.

3 **(1) Theorem 2 in [Maurer 2005] is totally different for our Theorem 2**, although they may look similar. Theorem 2
4 in [Maurer 2005] is the generalization bound for *Single Task Learning* (ordinary supervised learning), which certainly
5 depends on the sample size m (number of samples in the task). Actually, this bound is from [Bousquet and Elisseeff
6 2002]. In contrast, Theorem 2 in our paper is the generalization bound for *Meta-Learning* with S/Q training, which is
7 independent of the sample size m of each task but depends on the total number of tasks n . To our knowledge, this is the
8 first sample-size-free bound for meta-learning.

9 **(2) The empirical evaluation corroborates our theoretical claims.** By “specific learning rate schedule”, we think
10 you meant the learning rate should satisfy $\zeta_t \leq c/t$ for a constant c and $t = \{1, \dots, T\}$ where T is the total number of
11 training steps, as stated in Theorem 3 in the Appendix. Notice that this condition can be easily satisfied with a fixed
12 learning rate in practice. For example, ProtoNets with a fixed learning rate $\zeta_t = 1e - 3$ converges in 24, 000 episodes
13 ($T = 24, 000$) and satisfies the condition with $c = 240$. Actually, Hardt et al. (2016) also conducted experiments with a
14 fixed learning rate $1e - 2$ and a constant number of training steps to verify their theory.

15 We would like to reiterate that our empirical evaluation is conducted with most popular meta-algorithms including
16 MAML [Finn et al. 2017], ProtoNets [Snell et al. 2017] and Bilevel programming [Franceschi et al. 2018] by strictly
17 following their training details on standard benchmarks (few-shot classification on *mini*Imagenet and sinusoidal few-shot
18 regression). We think the empirical evidence is sufficient to verify our theoretical claims.

19 **(3) Our results are not contradictory to those in [Triantafillou et al. 2020].** Notice that generalization gap \neq
20 test error. In fact, generalization gap = test error – training error (see [1] [2] for further reading). It is entirely possible
21 that test error keeps decreasing while generalization gap remains unchanged, because training error can also be
22 decreasing. This is exactly the case here. Fig. 2(b) in [Triantafillou et al. 2020] shows that the increase of shots
23 (inner-task sample size) reduces test error, which is evidently true. However, the increase of shots also reduces training
24 error, and both our theoretical bound and empirical evaluation show that the generalization gap keeps unchanged for
25 S/Q training.

26 [1] Understanding Machine Learning: From Theory to Algorithms. Shai Shalev-Shwartz and Shai Ben-David. 2014.

27 [2] Predicting the Generalization Gap in Deep Networks with Margin Distributions. Yiding Jiang, Dilip Krishnan, Hossein Mobahi,
28 Samy Bengio. ICLR 2019.

29 **(4) For your other comments:** 1) The inner-task gap vanishes because the expectation of the loss function w.r.t.
30 a new sample $z \sim \mathcal{D}$ is the same as that w.r.t. a new sample set $S^{ts} \sim \mathcal{D}^q$. In particular, inner-task
31 gap of S/Q training: $\mathbb{E}_{\mathcal{D}, S^{tr}, S^{ts}} [\mathbb{E}_z l(h, z) - \hat{L}(h, S^{ts})] = \mathbb{E}_{\mathcal{D}, S^{tr}} [\mathbb{E}_{S^{ts}} [\mathbb{E}_z l(h, z)] - \mathbb{E}_{S^{ts}} [\frac{1}{q} \sum_{z_j \in S^{ts}} l(h, z_j)]] =$
32 $\mathbb{E}_{\mathcal{D}, S^{tr}} [\mathbb{E}_z l(h, z) - \frac{1}{q} \sum_{z_j \in S^{ts}} \mathbb{E}_{z_j} l(h, z_j)] = \mathbb{E}_{\mathcal{D}, S^{tr}} [\mathbb{E}_z l(h, z) - \mathbb{E}_z l(h, z)] = 0.$

33 2) The statements regarding the generalization bounds of LOO loss $\epsilon(n, \beta, \tilde{\beta})$ and S/Q loss $\epsilon(n, \beta)$ are not contradictory.
34 When we say both of them are determined by the uniform stability β of the meta-algorithm, we did not mean “solely
35 determined”. The former also depends on the uniform stability $\tilde{\beta}$ of the inner-task algorithm but the latter does not.

36 3) Chen et al. (2019) did not use “batch multi-task training” as you mentioned, which is a traditional way for training
37 meta-algorithms as used in [Maurer 2005]. They simply trained an ordinary supervised classifier and compared it with
38 S/Q trained meta-algorithms. They never compared or discussed different training schemes for meta-algorithms.

39 4) The notation \bar{w}_i is defined in line 157.

40 **To Reviewer 2** Thanks for your comments and suggestion. The results of LOO training were previously put in a
41 separate section, but due to space limitation, we merged them with the results of S/Q training. We will reorganize the
42 paper in the final version where more space will be given.

43 **To Reviewer 3** Thanks for your comments and feedback. Reptile is an inspiring meta-algorithm which does not need a
44 S/Q split for training but still achieves comparable performance with MAML. To make the traditional generalization
45 bound apply to Reptile, we may first need to derive the randomized uniform stability of Reptile w.r.t. its update rule,
46 which is not equivalent to “meta-level SGD”. We think it would be very interesting to study the generalization of Reptile
47 and will add more discussion in the revised version. Thank you for bringing that up.

48 **To Reviewer 4** Thanks for your comments and feedback. Indeed the discussion of LOO meta-training is not related to
49 Theorem 4 in our paper, but we introduce LOO meta-training because it is “a surrogate to the traditional scheme that is
50 compatible with gradient-based and metric-based meta-learning algorithms” (Reviewer 3 said it nicely and we quote).
51 Besides, it is a nice comparison to S/Q training in terms of generalization bounds. We will further clarify this in the
52 final version.