

1 We thank all the reviewers for their thoughtful comments and positive feedback. We will first address two major points
2 raised by the reviewers and then answer individual questions.

3 Major Points:

- 4 1. **(Reviewers 1 and 3)** - “Providing more empirical results and refined baselines will improve the experiments section
5 and be useful to study the strength of this algorithm in practice.” (paraphrase)
 - 6 • We agree and will have more experiments in revision. Specifically, we plan to include few-shot learning evaluations
7 on MiniImageNet [47] and detailed numerical results for both the meta-learning and federated learning settings.
8 Furthermore, following the suggestion of Reviewer 1, we will investigate a comparison with Reptile+Meta-SGD
9 [39], which would also learn a per-coordinate learning rate. One caveat here is that Meta-SGD was designed for a
10 MAML-like approach of only taking one gradient step within-task, whereas Reptile takes many iterations; since
11 Meta-SGD uses higher-order differentiation, repeated application may slow it down. Note that we do discuss
12 (MAML+) Meta-SGD briefly in lines 305 -> 307.
- 13 2. **(Reviewer 2)** - “The paper is missing a detailed discussion/examples of the behaviour of V_Ψ , which makes it hard to
14 judge the sensibility of the proven bounds.”
 - 15 • We agree that examples of V_Ψ are needed to understand the results, as it is the main measure of task-similarity
16 and if V_Ψ is small then so is the average regret. For the case of a fixed comparator (Theorem 3.2), we do give a
17 simple example in lines 136 -> 137; here V_Ψ is proportional to the empirical standard deviation of the optimal
18 task-parameters, so if they are close (i.e. tasks are similar) then V_Ψ is small. As we will describe in more detail in
19 revision, the case when the comparator is varying is similar, as V_Ψ is now proportional to the average deviation of
20 the optimal task-parameters from a shifting sequence of vectors. For example, if we see one task every day for
21 year, and Ψ is a sequence that fixes a single comparator for each month, then V_Ψ^2 is roughly the average over the
22 months of the empirical variance of each month’s thirty or so optimal task-parameters from that month’s fixed
23 comparator. This can be very small if the variation between tasks is well-described by a seasonal trend.

24 Reviewer #1:

- 25 1. “It would be useful to include a reference for the regret of OGD.” (paraphrase)
 - 26 • We agree and will add a pointer to the Shalev-Shwartz survey [49, Theorem 2.15 for $R(w) = \frac{1}{2\eta}\|w - \phi\|_2^2$].
- 27 2. “Is the form of the regret-upper-bound $\hat{R}_t(x_t)$ nice in cases more general than just that of OGD?” (paraphrase)
 - 28 • Yes! This is a main reason we expect this framework to be broadly applicable. In our paper, we show that several
29 results (specifically Theorems 3.1 and 3.2) hold for any algorithm in the OMD/FTRL family, which includes not
30 just OGD but also other classical methods such as exponentiated gradient/multiplicative weights. Even more
31 generally, regret guarantees often include terms that depend on some measure of distance from an initial state,
32 which are often amenable to study (e.g. because norms are convex). We will elaborate on this in the revision.
- 33 3. “On the line 108 -> 109, it’s not clear how $\hat{R}_T \leq o(T) + \min_x \frac{1}{T} \sum_t \hat{R}_t$.”
 - 34 • We believe this asking how the right-hand-side goes to zero. This is a typo - the first term should be $o_T(1)$, i.e.
35 sub-constant, not sub-linear. Similarly, the last term in the statement of Theorem 3.1 should be divided by T (the
36 expressions in the proof are correct as-is). We apologize for both errors and will correct them in revision.
- 37 4. “For FedAvg case, does the improvement comes from the meta-learning treatment (where we optimize for the
38 initialization) or the ARUBA algorithm itself (e.g. the fact that the learning rate is adaptable)?” (paraphrase)
 - 39 • In fact to get FedAvg+ARUBA we do not modify FedAvg except to adapt the learning rate - the global model is
40 still learned in the same way. This is possible because FedAvg is equivalent to Reptile with the outer-loop update
41 coefficient set to 1.0. So the improvement is indeed coming from the adaptivity.

42 Reviewer #2:

- 43 1. “Complicated and overloaded notations: too many versions of regret and bounds with very similar symbols, the
44 sequence of ψ_t for dynamic regret is not defined for notations of Theorem 3.1.”
 - 45 • We will make sure that the mathematical presentation is as clean as possible in the revision. One way of reducing
46 the many variations on the regret notations is to depend less on accents and represent regret-upper-bounds by a
47 different capital letter (e.g. \mathbf{U}) to better distinguish from regret terms (\mathbf{R}). As for the sequence of ψ_t , this can be
48 any arbitrary sequence of vectors in the action space Θ .

49 Reviewer #3:

- 50 1. “What factor leads to the big experimental difference between meta-learning, where the improvement over Adam is
51 relatively small, and federated learning, where the improvement over FedAvg is a lot more significant?” (paraphrase)
 - 52 • Our explanation, which we will add in revision, is based on the nature of data. Whereas standard evaluation
53 datasets in few-shot learning consist of tasks with identical amounts of i.i.d. data, the Shakespeare benchmark
54 we use for federated learning has tasks with highly variable amounts of data (two different roles can have very
55 different numbers of lines) and data that is not i.i.d. (lines are not shuffled but split in their order of appearance
56 in the play). Our method may be better able to handle such data, for example by tempering noisy directions in
57 low-data tasks by learning which directions are important based on the distance traveled in high-data tasks.