1 **Reviewer 1:** **Unclear about the evaluation for outer iterations; Does the number of aggregated tasks affect**
2 **convergence:** *Great question! Yes, the total complexity is proportional to the number of aggregated tasks. In addition,*
3 *in terms of updating task-specific parameters, ANIL takes the same steps as MAML, and the outer-loop gradient (line*
4 *10 of Alg. 1) also depends on the inner-loop outputs $w_{k,N}^i$ of tasks in $\mathcal{B}_k$. We will clarify it in the revision.*

5 **Add experiments to compare ANIL and MAML and w.r.t. the size $B$ of samples:** *Thanks for the suggestion! We*
6 *will absolutely follow these suggestions to add experiments in the revision.*

7 **Why sample size in inner-loop is not taken into analysis, as Fallah et al. [4] does:** *Great question! In our setting,*
8 *the inner-loop loss functions take a finite-sum form over pre-assigned samples. As a result, the inner-loop updates take*
9 *full gradient descent without data sampling, and hence gradient estimation bias (which can introduce sample size) does*
10 *not exist in convergence bound. This setting has also been considered in Rajeswaran et al. [24], Ji et al. [13]. As a*
11 *comparison, Fallah et al. [4] considered a different setting, where loss functions take the form in expectation and fresh*
12 *data are sampled as the algorithm runs. As a result, their analysis involves an estimation bias, which introduces the*
13 *dependence on the number of samples.*

14 **Experiments for non-convex and strongly convex cases with the same stepsize:** *Great point! We have run more*
15 *experiments on FC100 with the same stepsize $0.03, 0.05, 0.1$ for both cases, and the nature of results remain the same.*

16 **Elaborate more for line 170:** *The statement specifically refers to Theorem 1, where increasing $N$ leads to larger*
17 *stepsize $\beta_w$, which yields faster convergence rate $\mathcal{O}(\frac{1}{K\beta_w})$. We will clarify it in the revision.*

18 **Reviewer 2:** **Dependence on $\kappa$. iMAML depends on $\sqrt{\kappa}$ in contrast to poly$(\kappa)$ of this work:** *Great question!*
19 *High-level speaking, better dependence on $\kappa$ for iMAML is based on an ideal solution of an inner-loop optimization*
20 *problem, which can take many iterations. ANIL takes only a few inner-loop iterations (thus a lower cost), but has*
21 *worse outer-loop convergence (in terms of $\kappa$). Technically speaking, smoothness analysis of iMAML upper-bounds*
22 *the distance between two optimal points $w_*^i(w_1)$ and $w_*^i(w_2)$, each obtained by solving an inner-loop optimization*
23 *problem. As a comparison, analysis of ANIL upper-bounds the distance between two inner-loop paths, which sums up*
24 *the distances between all corresponding points on the two paths (see eq. $(21)$). This results in a worse dependence in $\kappa$.*

25 **Add an experiment to verify the tightness:** *Great point! We will definitely add such an experiment in the revision.*

26 **Extra assumption on Lipschitzness of the objective, which is not for iMAML:** *We take this assumption to ensure*
27 *the meta gradient to be bounded. As a comparison, iMAML alternatively assumes the search space of parameters to be*
28 *bounded (see Theorem 1 therein) so that the meta gradient (eq. (5) therein) can be bounded.*

29 **The role of $N$ in the theory seems to make convergence only slower:** *The exponential term has a worse dependence*
30 *on constants and $\tau, M$ than the linear term (we will add explicit forms in the revision), and hence the choice of $N$*
31 *depends on how large $\kappa$ is. For large $\kappa$, as the reviewer also pointed out, a small $N = 2$ is a better choice. However,*
32 *when $\kappa$ is not very large, e.g., in our experiments (in which increasing $N$ accelerates the iteration rate), the exponential*
33 *term dominates for a small $N$, and hence a larger $N$ is preferred. We will clarify it in the revision.*
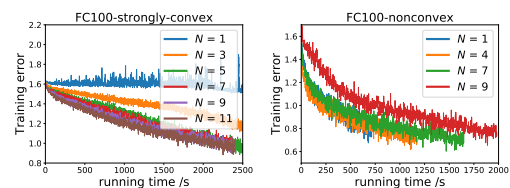
34 **Optimality of $N = 1$ contradicts the experiments where $N = 4, 7$ are the best:** *We assume the reviewer refers*
35 *to the experiments in left plot of Figure 2(a). This can be due to the fact that the influence of $N$ w.r.t. the number*
36 *of outer-loop iterations is offset by other constant-level parameters for small $N$. Evidently, right plot of Figure 2(a)*
37 *indicates that $N = 1$ is optimal w.r.t. the running time, which agrees with our result on computational complexity.*

38 **Suggestions on presentation and references:** *Many thanks! We will follow these suggestions to improve our paper.*

39 **Reviewer 3:** *We thank the reviewer for the positive comments!*

40 **Reviewer 4:** **Comments on insight of theoretical results:** *Our results theoretically characterize the order-level*
41 *computational complexity for ANIL and its comparison to MAML. In addition, our analysis techniques can be useful for*
42 *developing guarantee for other meta-learning and more broadly bi-level optimization algorithms.*

43

44 **Convergence analysis is done with vanilla gradient descent but**
45 **all experiments are done with Adam; Experiments with purely**
46 **first-order methods:** *Great point! We have done new experiments*
47 *on FC100 dataset using mini-batch SGD with a learning rate of*
48 *0.05, and the results are given in the figures to the right. It can be*
49 *seen that the nature of the results remains the same as those in our*
50 *paper. More results will be added in the revision.*



51 **Run experiments over different random seeds and over different hyper-parameter settings:** *Many thanks! We*
52 *will definitely provide these experimental results in the revision.*