

1 We thank the reviewers for valuable feedbacks. In below response, the reference order is the same as submitted work.

2 **Reviewer # 1:**

3 1) *The convergence theory should be w.r.t $\sum_i f_i(\theta_i)$ and/or $\sum_i f_i(w)$, why care about $\|w - \theta_i\|$?* → pFedMe can
4 be written compactly as $\min_w \min_{\theta_i} (1/N) \sum_{i=1}^N \{f_i(\theta_i) + (\lambda/2)\|\theta_i - w\|^2\}$, which shows that there is dependency
5 between solution components $\{\theta_i\}_i$ and w . Thus, we cannot separately evaluate each objective term. Similar to, e.g.,
6 l_2 -regularization where $f_i(w) = l_i(w) + (\lambda/2)\|w\|^2$ for some loss function $l_i(w)$ and the overfitting regularization term
7 $\|w\|^2$ plays a critical role in convergence of f_i , we believe that it makes sense to consider the personalized regularization
8 $\|\theta_i - w\|^2$ in pFedMe’s convergence analysis. If we want to move far away from w , setting λ small is one option.

9 2) *To what extent of λ to leverage client aggregation? From Theorem 1, $\lambda = 0$ to minimize $F(w) - F^*$?* → The
10 trade-off between personalization and exploiting data aggregation by varying λ is skilfully phrased by the reviewer.
11 Note that $\lambda \in (0, \infty)$ avoids extreme cases of $\lambda = 0$ (no FL) and $\lambda = \infty$ (no personalization). The reviewer is correct
12 in that setting $\lambda = 0$ reduces to simple ERM on each local client, but it is not true that in order to minimize $F(w) - F^*$,
13 we should just set $\lambda = 0$, because Theorem 1 is based on using the Moreau envelope F_i , which implicitly requires $\lambda > 0$
14 by its definition [39]. In other words, if $\lambda = 0$, there would be no local update in line 8, and hence no Theorem 1 (and
15 2). On the other hand, similar to the standard way to find a “sweet spot” for hyper-parameters in l_1 - or l_2 -regularization,
16 we fine-tune λ based on each dataset for each task in our experiments.

17 3) *Are personalized model and global model plottings w.r.t θ_i and w , respectively?* → Yes.

18 4) *Use personalized models for Per-FedAvg’s performance?* → We did use $\theta_i(w) = w - \alpha \nabla f_i(w)$ for Per-FedAvg.

19 5) *How to compare convergence rates with prior methods when they have different objectives?* → The standard
20 convergence analysis is performed according to the **loss functions’ properties such as strong/non-convexity and/or**
21 **smoothness** [49], regardless of different learning tasks/objectives. Therefore, the quadratic speedup of pFedMe is
22 compared with other linear-speedup methods having the same strong convexity and smoothness properties [3, 6, 19].

23 **Reviewer #2:**

24 1) *Discuss with decentralized frameworks of Vanhaesebrouck, Bellet, Tommasi (2017)* → We will update it.

25 2) *How to choose λ to ensure strong convexity of h , considering the assumption on f_i is in expectation?* →
26 Assumption 1 can apply a stronger L -smoothness assumption on $\tilde{f}_i(\theta_i; \xi_i), \forall \xi_i$, instead of f_i . In this case,
27 we can choose arbitrary $\lambda > 0$ for strongly convex and $\lambda > L$ for nonconvex loss functions, respectively.

28 **Reviewer #4:**

29 1) *β larger than 1 could yield some instabilities (e.g., Fig 5) and inducing one more*
30 *costly hyperparameter to tune?* → In Theorem 1, we need $\eta \leq \frac{\eta_1}{\beta R}$, which means
31 larger β requires **smaller** η for the stability. In practice, we can fix β to an arbitrary
32 constant, say 3, and perform the simple fine-tuning task for η .

33 2) *Run $K * R$ local steps for Per-FedAvg and FedAvg.* → It would be unfair to
34 compare pFedMe with these using $K * R$ local steps. This is because for every K
35 steps in a local round r , pFedMe only uses a single mini-batch, and thus only R
36 mini-batches during R local rounds. On the other hand, Per-FedAvg uses 2 different
37 mini-batches for 1 local update, and thus after $K * R$ local rounds, it will perform
38 $2 * K * R$ mini-batch updates. Similarly, FedAvg uses one mini-batch for 1 local
39 update, so in total it will have $K * R$ mini-batch updates. However, we still report an additional comparison when
40 FedAvg and Per-FedAvg were trained over $K * R$ local steps in Fig. 1 here (in a strongly convex setting), which shows
41 the personalized model (PM) of pFedMe still outperforms others.

42 3) *Use a more realistic FL dataset than MNIST (LEAF benchmark)?* → For synthetic data, we used a similar method to
43 LEAF to generate the data. For real data, in LEAF, FEMNIST is distributed to a large number of clients, each with a
44 small local dataset (of size around 226), which is not suitable for pFedMe because it needs a larger local dataset on
45 each client (due to fresh mini-batch sampling in every round in line 7). Thus we use MNIST with a smaller number of
46 clients, so each client can have up to 3,834 data samples. We note that the way we generate non-i.i.d and heterogeneous
47 data using MNIST is similar to that using FEMNIST in LEAF.

48 4) *How do (local) test performances computed and distributed? Any weighting?* → W.r.t test accuracy of all algorithms,
49 we sum all correctly classified samples over clients and divide it by the total number of samples of all clients without
50 weighting. Their histograms are quite similar and we reported their means and standard deviation in Table 1.

51 5) *The hyper parameters were optimised on the local test sets or on local validation sets?* → On local test set. We have
52 additionally fine-tuned hyper-parameters on a validation set, and still obtain the same values as those on the test set.

53 6) *Why is global model of pFedMe is lagging in Fig. 2? Due to varying number of clients or lower sampling rate?* →
54 The **non-fine-tuned hyper-parameter** is the main reason for the lagging performance of the global model (GM) in
55 Fig. 2. In cases where hyperparameters are fine-tuned, the GM is consistent and performs well compared to FedAvg
56 (see Table 1). We have run more experiments and observed that the GM performance is still consistent with varying
57 numbers of clients and sampling rates. The result is not presented here due to the lack of space.

58 7) Finally, we thank the reviewer for correcting several typos.

Algorithm	MNIST	Synthetic
FedAvg	94.11 ± 0.05	77.69 ± 0.2
Per-FedAvg	94.22 ± 0.02	79.79 ± 0.09
pFedMe-GM	94.18 ± 0.06	78.65 ± 0.25
pFedMe-PM	95.62 ± 0.04	83.20 ± 0.06

Figure 1: Accuracy comparison with fine-tuned hyperparameters. $R = 20, K = 5$ for pFedMe, $R = 100$ for others.